



Московский государственный
университет им. М. В. Ломоносова
Факультет вычислительной математики и кибернетики

И.А. Волкова, А.А. Вылиток, Т.В. Руденко

Формальные грамматики и языки. Элементы теории трансляции

Учебное пособие для студентов II курса

(издание третье, переработанное и дополненное)

Москва
2009

УДК 519.682.1(075.8)

ББК 22.19я73

В67

*Печатается по решению Редакционно-издательского совета факультета
вычислительной математики и кибернетики
МГУ им. М. В. Ломоносова*

Рецензенты:

проф., д.ф.-м.н. Машечкин И. В.

доцент, к.ф.-м.н. Терехин А.Н.

И. А. Волкова, А. А. Вылиток, Т. В. Руденко

Формальные грамматики и языки. Элементы теории трансляции: Учебное пособие для студентов II курса (издание третье, переработанное и дополненное). - М.: Издательский отдел факультета ВМиК МГУ им. М.В.Ломоносова (лицензия ИД № 05899 от 24.09.2001), 2009 – 115 с.

ISBN 978-5-89407-395-8

ISBN 978-5-317-03085-8

Приводятся основные определения, понятия и алгоритмы теории формальных грамматик и языков, некоторые методы трансляции, а также наборы задач по рассматриваемым темам. Излагаемые методы трансляции проиллюстрированы на примере модельного языка.

В третьем издании все примеры реализации методов и алгоритмов приводятся на языке Си++, изменен и расширен набор задач по всем разделам.

Для студентов факультета ВМК в поддержку основного лекционного курса «Системы программирования» и для преподавателей, ведущих практические занятия по этому курсу.

УДК 519.682.1(075.8)

ББК 22.19я73

ISBN 978-5-89407-395-8

ISBN 978-5-317-03085-8

© Факультет вычислительной математики и кибернетики МГУ им. М. В. Ломоносова, 2009

© И. А. Волкова, А. А. Вылиток, Т. В. Руденко, 2009

Элементы теории формальных языков и грамматик

Введение

В этом разделе изложены некоторые аспекты теории формальных языков, существенные с точки зрения трансляции. Здесь введены базовые понятия и даны определения, связанные с одним из основных механизмов определения языков — грамматиками, приведена наиболее распространенная классификация грамматик (по Хомскому). Особое внимание уделяется контекстно-свободным и регулярным грамматикам. Грамматики этих классов широко используются при трансляции языков программирования. В пособии не приводятся доказательства корректности алгоритмов и большинства сформулированных фактов, свойств, утверждений — их можно найти в книгах, указанных в списке литературы.

Основные понятия и определения

Определение: *алфавит* — это конечное множество символов.

Предполагается, что термин «символ» имеет достаточно ясный интуитивный смысл и не нуждается в дальнейшем уточнении.

Определение: *цепочкой символов в алфавите V* называется любая конечная последовательность символов этого алфавита.

Определение: цепочка, которая не содержит ни одного символа, называется *пустой цепочкой*. Для ее обозначения будем использовать греческую букву ε .

Предполагается, что сама буква ε в алфавит V не входит; она лишь помогает обозначить пустую последовательность символов.

Определение: если α и β — цепочки, то цепочка $\alpha\beta$ (результат приписывания цепочки β в конец цепочки α), называется *конкатенацией* (или *сцеплением*) цепочек α и β . Конкатенацию можно считать двуместной операцией над цепочками: $\alpha\beta = \alpha\beta$.

Например, если $\alpha = ab$ и $\beta = cd$, то $\alpha\beta = abcd$.

Для любой цепочки α справедливы равенства: $\alpha\varepsilon = \varepsilon\alpha = \alpha$.

Для любых цепочек α , β , γ справедливо $(\alpha\beta)\gamma = \alpha(\beta\gamma) = \alpha\beta\gamma$ (свойство ассоциативности операции конкатенации).

Определение: *обращением* (или *реверсом*) цепочки α называется цепочка, символы которой записаны в обратном порядке.

Обращение цепочки α будем обозначать α^R .

Например, если $\alpha = abcdef$, то $\alpha^R = fedcba$.

Для пустой цепочки: $\varepsilon^R = \varepsilon$.

Определение: n -ой степенью цепочки α (будем обозначать α^n) называется конкатенация n цепочек α :

$$\alpha^n = \underbrace{\alpha\alpha \dots \alpha\alpha}_n.$$

Свойства степени: $\alpha^0 = \varepsilon$; $\alpha^n = \alpha\alpha^{n-1} = \alpha^{n-1}\alpha$

Определение: *длина цепочки* — это число составляющих ее символов (или длина последовательности символов).

Например, если $\alpha = abbcaad$, то длина α равна 7.

Длину цепочки α будем обозначать $|\alpha|$. Длина ε равна 0.

Определение: через $|\alpha|_s$ обозначают *число вхождений* символа s в цепочку α .

Например, $|babb|_a = 1$, $|babb|_b = 3$, $|babb|_c = 0$.

Определение: обозначим через V^* множество, содержащее все цепочки в алфавите V , включая пустую цепочку ε .

Например, если $V = \{0, 1\}$, то $V^* = \{\varepsilon, 0, 1, 00, 11, 01, 10, 000, 001, 011, \dots\}$.

Определение: обозначим через V^+ множество, содержащее все цепочки в алфавите V , исключая пустую цепочку ε .

Следовательно, $V^* = V^+ \cup \{\varepsilon\}$.

Определение: *язык* в алфавите V — это подмножество множества всех цепочек в этом алфавите. Для любого языка L справедливо $L \subseteq V^*$.

Известны различные способы описания языков [3]. Конечный язык можно описать простым перечислением его цепочек. Поскольку формальный язык может быть и бесконечным, требуются механизмы, позволяющие конечным образом представлять бесконечные языки. Можно выделить два основных подхода для такого представления: механизм распознавания и механизм порождения (генерации).

Механизм распознавания (*распознаватель*), по сути, является процедурой специального вида, которая по заданной цепочке определяет, принадлежит ли она языку. Если принадлежит, то процедура останавливается с ответом «да», т. е. *допускает* цепочку; иначе — останавливается с ответом «нет» или заикливается. Язык, определяемый распознавателем — это множество всех цепочек, которые он допускает.

Примеры распознавателей:

- Машина Тьюринга (МТ). Язык, который можно задать с помощью МТ, называется *рекурсивно перечислимым*.¹⁾ Вместо МТ можно использовать эквивалентные алгоритмические схемы: нормальный алгоритм Маркова (НАМ), машину Поста и др.
- Линейно ограниченный автомат (ЛОА). Представляет собой МТ, в которой лента не бесконечна, а ограничена длиной входного слова²⁾. Определяет *контекстно-зависимые языки*.

1) Термин «рекурсивно перечислимый» происходит из теории рекурсивных функций, являющейся, также как НАМ и МТ, одной из формализаций понятия вычислимости (алгоритма). Смысл остальных названий, характеризующих распознаваемые языки, будет пояснен ниже.

2) ЛОА является недетерминированной машиной: в какой-то момент вычисления (во время работы ЛОА над заданной цепочкой) может возникнуть ситуация, когда есть две или более подходящие для исполнения команды, то есть выбор исполняемой команды не детерминирован. С этого момента возникают две или более копии ЛОА, соответствующие каждому возможному выбору; эти копии продолжают вычисления независимо (параллельно). Ответ «да» ЛОА дает, если хотя бы одно из вычислений успешно завершается. Известен алгоритм, позволяющий по любой недетерминированной МТ построить эквивалентную детерминированную

- Автомат с магазинной (внешней) памятью (МП-автомат). В отличие от ЛОА, головка не может изменять входное слово и не может сдвигаться влево; имеется дополнительная бесконечная память (*магазин*, или стек), работающая по дисциплине LIFO. Определяет *контекстно-свободные языки*.
- Конечный автомат (КА). Отличается от МП-автомата отсутствием магазина. Определяет *регулярные языки*.

Основной способ реализации механизма порождения — использование порождающих грамматик, которые иногда называют грамматиками Хомского. На изучении порождающих грамматик мы и остановимся подробно, и именно этот способ описания языков чаще всего будем использовать в дальнейшем.

Отметим, что не каждый формальный язык можно задать с помощью конечного описания. Действительно, само описание можно рассматривать как цепочку в некотором расширенном алфавите. Следовательно, множество описаний языков счётно, так как множество всех цепочек в заданном алфавите счётно. Каждый формальный язык в алфавите V является подмножеством счётного множества V^* . Из теории множеств известно, что множество всех подмножеств счётного множества является несчётным. Таким образом, мощность множества формальных языков больше мощности множества конечных описаний и, следовательно, не каждый язык представим в виде конечного описания.

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

Определение: *порождающая грамматика* G — это четверка $\langle T, N, P, S \rangle$, где

- T — алфавит терминальных символов (терминалов);
- N — алфавит нетерминальных символов (нетерминалов), $T \cap N = \emptyset$;
- P — конечное подмножество множества $(T \cup N)^+ \times (T \cup N)^*$; элемент (α, β) множества P называется *правилом вывода* и записывается в виде $\alpha \rightarrow \beta$; α называется *левой частью* правила, β — *правой частью*; левая часть любого правила из P обязана содержать хотя бы один нетерминал;
- S — начальный символ (цель) грамматики, $S \in N$.

Для записи правил вывода с одинаковыми левыми частями

$$\alpha \rightarrow \beta_1 \quad \alpha \rightarrow \beta_2 \quad \dots \quad \alpha \rightarrow \beta_n$$

будем пользоваться сокращенной записью

$$\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n.$$

Каждое β_i ($i = 1, 2, \dots, n$) будем называть *альтернативой* правила вывода из цепочки α .

Пример грамматики:

$$G_{example} = \langle \{0, 1\}, \{A, S\}, P, S \rangle, \quad \text{где } P \text{ состоит из правил:}$$

МТ. Однако до сих пор открыт вопрос, существует ли для любого ЛОА эквивалентный ему детерминированный ЛОА.

$$\begin{aligned} S &\rightarrow 0A1 \\ 0A &\rightarrow 00A1 \\ A &\rightarrow \varepsilon \end{aligned}$$

Определение: цепочка $\beta \in (T \cup N)^*$ непосредственно выводима из цепочки $\alpha \in (T \cup N)^+$ в грамматике $G = \langle T, N, P, S \rangle$ (обозначается $\alpha \rightarrow_G \beta$), если $\alpha = \xi_1 \gamma \xi_2$, $\beta = \xi_1 \delta \xi_2$, где $\xi_1, \xi_2, \delta \in (T \cup N)^*$, $\gamma \in (T \cup N)^+$ и правило вывода $\gamma \rightarrow \delta$ содержится в P . Индекс G в обозначении \rightarrow_G обычно опускают, если понятно, о какой грамматике идет речь.

Например, цепочка $00A11$ непосредственно выводима из $0A1$ в грамматике $G_{example} : \underline{0A1} \rightarrow \underline{00A1}1$. Здесь цепочка $\underline{0A1}$, подчеркнутая двойной чертой, играет роль подцепочки γ из определения, цепочка $\underline{00A1}$ играет роль подцепочки δ , $\xi_1 = \varepsilon$, $\xi_2 = 1$.

Определение: цепочка $\beta \in (T \cup N)^*$ выводима из цепочки $\alpha \in (T \cup N)^+$ в грамматике $G = \langle T, N, P, S \rangle$ (обозначается $\alpha \Rightarrow_G \beta$), 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 . Индекс G в обозначении \Rightarrow_G опускают, если понятно, какая грамматика подразумевается.

Например, $S \Rightarrow 000A111$ в грамматике $G_{example}$ (см. пример выше), т.к. существует вывод $S \rightarrow 0A1 \rightarrow 00A11 \rightarrow 000A111$. Длина вывода равна 3.

Определение: языком, порождаемым грамматикой $G = \langle T, N, P, S \rangle$, называется множество $L(G) = \{\alpha \in T^* \mid S \Rightarrow \alpha\}$.

Другими словами, $L(G)$ — это все цепочки в алфавите T , которые выводимы из S с помощью правил P . Например, $L(G_{example}) = \{0^n 1^n \mid n > 0\}$.

Определение: цепочка $\alpha \in (T \cup N)^*$, для которой $S \Rightarrow \alpha$, называется *сентенциальной формой* в грамматике $G = \langle T, N, P, S \rangle$.

Таким образом, язык, порождаемый грамматикой, можно определить как множество терминальных сентенциальных форм.

Определение: грамматики G_1 и G_2 называются *эквивалентными*, если $L(G_1) = L(G_2)$.

Пример. Грамматики $G_1 = \langle \{0,1\}, \{A,S\}, P_1, S \rangle$ и $G_2 = \langle \{0,1\}, \{S\}, P_2, S \rangle$ с правилами

$$\begin{array}{ll} P_1: & P_2: \\ S & \rightarrow 0A1 \\ 0A & \rightarrow 00A1 \\ A & \rightarrow \varepsilon \end{array} \quad \begin{array}{l} S \\ S \end{array} \rightarrow 0S1 \mid 01$$

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

Определение: грамматики G_1 и G_2 *почти эквивалентны*, если $L(G_1) \cup \{\varepsilon\} = L(G_2) \cup \{\varepsilon\}$.

Другими словами, грамматики почти эквивалентны, если языки, ими порождаемые, отличаются не более, чем на ε .

Например, почти эквивалентны грамматики

$$G_1 = \langle \{0,1\}, \{A, S\}, P_1, S \rangle \quad \text{и} \quad G_2 = \langle \{0, 1\}, \{S\}, P_2, S \rangle \quad \text{с правилами}$$

$$\begin{array}{l}
 P_1: \\
 S \rightarrow 0A1 \\
 0A \rightarrow 00A1 \\
 A \rightarrow \varepsilon
 \end{array}
 ,
 \quad
 \begin{array}{l}
 P_2: \\
 S \rightarrow 0S1 \mid \varepsilon
 \end{array}$$

так как $L(G_1) = \{0^n 1^n \mid n > 0\}$, а $L(G_2) = \{0^n 1^n \mid n \geq 0\}$, т. е. $L(G_2)$ состоит из всех цепочек языка $L(G_1)$ и пустой цепочки, которая в $L(G_1)$ не входит.

Классификация грамматик и языков по Хомскому

Определим с помощью ограничений на вид правил вывода четыре типа грамматик: тип 0, тип 1, тип 2, тип 3. Каждому типу грамматик соответствует свой класс³⁾ языков. Если язык порождается грамматикой типа i (для $i = 0, 1, 2, 3$), то он является языком типа i .

Тип 0

Любая порождающая грамматика является грамматикой *типа 0*. На вид правил грамматик этого типа не накладывается никаких дополнительных ограничений. Класс языков типа 0 совпадает с классом рекурсивно перечислимых языков.

Грамматика с ограничениями на вид правил вывода

Тип 1

Грамматика $G = \langle T, N, P, S \rangle$ называется *неукорачивающей*, если правая часть каждого правила из P не короче левой части (т. е. для любого правила $\alpha \rightarrow \beta \in P$ выполняется неравенство $|\alpha| \leq |\beta|$). В виде исключения в неукорачивающей грамматике допускается наличие правила $S \rightarrow \varepsilon$, при условии, что S (начальный символ) не встречается в правых частях правил.

Грамматикой *типа 1* будем называть неукорачивающую грамматику.

Тип 1 в некоторых источниках определяют с помощью так называемых контекстно-зависимых грамматик.

Грамматика $G = \langle T, N, P, S \rangle$ называется *контекстно-зависимой (КЗ)*, если каждое правило из P имеет вид $\alpha \rightarrow \beta$, где $\alpha = \xi_1 A \xi_2$, $\beta = \xi_1 \gamma \xi_2$, $A \in N$, $\gamma \in (T \cup N)^+$, $\xi_1, \xi_2 \in (T \cup N)^*$. В виде исключения в КЗ-грамматике допускается наличие правила с пустой правой частью $S \rightarrow \varepsilon$, при условии, что S (начальный символ) не встречается в правых частях правил.

Цепочку ξ_1 называют *левым контекстом*, цепочку ξ_2 называют *правым контекстом*. Язык, порождаемый контекстно-зависимой грамматикой, называется *контекстно-зависимым* языком.

Замечание

Из определений следует, что если язык, порождаемый контекстно-зависимой или неукорачивающей грамматикой $G = \langle T, N, P, S \rangle$, содержит пустую цепочку, то эта цепочка выводится в G за один шаг с помощью правила $S \rightarrow \varepsilon$. Других выводов для ε не существует. Если среди правил P нет правила $S \rightarrow \varepsilon$, то порождаемый контекстно-зависимой или неукорачивающей грамматикой G язык не содержит пустую цепочку ε .

³⁾ Классом обычно называют множество, элементы которого сами являются множествами.

Утверждение 1. Пусть L — формальный язык. Следующие утверждения эквивалентны:

- 1) существует контекстно-зависимая грамматика G_1 , такая что $L = L(G_1)$;
- 2) существует неукорачивающая грамматика G_2 , такая что $L = L(G_2)$.

Очевидно, что из (1) следует (2): любая контекстно-зависимая грамматика удовлетворяет ограничениям неукорачивающей грамматики (см. определения). Доказательство в обратную сторону основывается на том, что можно каждое неукорачивающее правило заменить эквивалентной серией контекстно-зависимых правил.

Из утверждения 1 следует, что язык, порождаемый неукорачивающей грамматикой, контекстно-зависим. Таким образом, неукорачивающие и КЗ-грамматики определяют один и тот же класс языков.

Тип 2

Грамматика $G = \langle T, N, P, S \rangle$ называется *контекстно-свободной (КС)*, если каждое правило из P имеет вид $A \rightarrow \beta$, где $A \in N$, $\beta \in (T \cup N)^*$.

Заметим, что в КС-грамматиках допускаются правила с пустыми правыми частями.

Язык, порождаемый контекстно-свободной грамматикой, называется *контекстно-свободным языком*.

Грамматикой *типа 2* будем называть контекстно-свободную грамматику.

КС-грамматика может являться неукорачивающей, т.е. удовлетворять ограничениям неукорачивающей грамматики.

Утверждение 2. Для любой КС-грамматики G существует неукорачивающая КС-грамматика G' , такая что $L(G) = L(G')$.

Согласно утверждению 2, любой контекстно-свободный язык порождается неукорачивающей контекстно-свободной грамматикой. Как следует из определений, такая КС-грамматика может содержать не более одного правила с пустой правой частью, причем это правило имеет вид $S \rightarrow \epsilon$, где S — начальный символ, и при этом никакое правило из P не содержит S в своей правой части.

В разделе «Преобразования грамматик» будет приведен алгоритм, позволяющий устранить из любой КС-грамматики все правила с пустыми правыми частями (в получившейся грамматике может быть правило $S \rightarrow \epsilon$ в случае, если пустая цепочка принадлежит языку, при этом начальный символ S не будет входить в правые части правил).

Тип 3

Грамматика $G = \langle T, N, P, S \rangle$ называется *праволинейной*, если каждое правило из P имеет вид $A \rightarrow wB$ либо $A \rightarrow w$, где $A, B \in N$, $w \in T^*$.

Грамматика $G = \langle T, N, P, S \rangle$ называется *леволинейной*, если каждое правило из P имеет вид $A \rightarrow Bw$ либо $A \rightarrow w$, где $A, B \in N$, $w \in T^*$.

Утверждение 3. Пусть L — формальный язык. Следующие два утверждения эквивалентны:

- существует праволинейная грамматика G_1 , такая что $L = L(G_1)$;
- существует леволинейная грамматика G_2 , такая что $L = L(G_2)$.

Из утверждения 3 следует, что праволинейные и левوليнейные грамматики определяют один и тот же класс языков. Такие языки называются *регулярными*. Право- и левوليнейные грамматики также будем называть регулярными.

Регулярная грамматика является грамматикой *типа 3*.

Автоматной грамматикой называется праволинейная (левوليнейная) грамматика, такая, что каждое правило с непустой правой частью имеет вид: $A \rightarrow a$ либо $A \rightarrow aB$ (для левوليнейной, соответственно, $A \rightarrow a$ либо $A \rightarrow Ba$), где $A, B \in N$, $a \in T$.⁴⁾

Автоматная грамматика является более простой формой регулярной грамматики. Существует алгоритм, позволяющий по регулярной (право- или левوليнейной) грамматике построить соответствующую автоматную грамматику. Таким образом, любой регулярный язык может быть порожден автоматной грамматикой. Кроме того, существует алгоритм, позволяющий устранить из регулярной (автоматной) грамматики все ε -правила (кроме $S \rightarrow \varepsilon$ в случае, если пустая цепочка принадлежит языку; при этом S не будет встречаться в правых частях правил)⁵⁾. Поэтому справедливо:

Утверждение 4. Для любой регулярной (автоматной) грамматики G существует неукорачивающая регулярная (автоматная) грамматика G' , такая что $L(G) = L(G')$.

Иерархия Хомского

Утверждение 5. Справедливы следующие соотношения:

- любая регулярная грамматика является КС-грамматикой;
- любая неукорачивающая КС-грамматика является КЗ-грамматикой;
- любая неукорачивающая грамматика является грамматикой типа 0.

Утверждение 5 следует непосредственно из определений.

Рассматривая только неукорачивающие регулярные и неукорачивающие КС-грамматики, что согласно утверждениям 2 и 4 не ограничивает классы соответствующих языков, получаем следующую иерархию классов грамматик:

$$\text{неукорачивающие Регулярные} \subset \text{неукорачивающие КС} \subset \text{КЗ} \subset \text{Тип 0}$$

(Запись $A \subset B$ означает, что A является собственным подклассом класса B .)

Вместо класса КЗ-грамматик в данной иерархии можно указать более широкий класс неукорачивающих грамматик, которые, как известно, порождают те же языки, что и КЗ-грамматики.

⁴⁾ В разделе «Разбор по регулярным грамматикам» будет рассмотрен алгоритм построения по конечному автомату эквивалентной грамматики. Построенная алгоритмом грамматика будет автоматной.

⁵⁾ Алгоритм устранения правил с пустой правой частью для КС-грамматик, приведенный в разделе «Преобразования грамматик», пригоден также и для устранения ε -правил из регулярных и автоматных грамматик.

Утверждение 6. Справедливы следующие соотношения:

- каждый регулярный язык является КС-языком, но существуют КС-языки, которые не являются регулярными, например:

$$L = \{a^n b^n \mid n > 0\};$$

- каждый КС-язык является КЗ-языком, но существуют КЗ-языки, которые не являются КС-языками, например:

$$L = \{a^n b^n c^n \mid n > 0\};$$

- каждый КЗ-язык является языком типа 0 (т. е. рекурсивно перечислимым языком), но существуют языки типа 0, которые не являются КЗ-языками, например: язык, состоящий из записей самоприменимых алгоритмов Маркова в некотором алфавите.⁶⁾

Из утверждения 6 следует иерархия классов языков:

$$\text{Тип 3 (Регулярные)} \subset \text{Тип 2 (КС)} \subset \text{Тип 1 (КЗ)} \subset \text{Тип 0}$$



Рис. 1. Иерархия классов языков.

На рис. 1 иерархия Хомского для классов языков изображена в виде диаграммы Венна. Классы вкладываются друг в друга. Самый широкий класс языков (типа 0) содержит в себе все остальные классы.

Утверждение 7. Проблема «Можно ли язык, описанный грамматикой типа k ($k = 0, 1, 2$), описать грамматикой типа $k + 1$?» является алгоритмически неразрешимой.

Заметим, что для $k = 1, 2, 3$ язык типа k является также и языком типа $k - 1$ (согласно иерархии на рис.1, класс языков типа k является подклассом класса языков типа $k - 1$). Несмотря на то, что нет алгоритма, позволяющего по заданному описанию языка L (например, по грамматике), определить максимальное k , такое что L является языком типа k , при ответе на вопрос «Какого типа заданный язык L ?» в примерах и задачах будем указывать, если не оговорено иное, максимально возможное k (0, 1, 2, или 3) для заданного языка L .

Рассмотрим, например, язык $L_{a,b} = \{a, b\}$, состоящий из двух однобуквенных цепочек. Какого типа язык $L_{a,b}$? Ответ: типа 3, так как, во-первых, он порождается грамматикой

⁶⁾ Понятие записи нормального алгоритма Маркова определяется в [10].

$S \rightarrow a / b$ типа 3; во-вторых, не существует грамматики типа $k > 3$, которая бы порождала $L_{a,b}$ (т. е. 3 — это максимально возможное k , такое, что $L_{a,b}$ является языком типа k).

Отметим, что согласно иерархии языков, следующие утверждения также справедливы: « $L_{a,b}$ является языком типа 2», « $L_{a,b}$ является языком типа 1», « $L_{a,b}$ является языком типа 0». Но, с учетом вышесказанного, эти утверждения не считаются исчерпывающими ответами на вопрос «Какого типа язык $L_{a,b}$?».

Соглашение: в примерах и задачах для краткости вместо полной четверки $G = \langle T, N, P, S \rangle$ будем иногда выписывать только так называемую «схему» грамматики — правила вывода P , считая, что большие латинские буквы обозначают нетерминальные символы, начальный символ (цель) грамматики обязательно стоит в левой части первого правила, ε — пустая цепочка, все остальные символы — терминальные.

Задача.

Даны грамматики G_1 и G_2 :

G_1 :

$S \rightarrow 0A1$

$A \rightarrow 0A0$

$A \rightarrow \varepsilon$

G_2 :

$S \rightarrow aSb / \varepsilon$

(1) Какого типа язык $L(G_1)$?

(2) Какого типа язык $L(G_2)$?

Решение

1) Правила грамматики G_1 удовлетворяют ограничениям на правила КС-грамматик, но не удовлетворяют ограничениям регулярных грамматик, т. е. это грамматика типа 2, и поэтому $L(G_1)$ является языком типа 2 (следовательно, это язык и типа 1, и типа 0). Заметим, что $L(G_1) = \{00^n0^n1 \mid n \geq 0\} = \{0^{(2n+1)}1 \mid n \geq 0\}$ и этот язык является также языком типа 3, поскольку описывается следующей регулярной грамматикой:

$S \rightarrow 0A$

$A \rightarrow 0B / 1$

$B \rightarrow 0A$

Итак, максимальное k , для которого $L(G_1)$ является языком типа k , равно 3.

2) По виду правил G_2 является грамматикой типа 2 и не является грамматикой типа 3. Следовательно, $L(G_2)$ является языком типа 2. Является ли $L(G_2)$ языком типа 3? Ответ на этот вопрос отрицательный, так как не существует регулярной грамматики (т. е. грамматики типа 3), порождающей язык $L(G_2) = \{a^n b^n \mid n \geq 0\}$ ⁷⁾.

Итак, максимальное k , для которого $L(G_2)$ является языком типа k , равно 2.

Примеры грамматик и языков

Примеры грамматик и языков, иллюстрирующие классы иерархии Хомского, приведем в порядке, обратном классификации, т. е. будем идти от простого к сложному.

⁷⁾ Нерегулярность языка $\{a^n b^n \mid n \geq 0\}$ будет доказана в разделе «Разбор по регулярным грамматикам».

Регулярные

1) Грамматика $S \rightarrow aS \mid a$ является праволинейной (неукорачивающей) грамматикой. Она порождает регулярный язык $\{a^n \mid n > 0\}$. Этот язык может быть порожден и леволинейной грамматикой: $S \rightarrow Sa \mid a$. Заметим, что обе эти грамматики являются автоматными.

2) Грамматика $S \rightarrow aS \mid \varepsilon$ является праволинейной и порождает регулярный язык $\{a^n \mid n \geq 0\}$. Для любого регулярного языка существует неукорачивающая регулярная грамматика (см. утверждение 4). Построим неукорачивающую регулярную грамматику для данного языка:

$$\begin{aligned} S &\rightarrow aA \mid a \mid \varepsilon \\ A &\rightarrow a \mid aA \end{aligned}$$

В самом деле, правило с пустой правой частью может применяться только один раз и только на первом шаге вывода; остальные правила таковы, что их правая часть не короче левой, т. е. грамматика неукорачивающая.

3) Грамматика

$$\begin{aligned} S &\rightarrow A\perp \mid B\perp \\ A &\rightarrow a \mid Ba \\ B &\rightarrow b \mid Bb \mid Ab \end{aligned}$$

леволинейная; она порождает регулярный язык, состоящий из всех непустых цепочек в алфавите $\{a, b\}$, заканчивающихся символом \perp (маркер конца) и не содержащих подцепочку aa . То есть в цепочках этого языка символ a не может встречаться два раза подряд, хотя бы один символ b обязательно присутствует между любыми двумя a . С помощью теоретико-множественной формулы данный язык описывается так: $\{\omega\perp \mid \omega \in \{a, b\}^+, aa \not\subseteq \omega\}$.

Контекстно-свободные

4) Грамматика

$$\begin{aligned} S &\rightarrow aQb \mid accb \\ Q &\rightarrow cSc \end{aligned}$$

является контекстно-свободной (неукорачивающей) и порождает КС-язык $\{(ac)^n (cb)^n \mid n > 0\}$, который, как и встречавшийся нам ранее язык $\{a^n b^n \mid n \geq 0\}$, не является регулярным, т. е. не может быть порожден ни одной регулярной грамматикой.

5) Грамматика

$$S \rightarrow aSa \mid bSb \mid \varepsilon$$

порождает КС-язык $\{xx^R, x \in \{a, b\}^*\}$. Данный язык не является регулярным. Грамматика не удовлетворяет определению неукорачивающей, но для нее существует эквивалентная неукорачивающая грамматика (см. утверждение 2). Приведем такую грамматику:

$$\begin{aligned} S &\rightarrow A \mid \varepsilon \\ A &\rightarrow aAa \mid bAb \mid aa \mid bb \end{aligned}$$

Неукорачивающие и контекстно-зависимые

6) Грамматика:

$$\begin{aligned} S &\rightarrow aSBC / abC \\ CB &\rightarrow BC \\ bB &\rightarrow bb \\ bC &\rightarrow bc \\ cC &\rightarrow cc \end{aligned}$$

неукорачивающая и порождает язык $\{a^n b^n c^n \mid n > 0\}$, который является языком типа 1, но не является языком типа 2 (этот факт доказывается с помощью «леммы о разрастании цепочек КС-языка» [3, т.1]).

Правило $CB \rightarrow BC$ не удовлетворяет определению КЗ-грамматики. Однако, заменив это правило тремя новыми $CB \rightarrow CD$, $CD \rightarrow BD$, $BD \rightarrow BC$ (добавляем новый символ D в множество нетерминалов N), мы получим эквивалентную серию контекстно-зависимых правил, которые меняют местами символы C и B . Таким образом, получаем эквивалентную КЗ-грамматику:

$$\begin{aligned} S &\rightarrow aSBC / abC \\ CB &\rightarrow CD \\ CD &\rightarrow BD \\ BD &\rightarrow BC \\ bB &\rightarrow bb \\ bC &\rightarrow bc \\ cC &\rightarrow cc \end{aligned}$$

Без ограничений на вид правил (тип 0)

7) В грамматике типа 0

$$\begin{aligned} S &\rightarrow SS \\ SS &\rightarrow \varepsilon \end{aligned}$$

второе правило не удовлетворяет ограничениям неукорачивающей грамматики. Существует бесконечно много выводов в данной грамматике, однако порождаемый язык конечен и состоит из единственной цепочки: $\{\varepsilon\}$.

Следующая грамматика также не является неукорачивающей и порождает пустой язык, так как ни одна терминальная цепочка не выводится из S (для обозначения пустого языка используется \emptyset — знак пустого множества):

$$\begin{aligned} S &\rightarrow SS \\ SS &\rightarrow Sa / S \end{aligned}$$

Заметим, что языки $L_\varepsilon = \{\varepsilon\}$ и $L_\emptyset = \emptyset$ (пустой язык) могут быть описаны грамматиками типа 3. Кроме того, отметим, что $L_\varepsilon \neq L_\emptyset$, так как L_\emptyset не содержит цепочек вообще, а L_ε — язык, состоящий из одной цепочки (пустая цепочка ε по определению равноправна с остальными цепочками в алфавите).

8) Содержательные примеры грамматик, порождающих языки, не принадлежащие классу контекстно-зависимых (тип 1), не приводятся из-за их громоздкости. Ниже обсуждается лишь возможность построения такой грамматики для одного языка типа 0.

Как известно, языки типа 0 (рекурсивно перечислимые множества) можно задавать с помощью распознавателей: НАМ или МТ. Существуют алгоритмы, позволяющие по НАМ

или МТ построить эквивалентную грамматику типа 0, задающую тот же язык, что и исходный распознаватель.

Пусть задан некоторый алфавит. Из теории алгоритмов известно, что можно построить НАМ A , на вход которому подается запись любого алгоритма (НАМ) X в заданном алфавите, и алгоритм A эмулирует работу алгоритма X в применении к записи X . Можно также добиться, чтобы A заикливался, если ему подали на вход цепочку, не являющуюся правильной записью какого-либо алгоритма. Таким образом, алгоритм A останавливается только на записях самоприменимых алгоритмов, а на всех других входах — заикливается. Другими словами, A задает язык записей самоприменимых алгоритмов, — обозначим этот язык L_{self} . Следовательно, L_{self} можно описать с помощью грамматики типа 0.

Покажем, что не существует грамматики типа 1, которая порождает язык L_{self} . Известно, что проблема распознавания для грамматик типа 1 алгоритмически разрешима. То есть существует алгоритм, который определяет, принадлежит ли заданная цепочка языку, порождаемому грамматикой. Предположим, существует грамматика типа 1, порождающая язык L_{self} . Тогда алгоритм распознавания дает ответ для любой цепочки, является ли она записью самоприменимого алгоритма или нет. Имеем противоречие с известным фактом об алгоритмической неразрешимости проблемы самоприменимости. Следовательно, грамматики типа 1, порождающей язык L_{self} , не существует.

Замечание о связи между языком и грамматикой

В приведенных ранее примерах мы утверждали, что заданная грамматика порождает определенный язык исходя из нестрогих рассуждений, интуитивно, по совокупности правил вывода. Вообще говоря, следует доказывать, что заданная грамматика порождает нужный язык. Для этого требуется доказать, что в данной грамматике:

- I) выводится любая цепочка, принадлежащая языку,
- II) не выводятся никакие другие цепочки.

Задача.

Доказать, что грамматика G с правилами вывода:

- (1, 2) $S \rightarrow aSBC \mid abC$
- (3) $CB \rightarrow BC$
- (4) $bB \rightarrow bb$
- (5) $bC \rightarrow bc$
- (6) $cC \rightarrow cc$

порождает язык $L(G) = \{ a^n b^n c^n \mid n \geq 1 \}$.

(В скобках слева приведена нумерация для удобства ссылок на правила).

Решение

I) Приведем схемы порождения цепочек вида $a^n b^n c^n$, $n \geq 1$ с указанием номера правила на каждом шаге вывода.

Для $n = 1$: $S \xrightarrow{(2)} abC \xrightarrow{(5)} abc$.

Для $n > 1$: $S \xrightarrow{(1)} aSBC \xrightarrow{(1)} aaSBCBC \rightarrow \dots \xrightarrow{(1)} a^{n-1}S(BC)^{n-1} \xrightarrow{(2)} a^n bC(BC)^{n-1} \xrightarrow{(3)} \dots \xrightarrow{(3)} a^n bB^{n-1}C^n \xrightarrow{(4)} a^n bbB^{n-2}C^n \rightarrow \dots \xrightarrow{(4)} a^n b^n C^n \xrightarrow{(5)} a^n b^n cC^{n-1} \xrightarrow{(6)} a^n b^n ccC^{n-2} \rightarrow \dots \xrightarrow{(6)} a^n b^n c^n$.

II) Из правил следует:

1. Новые символы a , $[b \mid B]$ и C появляются только при применении правил (1), (2) в равных количествах, т. е. в любой сентенциальной форме всегда **равное количество** a , $[b \mid B]$ и $[c \mid C]$.
2. Символ B заменяется только на b , а C — только на c .
3. Появившись, терминальные символы уже не меняют своей позиции, т. е. в любой сентенциальной форме символ a всегда **левее** любых $[b \mid B]$ и $[c \mid C]$.
4. Первый символ b появляется только после применения правила (2).
5. Символ B заменяется на b , только если слева от B стоит b , т. е. второй символ b появляется только непосредственно справа от первого b , третий b — непосредственно справа от второго b и т. д. Правило (5) применяется только после того, как исчерпана возможность применять (3), иначе вывод не будет завершен из-за наличия подцепочки cB в сентенциальной форме.

Из пунктов 3, 4 и 5 следует, что любой символ b расположен всегда **левее** любого $[c \mid C]$. Следовательно, в любой выводимой цепочке равное количество a , b и c ; a всегда стоит левее, чем b и c , b всегда стоит левее, чем c , т. е. любая цепочка имеет вид $a^n b^n c^n$, что и требовалось доказать.

Разбор цепочек

Цепочка в алфавите T принадлежит языку, порождаемому грамматикой $\langle T, N, P, S \rangle$, только в том случае, если существует ее вывод из начального символа S этой грамматики. Процесс построения такого вывода (а, следовательно, и определения принадлежности цепочки языку) называется *разбором*⁸⁾. Построение вывода можно осуществлять и в обратном порядке: в исходной цепочке ищем вхождение правой части некоторого правила и заменяем его на левую часть (это называется *сверткой*), в результате исходная цепочка «сворачивается» к некоторой сентенциальной форме, затем идет следующая свертка и т. д., пока не придем к цели грамматики — S . Процесс разбора называют также *анализом*.

С практической точки зрения наибольший интерес представляет разбор по **контекстно-свободным грамматикам**. Их порождающей мощности достаточно для описания большей части синтаксической структуры языков программирования, для различных подклассов КС-грамматик имеются хорошо разработанные практически приемлемые способы решения задачи разбора.

Рассмотрим основные понятия и определения, связанные с разбором по КС-грамматике.

Определение: вывод цепочки $\beta \in T^*$ из $S \in N$ в КС-грамматике $G = \langle T, N, P, S \rangle$, называется *левым (левосторонним)*, если в этом выводе каждая очередная сентенциальная форма получается из предыдущей заменой самого левого нетерминала.

⁸⁾ Разбором также называют и результат этого процесса, т. е. вывод цепочки, представленный (зафиксированный) каким-нибудь способом.

Определение: вывод цепочки $\beta \in T^*$ из $S \in N$ в КС-грамматике $G = \langle T, N, P, S \rangle$, называется *правым (правосторонним)*, если в этом выводе каждая очередная сентенциальная форма получается из предыдущей заменой самого правого нетерминала.

В грамматике для одной и той же цепочки может быть несколько выводов, эквивалентных в том смысле, что в них в одних и тех же местах применяются одни и те же правила вывода, но в различном порядке.

Например, для цепочки $a + b + a$ в грамматике:

$$G_{expr} = \langle \{a, b, +\}, \{S, T\}, \{S \rightarrow T \mid T + S; T \rightarrow a \mid b\}, S \rangle$$

можно построить выводы:

- (1) $S \rightarrow T + S \rightarrow T + T + S \rightarrow T + T + T \rightarrow a + T + T \rightarrow a + b + T \rightarrow a + b + a$
- (2) $S \rightarrow T + S \rightarrow a + S \rightarrow a + T + S \rightarrow a + b + S \rightarrow a + b + T \rightarrow a + b + a$
- (3) $S \rightarrow T + S \rightarrow T + T + S \rightarrow T + T + T \rightarrow T + T + a \rightarrow T + b + a \rightarrow a + b + a$

Здесь (2) — левосторонний вывод, (3) — правосторонний, а (1) не является ни левосторонним, ни правосторонним, но все эти выводы являются эквивалентными в указанном выше смысле.

Для КС-грамматик можно ввести удобное графическое представление вывода, называемое деревом вывода, причем для всех эквивалентных выводов дерева вывода совпадают.

Определение: ориентированное упорядоченное⁹⁾ дерево называется *деревом вывода* (или *деревом разбора*) в КС-грамматике $G = \langle T, N, P, S \rangle$, если выполнены следующие условия:

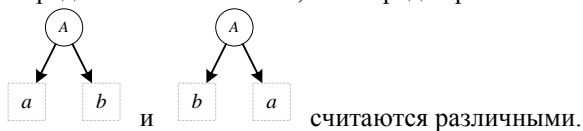
- (1) каждая вершина дерева помечена символом из множества $N \cup T \cup \{\varepsilon\}$, при этом корень дерева помечен символом S ; листья — символами из $T \cup \{\varepsilon\}$;
- (2) если вершина дерева помечена символом A , а ее непосредственные потомки — символами a_1, a_2, \dots, a_n , где каждое $a_i \in T \cup N$, то $A \rightarrow a_1 a_2 \dots a_n$ — правило вывода в этой грамматике;
- (3) если вершина дерева помечена символом A , а ее непосредственный потомок помечен символом ε , то этот потомок единственный и $A \rightarrow \varepsilon$ — правило вывода в этой грамматике.

На рисунке 2 изображен пример дерева для цепочки $a + b + a$ в грамматике G_{expr}

Определение: КС-грамматика G называется *неоднозначной*, если существует хотя бы одна цепочка $\alpha \in L(G)$, для которой может быть построено два или более различных деревьев вывода.

В противном случае грамматика называется *однозначной*.

⁹⁾ Упорядоченность означает, что порядок расположения потомков вершины существен. Например, деревья



Наличие двух или более деревьев вывода эквивалентно тому, что цепочка α имеет два или более разных левосторонних (или правосторонних) выводов.

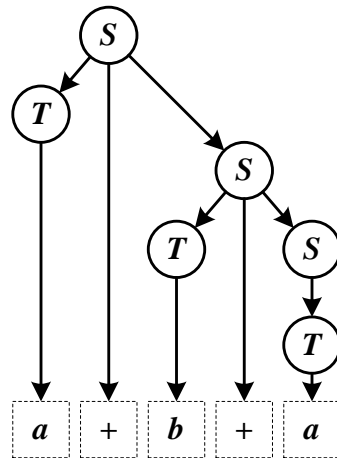


Рис. 2. Пример дерева вывода в грамматике G_{expr} .

Определение: язык, порождаемый грамматикой, называется *неоднозначным*, если он не может быть порожден никакой однозначной грамматикой.

Пример неоднозначной грамматики:

$$G_{if-then} = \langle \{if, then, else, a, b\}, \{S\}, P, S \rangle,$$

где $P = \{S \rightarrow if\ b\ then\ S\ else\ S \mid if\ b\ then\ S \mid a\}$.

В этой грамматике для цепочки *if b then if b then a else a* можно построить два различных дерева вывода, изображенных на рисунке 3 (а, б).

Однако это не означает, что язык $L(G_{if-then})$ обязательно неоднозначный. Обнаруженная в $G_{if-then}$ неоднозначность — это свойство грамматики, а не языка. Для некоторых неоднозначных грамматик существуют эквивалентные им однозначные грамматики.

Если грамматика используется для определения языка программирования, то она должна быть однозначной.

В приведенном выше примере разные деревья вывода предполагают соответствие *else* разным *then*. Если договориться, что *else* должно соответствовать ближайшему к нему *then*, и подправить грамматику $G_{if-then}$, то неоднозначность будет устранена:

$$S \rightarrow if\ b\ then\ S \mid if\ b\ then\ S' \ else\ S \mid a$$

$$S' \rightarrow if\ b\ then\ S' \ else\ S' \mid a$$

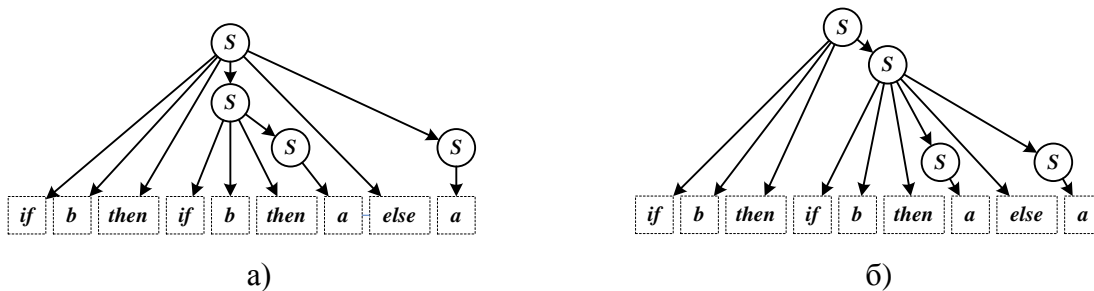


Рис. 3. Деревья вывода для «if b then if b then a else a» в грамматике $G_{if-then}$.

Утверждение 8. Проблема, порождает ли данная КС-грамматика однозначный язык (т.е. существует ли эквивалентная ей однозначная грамматика), является алгоритмически неразрешимой.

Более того, справедливо следующее.

Утверждение 9. Проблема, является ли данная КС-грамматика однозначной, алгоритмически неразрешима.

Поиск ответа на вопрос, неоднозначна или однозначна заданная грамматика — это искусство поиска цепочки с двумя различными деревьями вывода, или доказательство того, что таких цепочек не существует. Универсального способа решения этой задачи, к сожалению, нет.

Однако можно указать некоторые виды правил вывода, которые приводят к неоднозначности (при условии, что эти правила не являются тупиковыми¹⁰⁾, т. е. действительно используются на каком-нибудь шаге вывода терминальной цепочки из начального символа):

в приводимых схемах $\alpha, \beta, \gamma \in (T \cup N)^*$

$$(1) A \rightarrow AA \mid \alpha$$

$$(2) A \rightarrow A\alpha A \mid \beta$$

$$(3) A \rightarrow \alpha A \mid A\beta \mid \gamma \quad (\text{здесь хотя бы одна из цепочек } \alpha \text{ или } \beta \text{ не пуста})$$

$$(4) A \rightarrow \alpha A \mid \alpha A \beta A \mid \gamma$$

Отметим, что это всего лишь некоторые шаблоны. Все ситуации, приводящие к неоднозначности, перечислить невозможно в силу утверждения 9.

Пример неоднозначного КС-языка:

$$L = \{a^i b^j c^k \mid i, j, k \geq 0, i = j \text{ или } j = k\}.$$

Интуитивно неоднозначность L объясняется тем, что цепочки с $i = j$ должны порождаться группой правил вывода, отличных от правил, порождающих цепочки с $j = k$. Но тогда, по крайней мере, некоторые из цепочек с $i = j = k$ будут порождаться обеими группами правил и, следовательно, будут иметь по два разных дерева вывода. Доказательство того, что КС-язык L неоднозначный, приведено в [3, т.1, стр. 235–236].

Одна из грамматик, порождающих L , такова:

$$S \rightarrow AB \mid DC$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow bBc \mid \varepsilon$$

$$C \rightarrow cC \mid \varepsilon$$

$$D \rightarrow aDb \mid \varepsilon$$

Она неоднозначна; однозначных грамматик для L не существует.

Дерево вывода можно строить *нисходящим* либо *восходящим* способом.

При нисходящем разборе дерево вывода формируется от корня к листьям; на каждом шаге для вершины, помеченной нетерминальным символом, пытаются найти такое правило

¹⁰⁾ Как избавиться от правил, не участвующих в построении выводов, показано в разделе «Преобразования грамматик».

вывода, чтобы имеющиеся в нем терминальные символы проецировались на символы исходной (анализируемой) цепочки.

Метод восходящего разбора основан на обратном построении вывода с помощью сверток от исходной цепочки к цели грамматики S . При этом дерево растет снизу вверх — от листьев (символов анализируемой цепочки) к корню S . Если грамматика однозначная, то при любом способе построения будет получено одно и то же дерево разбора.

Преобразования грамматик

В некоторых случаях КС-грамматика может содержать бесполезные символы, которые не участвуют в порождении цепочек языка и поэтому могут быть удалены из грамматики.

Определение: символ $x \in T \cup N$ называется *недостижимым* в грамматике $G = \langle T, N, P, S \rangle$, если он не появляется ни в одной сентенциальной форме этой грамматики.

Алгоритм удаления недостижимых символов

Вход: КС-грамматика $G = \langle T, N, P, S \rangle$,

Выход: КС-грамматика $G' = \langle T', N', P', S \rangle$, не содержащая недостижимых символов, для которой $L(G) = L(G')$.

Метод:

1. $V_0 := \{S\}; i := 1$.

2. $V_i := V_{i-1} \cup \{x \mid x \in T \cup N, A \rightarrow \alpha x \beta \in P, A \in V_{i-1}, \alpha, \beta \in (T \cup N)^*\}$.

Если $V_i \neq V_{i-1}$, то $i := i + 1$ и переходим к шагу 2, иначе $N' := V_i \cap N$; $T' := V_i \cap T$; P' состоит из правил множества P , содержащих только символы из V_i ; $G' := \langle T', N', P', S \rangle$.

Определение: символ $A \in N$ называется *бесплодным* в грамматике $G = \langle T, N, P, S \rangle$, если множество $\{\alpha \in T^* \mid A \Rightarrow \alpha\}$ пусто.

Алгоритм удаления бесплодных символов

Вход: КС-грамматика $G = \langle T, N, P, S \rangle$.

Выход: КС-грамматика $G' = \langle T, N', P', S \rangle$, не содержащая бесплодных символов, для которой $L(G) = L(G')$.

Метод:

Строим множества N_0, N_1, \dots

1. $N_0 := \emptyset, i := 1$.

2. $N_i := N_{i-1} \cup \{A \mid A \rightarrow \alpha \in P \text{ и } \alpha \in (N_{i-1} \cup T)^*\}$.

Если $N_i \neq N_{i-1}$, то $i := i + 1$ и переходим к шагу 2, иначе $N' := N_i$; P' состоит из правил множества P , содержащих только символы из $N' \cup T$; $G' := \langle T, N', P', S \rangle$.

Определение: КС-грамматика G называется *приведенной*, если в ней нет недостижимых и бесплодных символов.

Алгоритм приведения грамматики

1. Обнаруживаются и удаляются все бесплодные нетерминалы.
2. Обнаруживаются и удаляются все недостижимые символы.

Удаление символов сопровождается удалением правил вывода, содержащих эти символы.¹¹⁾

Замечание

Если в этом алгоритме приведения поменять местами шаги (1) и (2), то не всегда результатом будет приведенная грамматика.

Для описания синтаксиса языков программирования стараются использовать однозначные приведенные КС-грамматики.

Некоторые применяемые на практике алгоритмы разбора по КС-грамматикам требуют, чтобы в грамматиках не было правил с пустой правой частью, т. е. чтобы КС-грамматика была неукорачивающей. Для любой КС-грамматики существует эквивалентная неукорачивающая (см. утверждение 2).

Ниже приводится алгоритм, позволяющий преобразовать любую КС-грамматику в неукорачивающую. На первом шаге алгоритма строится множество X , состоящее из нетерминалов грамматики, из которых выводима пустая цепочка. Построение этого множества можно провести по аналогии с шагами алгоритма удаления бесплодных символов: (1) $X_1 := \{ A \mid (A \rightarrow \varepsilon) \in P \}$; $i := 2$; (2) $X_i := \{ A \mid (A \rightarrow \alpha) \in P \text{ и } \alpha \in X_{i-1}^* \} \cup X_{i-1}$. Далее, пока $X_i \neq X_{i-1}$ увеличиваем i на единицу и повторяем (2). Последнее X_i — искомое множество X .

Алгоритм устранения правил с пустой правой частью

Вход: КС-грамматика $G = \langle T, N, P, S \rangle$.

Выход: КС-грамматика $G' = \langle T, N', P', S' \rangle$, G' — неукорачивающая, $L(G') = L(G)$.

Метод:

1. Построить множество $X = \{ A \in N \mid A \Rightarrow \varepsilon \}$; $N' := N$.
2. Построить P' , удалив из множества правил P все правила с пустой правой частью.
3. Если $S \in X$, то ввести новый начальный символ S' , добавив его в N' , и в множество правил P' добавить правило $S' \rightarrow S \mid \varepsilon$. Иначе просто переименовать S в S' .
4. Изменить P' следующим образом. Каждое правило вида $B \rightarrow \alpha_1 A_1 \alpha_2 A_2 \dots \alpha_n A_n \alpha_{n+1}$, где $A_i \in X$ для $i = 1, \dots, n$, $\alpha_i \in ((N' - X) \cup T)^*$ для $i = 1, \dots, n + 1$ (т. е. α_i — цепочка, не содержащая символов из X), заменить 2^n правилами, соответствующими всем возможным комбинациям вхождений A_i между α_i :

$$\begin{aligned}
 B &\rightarrow \alpha_1 \alpha_2 \dots \alpha_n \alpha_{n+1} \\
 B &\rightarrow \alpha_1 \alpha_2 \dots \alpha_n A_n \alpha_{n+1} \\
 &\dots \\
 B &\rightarrow \alpha_1 \alpha_2 A_2 \dots \alpha_n A_n \alpha_{n+1} \\
 B &\rightarrow \alpha_1 A_1 \alpha_2 A_2 \dots \alpha_n A_n \alpha_{n+1}
 \end{aligned}$$

¹¹⁾ Если начальный символ S — бесполезный в грамматике, то грамматика порождает пустой язык. Множество P после приведения такой грамматики будет пустым, однако S не следует удалять из N , так как алфавит N не может быть пустым и на последнем месте в четверке, задающей грамматику, должен стоять нетерминал — начальный символ.

Замечание

Если $\alpha_i = \varepsilon$ для всех $i = 1, \dots, n + 1$, то получившееся на данном шаге правило $B \rightarrow \varepsilon$ не включать в множество P' .

5. Удалить бесплодные и недостижимые символы и правила, их содержащие. (Кроме изначально имеющихся (в неприведенной грамматике), бесполезные символы могут возникнуть в результате шагов 2–4).

Замечание

Если применить данный алгоритм к регулярной (автоматной) грамматике, то результатом будет неукорачивающая регулярная (автоматная) грамматика.

Далее везде, если не оговорено иное, будем рассматривать только приведенные грамматики.

Элементы теории трансляции

Введение

В этом разделе будут рассмотрены некоторые алгоритмы и технические приемы, применяемые при построении трансляторов. Практически во всех трансляторах (и в компиляторах, и в интерпретаторах) в том или ином виде присутствует большая часть перечисленных ниже процессов:

- лексический анализ,
- синтаксический анализ,
- семантический анализ,
- генерация внутреннего представления программы,
- оптимизация,
- генерация объектной программы.

В конкретных компиляторах порядок этих процессов может быть несколько иным, какие-то из них могут объединяться в одну фазу, другие могут выполняться в течение всего процесса компиляции. В интерпретаторах и при смешанной стратегии трансляции часть этих процессов может вообще отсутствовать.

В данном разделе мы рассмотрим некоторые методы, используемые для построения анализаторов (лексического, синтаксического и семантического), язык промежуточного представления программы, способ генерации промежуточной программы, ее интерпретацию. Излагаемые алгоритмы и методы иллюстрируются на примере модельного паскалеподобного языка (М-языка). Приводится реализация в виде программы на Си++.

Информацию о других методах, алгоритмах и приемах, используемых при создании трансляторов, можно найти в [1, 2, 3, 4, 5, 8].

Разбор по регулярным грамматикам

Рассмотрим методы и средства, которые обычно используются при построении лексических анализаторов. В основе таких анализаторов лежат регулярные грамматики, поэтому рассмотрим грамматики этого класса более подробно.

Соглашение: в дальнейшем, если особо не оговорено, под регулярной грамматикой будем понимать левостороннюю автоматную грамматику $G = \langle T, N, P, S \rangle$ без пустых правых частей¹²⁾. Напомним, что в такой грамматике каждое правило из P имеет вид $A \rightarrow Bt$ либо $A \rightarrow t$, где $A, B \in N$, $t \in T$.

Соглашение: предположим, что анализируемая цепочка заканчивается специальным символом \perp — *признаком конца цепочки*.

Для грамматик этого типа существует алгоритм определения того, принадлежит ли анализируемая цепочка языку, порождаемому этой грамматикой (*алгоритм разбора*):

- (1) первый символ исходной цепочки $a_1a_2\dots a_n\perp$ заменяем нетерминалом A , для которого в грамматике есть правило вывода $A \rightarrow a_1$ (другими словами, производим свертку терминала a_1 к нетерминалу A)
- (2) затем многократно (до тех пор, пока не считаем признак конца цепочки) выполняем следующие шаги: полученный на предыдущем шаге нетерминал A и расположенный непосредственно справа от него очередной терминал a_i исходной цепочки заменяем нетерминалом B , для которого в грамматике есть правило вывода $B \rightarrow Aa_i$ ($i = 2, 3, \dots, n$);

Это эквивалентно построению дерева разбора методом снизу-вверх: на каждом шаге алгоритма строим один из уровней в дереве разбора, поднимаясь от листьев к корню.

При работе этого алгоритма возможны следующие ситуации:

- (1) прочитана вся цепочка; на каждом шаге находилась единственная нужная свертка; на последнем шаге свертка произошла к символу S . Это означает, что исходная цепочка $a_1a_2\dots a_n\perp \in L(G)$.
- (2) прочитана вся цепочка; на каждом шаге находилась единственная нужная свертка; на последнем шаге свертка произошла к символу, отличному от S . Это означает, что исходная цепочка $a_1a_2\dots a_n\perp \notin L(G)$.
- (3) на некотором шаге не нашлось нужной свертки, т. е. для полученного на предыдущем шаге нетерминала A и расположенного непосредственно справа от него очередного терминала a_i исходной цепочки не нашлось нетерминала B , для которого в грамматике было бы правило вывода $B \rightarrow Aa_i$. Это означает, что исходная цепочка $a_1a_2\dots a_n\perp \notin L(G)$.
- (4) на некотором шаге работы алгоритма оказалось, что есть более одной подходящей свертки, т. е. в грамматике разные нетерминалы имеют правила вывода с одинаковыми правыми частями, и поэтому непонятно, к какому из них производить свертку. Это говорит о *недетерминированности разбора*. Анализ этой ситуации будет дан ниже.

Допустим, что разбор на каждом шаге детерминированный.

¹²⁾ Полное отсутствие ϵ -правил в грамматике не позволяет описывать языки, содержащие пустую цепочку. Для наших целей в данном разделе это ограничение оправдано — мы будем применять автоматные грамматики для описания и разбора лексических единиц (лексем) языков программирования. Лексемы не могут быть пустыми.

Для того, чтобы быстрее находить правило с подходящей правой частью, зафиксируем все возможные свертки (это определяется только грамматикой и не зависит от вида анализируемой цепочки). Сделаем это в виде таблицы, столбцы которой помечены терминальными символами. Первая строка помечена символом H ($H \notin N$), а значение каждого элемента этой строки — это нетерминал, к которому можно свернуть помечающий столбец терминальный символ. Остальные строки помечены нетерминальными символами грамматики. Значение каждого элемента таблицы, начиная со второй строки — это нетерминальный символ, к которому можно свернуть пару «нетерминал-терминал», которыми помечены соответствующие строка и столбец.

Например, для левосторонней грамматики $G_{left} = \langle \{a, b, \perp\}, \{S, A, B, C\}, P, S \rangle$, где

P :

$$\begin{aligned} S &\rightarrow C\perp \\ C &\rightarrow Ab \mid Ba \\ A &\rightarrow a \mid Ca \\ B &\rightarrow b \mid Cb, \end{aligned}$$

такая таблица будет выглядеть следующим образом:

	a	b	\perp
H	A	B	—
C	A	B	S
A	—	C	—
B	C	—	—
S	—	—	—

Знак «—» ставится в том случае, если соответствующей свертки нет.

Но чаще информацию о возможных свертках представляют в виде *диаграммы состояний* (ДС) — неупорядоченного ориентированного помеченного графа, который строится следующим образом:

- (1) строим вершины графа, помеченные нетерминалами грамматики (для каждого нетерминала — одну вершину), и еще одну вершину, помеченную символом, отличным от нетерминальных, например, H . Эти вершины будем называть *состояниями*. H — начальное состояние.
- (2) соединяем эти состояния дугами по следующим правилам:
 - а) для каждого правила грамматики вида $W \rightarrow t$ соединяем дугой состояния H и W (от H к W) и помечаем дугу символом t ;
 - б) для каждого правила $W \rightarrow Vt$ соединяем дугой состояния V и W (от V к W) и помечаем дугу символом t ;

Диаграмма состояний для грамматики G_{left} (см. пример выше) изображена на рис.4:

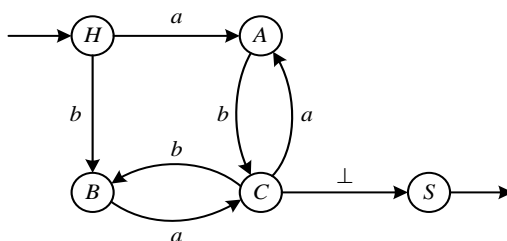


Рис. 4. Диаграмма состояний для грамматики G_{left} .

Алгоритм разбора по диаграмме состояний

- (1) объявляем текущим начальное состояние H ;
- (2) затем многократно (до тех пор, пока не считаем признак конца цепочки) выполняем следующие шаги: считываем очередной символ исходной цепочки и переходим из текущего состояния в другое состояние по дуге, помеченной этим символом. Состояние, в которое мы при этом попадаем, становится текущим.

При работе этого алгоритма возможны следующие ситуации (аналогичные ситуациям, которые возникают при разборе непосредственно по регулярной грамматике):

- 1) Прочитана вся цепочка; на каждом шаге находилась единственная дуга, помеченная очередным символом анализируемой цепочки; в результате последнего перехода оказались в состоянии S . Это означает, что исходная цепочка принадлежит $L(G)$.
- 2) Прочитана вся цепочка; на каждом шаге находилась единственная «нужная» дуга; в результате последнего шага оказались в состоянии, отличном от S . Это означает, что исходная цепочка не принадлежит $L(G)$.
- 3) На некотором шаге не нашлось дуги, выходящей из текущего состояния и помеченной очередным анализируемым символом. Это означает, что исходная цепочка не принадлежит $L(G)$.
- 4) На некотором шаге работы алгоритма оказалось, что есть несколько дуг, выходящих из текущего состояния, помеченных очередным анализируемым символом, но ведущих в разные состояния. Это говорит о *недетерминированности разбора*. Анализ этой ситуации будет приведен ниже.

Диаграмма состояний определяет конечный автомат, построенный по регулярной грамматике, который допускает множество цепочек, составляющих язык, определяемый этой грамматикой. Состояния и дуги ДС — это графическое изображение функции переходов конечного автомата из состояния в состояние при условии, что очередной анализируемый символ совпадает с символом-меткой дуги. Среди всех состояний выделяется начальное (считается, что в начальный момент своей работы автомат находится в этом состоянии) и заключительное (если автомат завершает работу переходом в это состояние, то анализируемая цепочка им допускается). На ДС эти состояния отмечаются короткими входящей и соответственно исходящей стрелками, не соединенными с другими вершинами.

Определение: *детерминированный конечный автомат (ДКА)* — это пятерка $\langle K, T, \delta, H, S \rangle$, где

K — конечное множество состояний;

T — конечное множество допустимых входных символов;

δ — отображение множества $K \times T$ в K , определяющее поведение автомата;

$H \in K$ — начальное состояние;

$S \in K$ — заключительное состояние (либо множество заключительных состояний $S \subset K$).

Замечания к определению ДКА:

- (1) Заключительных состояний в ДКА может быть более одного, однако для любого регулярного языка, все цепочки которого заканчиваются маркером конца (\perp), су-

существует ДКА с единственным заключительным состоянием. Заметим также, что ДКА, построенный по регулярной грамматике рассмотренным выше способом, всегда будет иметь единственное заключительное состояние S .¹³⁾

- (2) Отображение $\delta: K \times T \rightarrow K$ называют *функцией переходов* ДКА. $\delta(A, t) = B$ означает, что из состояния A по входному символу t происходит переход в состояние B . Иногда δ определяют лишь на подмножестве $K \times T$ (частичная функция). Если значение $\delta(A, t)$ не определено, то автомат не может дальше продолжать работу и останавливается в состоянии «ошибка».

Определение: ДКА допускает цепочку $a_1a_2\dots a_n$, если $\delta(H, a_1) = A_1$; $\delta(A_1, a_2) = A_2$; ... ; $\delta(A_{n-2}, a_{n-1}) = A_{n-1}$; $\delta(A_{n-1}, a_n) = S$, где $a_i \in T$, $A_j \in K$, $j = 1, 2, \dots, n-1$; $i = 1, 2, \dots, n$; H — начальное состояние, S — заключительное состояние.

Определение: множество цепочек, допускаемых ДКА, составляет определяемый им язык.

Для более удобной работы с диаграммами состояний введем несколько соглашений:

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

По диаграмме состояний легко написать анализатор для регулярной грамматики. Например, для грамматики $G_{left} = \langle \{a, b, \perp\}, \{S, A, B, C\}, P, S \rangle$ с правилами

P :

$$\begin{aligned} S &\rightarrow C\perp \\ C &\rightarrow Ab \mid Ba \\ A &\rightarrow a \mid Ca \\ B &\rightarrow b \mid Cb \end{aligned}$$

анализатор будет таким:

```
#include <iostream>
using namespace std;
char c; // текущий символ

void gc () { cin >> c; } // считать очередной символ из входной цепочки

bool scan_G ()
{
    enum state { H, A, B, C, S, ER }; // множество состояний
```

¹³⁾ Нетрудно указать и обратный способ — построения грамматики по диаграмме состояний автомата, — причем получившаяся грамматика будет автоматной. Каждой дуге из начального состояния H в состояние W , помеченной символом t , будет соответствовать правило $W \rightarrow t$; каждой дуге из состояния V в состояние W , помеченной символом t , будет соответствовать правило $W \rightarrow Vt$. Заключительное состояние S объявляется начальным символом грамматики.

Если в вершину H входит некоторая дуга (это возможно в произвольно построенном автомате), то алгоритм модифицируется так: каждой дуге из начального состояния H в состояние W , помеченной символом t , будет соответствовать правило $W \rightarrow Ht$, и в грамматику добавляется правило $H \rightarrow \varepsilon$; затем к построенной грамматике применяется алгоритм удаления правил с пустыми правыми частями.

```
state CS; // CS — текущее состояние
CS = H;

do
{ gc ();
  switch (CS)
  {
    case H: if ( c == 'a' )
            {
              CS = A;
            }
            else if ( c == 'b' )
            {
              CS = B;
            }
            else
              CS = ER;
            break;
    case A: if ( c == 'b' )
            {
              CS = C;
            }
            else
              CS = ER;
            break;
    case B: if ( c == 'a' )
            {
              CS = C;
            }
            else
              CS = ER;
            break;
    case C: if ( c == 'a' )
            {
              CS = A;
            }
            else if ( c == 'b' )
            {
              CS = B;
            }
            else if ( c == '⊥' )
              CS = S;
            else
              CS = ER;
            break;
  }
}

while ( CS != S && CS != ER);

return CS == S; // true, если CS != ER, иначе false
}
```

Пример разбора цепочки

Рассмотрим работу анализатора для грамматики G на примере цепочки $abba\perp$. При анализе данной цепочки получим следующую последовательность переходов в ДС:

$$H \xrightarrow{a} A \xrightarrow{b} C \xrightarrow{b} B \xrightarrow{a} C \xrightarrow{\perp} S$$

Вспомним, что каждый переход в ДС означает свертку сентенциальной формы путем замены в ней пары «нетерминал-терминал» Nt на нетерминал L , где $L \rightarrow Nt$ правило вывода в грамматике. Такое применение правила в обратную сторону будем записывать с помощью обратной стрелки $Nt \leftarrow L$ (обращение правила вывода). Тогда получим следующую последовательность сверток, соответствующую переходам в ДС:

$$abba\perp \leftarrow Abba\perp \leftarrow Cba\perp \leftarrow Ba\perp \leftarrow C\perp \leftarrow S$$

Эта последовательность не что иное, как обращение (правого) вывода цепочки $abba\perp$ в грамматике G . Она соответствует построению дерева снизу вверх (см. рис. 5).

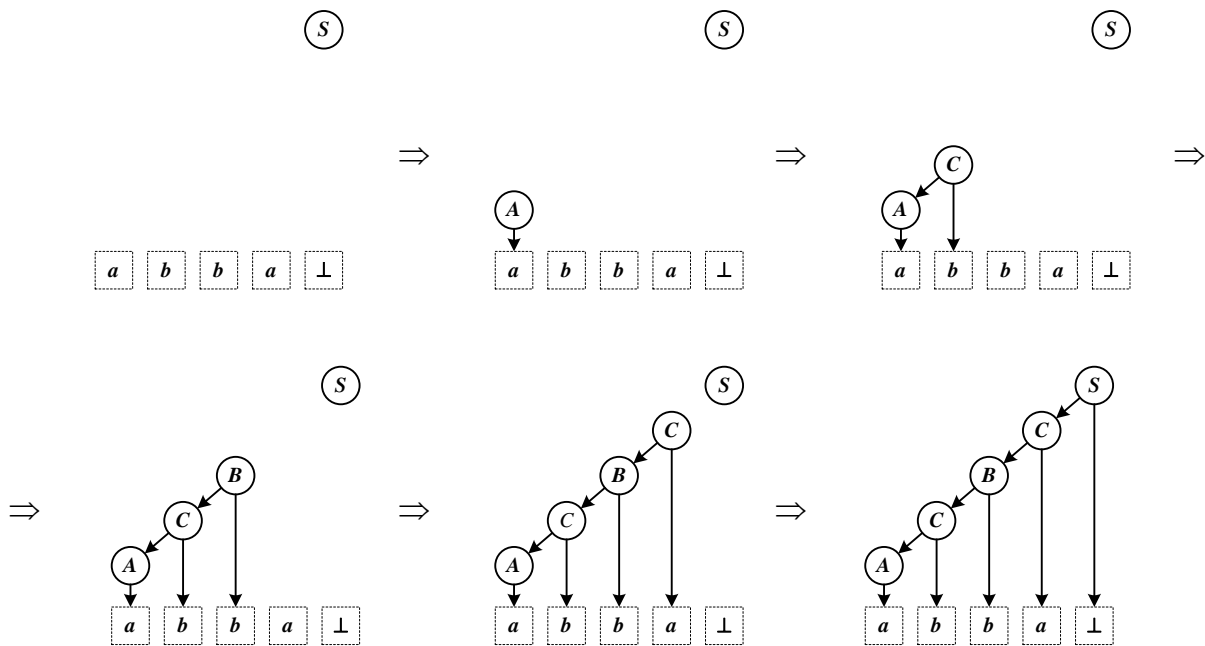


Рис. 5. Построение дерева вывода снизу вверх.

Разбор по праволинейной грамматике

По диаграмме состояний (см. рис. 4) построим праволинейную автоматную грамматику G_{right} следующим способом:

- нетерминалами будут состояния из ДС (кроме S);
- каждой дуге из состояния V в заключительное состояние S (помеченной знаком конца \perp) будет соответствовать правило $V \rightarrow \perp$;
- каждой дуге из состояния V в состояние W , помеченной символом t , будет соответствовать правило $V \rightarrow tW$;
- начальное состояние H объявляется начальным символом грамматики ¹⁴⁾.

¹⁴⁾ Нетрудно описать и обратный алгоритм для праволинейной автоматной грамматики, если все ее правила с односимвольной правой частью имеют вид $V \rightarrow \perp$. Состояниями ДС будут нетерминалы грамматики и еще одно специальное заключительное состояние S , в которое для каждого правила вида $V \rightarrow \perp$ проводится дуга

$$G_{right} = \langle \{a, b, \perp\}, \{H, A, B, C\}, P, H \rangle$$

P :

$$\begin{aligned} H &\rightarrow aA \mid bB \\ A &\rightarrow bC \\ C &\rightarrow bB \mid aA \mid \perp \\ B &\rightarrow aC \end{aligned}$$

Заметим, что $L(G_{right}) = L(G_{left})$, так как грамматики G_{right} и G_{left} соответствуют одной и той же ДС (см. рис. 4).

Рассмотрим разбор цепочки $abba\perp$ по праволинейной грамматике G_{right} . Последовательность переходов в ДС для этой цепочки такова:

$$H \xrightarrow{a} A \xrightarrow{b} C \xrightarrow{b} B \xrightarrow{a} C \xrightarrow{\perp} S$$

Каждый переход, за исключением последнего, означает теперь замену в сентенциальной форме нетерминала на пару «терминал-нетерминал» с помощью некоторого правила вывода грамматики G_{right} . В результате получаем следующий (левый) вывод, который соответствует последовательности переходов в ДС:

$$H \rightarrow aA \rightarrow abC \rightarrow abbB \rightarrow abbaC \rightarrow abba\perp$$

Такой вывод отражает построение дерева вывода сверху вниз (см. рис. 6).

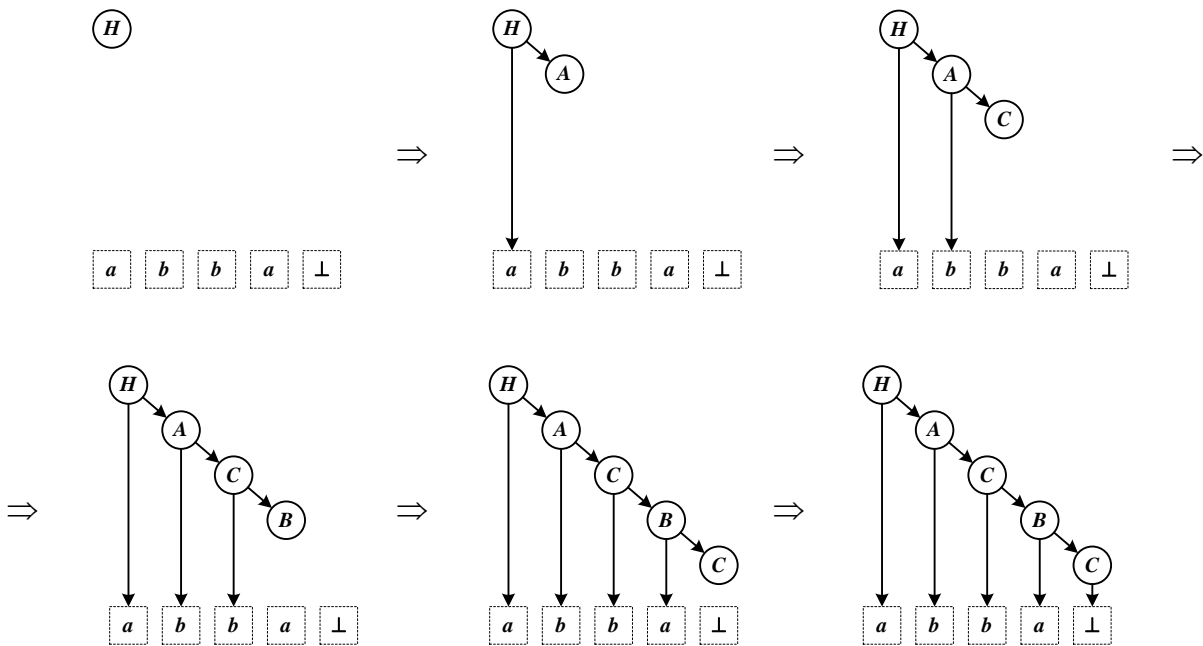


Рис. 6. Построение дерева вывода сверху вниз.

О недетерминированном разборе

При анализе по леволинейной грамматике может оказаться, что несколько нетерминалов имеют одинаковые правые части, и поэтому неясно, к какому из них делать свертку (см. ситуацию 4 в описании алгоритма). При анализе по праволинейной грамматике может ока-

из V , помеченная признаком конца \perp . Для каждого правила вида $V \rightarrow tW$ проводится дуга из V в W , помеченная символом t . Начальным состоянием в ДС будет начальный символ H .

заться, что нетерминал имеет две правые части с одинаковыми терминальными символами и поэтому неясно, на какую альтернативу заменить нетерминал. В терминах диаграммы состояний эти ситуации означают, что из одного состояния выходит несколько дуг, ведущих в разные состояния, но помеченных одним и тем же символом.

Например, для грамматики $G = \langle \{a, b, \perp\}, \{S, A, B\}, P, S \rangle$, где

P :

$S \rightarrow A\perp$

$A \rightarrow a | Bb$

$B \rightarrow b | Bb$

разбор будет недетерминированным (т.к. у нетерминалов A и B есть одинаковые правые части — Bb).

Такой грамматике будет соответствовать недетерминированный конечный автомат.

Определение: *недетерминированный конечный автомат (НКА)* — это пятерка $\langle K, T, \delta, H, S \rangle$, где

K — конечное множество состояний;

T — конечное множество допустимых входных символов;

δ — отображение множества $K \times T$ в множество подмножеств K ;

$H \subset K$ — конечное множество начальных состояний;

$S \subset K$ — конечное множество заключительных состояний.

$\delta(A, t) = \{B_1, B_2, \dots, B_n\}$ означает, что из состояния A по входному символу t можно осуществить переход в любое из состояний B_i , $i = 1, 2, \dots, n$. Также как и ДКА, любой НКА можно представить в виде таблицы (в одной ячейке такой таблицы можно указывать сразу несколько состояний, в которые возможен переход из заданного состояния по текущему символу) или в виде диаграммы состояний (ДС). В ДС каждому состоянию из множества K соответствует вершина; из вершины A в вершину B ведет дуга, помеченная символом t , если $B \in \delta(A, t)$. (В НКА из одной вершины могут исходить несколько дуг с одинаковой пометкой). *Успешный путь* — это путь из начальной вершины в заключительную; *пометка пути* — это последовательность пометок его дуг. *Язык*, допускаемый НКА, — это множество пометок всех успешных путей.

Если начальное состояние автомата (НКА или ДКА) одновременно является и заключительным, то автомат допускает пустую цепочку ε .

Замечание

Автомат, построенный по регулярной грамматике без пустых правых частей, не допускает ε .

Для построения разбора по регулярной грамматике в недетерминированном случае можно предложить алгоритм, который будет перебирать все возможные варианты сверток (переходов) один за другим; если цепочка принадлежит языку, то будет найден успешный путь; если каждый из просмотренных вариантов завершится неудачей, то цепочка языку не принадлежит. Однако такой алгоритм практически неприемлем, поскольку при переборе вариантов мы, скорее всего, снова окажемся перед проблемой выбора и, следовательно, будем иметь «дерево отложенных вариантов» и экспоненциальный рост сложности разбора.

Один из наиболее важных результатов теории конечных автоматов состоит в том, что класс языков, определяемых недетерминированными конечными автоматами, совпадает с классом языков, определяемых детерминированными конечными автоматами.

Утверждение 10. Пусть L — формальный язык. Следующие утверждения эквивалентны:

- (1) L порождается регулярной грамматикой;
- (2) L допускается ДКА;
- (3) L допускается НКА.

Эквивалентность пунктов (1) и (2) следует из рассмотренных выше алгоритмов построения конечного автомата по регулярной грамматике и обратно — грамматики по автомату. Очевидно, что из (2) следует (3): достаточно записать вместо каждого перехода ДКА $\delta(C, a) = b$ эквивалентный ему переход в НКА $\delta(C, a) = \{b\}$, начальное состояние ДКА поместить в множество начальных состояний НКА, а заключительное состояние ДКА поместить в множество заключительных состояний НКА. Приводимый ниже алгоритм построения ДКА, эквивалентного НКА, обосновывает то, что из (3) следует (2).

Алгоритм построения ДКА по НКА

Вход: НКА $M = \langle K, T, \delta, H, S \rangle$.

Выход: ДКА $M_1 = \langle K_1, T, \delta_1, H_1, S_1 \rangle$, допускающий тот же язык, что и автомат M : $L(M) = L(M_1)$.

Метод:

1. Элементами K_1 , т. е. состояниями в ДКА, будут некоторые подмножества множества состояний НКА. Заметим, что в силу конечности множества K , множество K_1 также конечно и имеет не более 2^s элементов, где s — мощность K .

Подмножество $\{A_1, A_2, \dots, A_n\}$ состояний из K будем для краткости записывать как $\underline{A_1A_2\dots A_n}$. Множество K_1 и переходы, определяющие функцию δ_1 , будем строить, начиная с состояния H_1 : $H_1 := \underline{A_1A_2\dots A_n}$, где $A_1, A_2, \dots, A_n \in H$. Другими словами, все начальные состояния НКА M объединяются в одно состояние H_1 для ДКА M_1 . Добавляем в множество K_1 построенное начальное состояние H_1 и пока считаем его нерассмотренным (на втором шаге оно рассматривается и строятся остальные состояния множества K_1 , а также переходы δ_1).

2. Пока в K_1 есть нерассмотренный элемент $\underline{A_1A_2\dots A_m}$, «рассматриваем» его и выполняем для каждого $t \in T$ следующие действия:

- Полагаем $\delta_1(\underline{A_1A_2\dots A_m}, t) = \underline{B_1B_2\dots B_k}$, где для $1 \leq j \leq k$ в НКА $\delta(A_i, t) = B_j$ для некоторых $1 \leq i \leq m$. Другими словами, $\underline{B_1B_2\dots B_k}$ — это множество всех состояний в НКА, куда можно перейти по символу t из множества состояний $\underline{A_1A_2\dots A_m}$. В ДКА M_1 получается детерминированный переход по символу t из состояния $\underline{A_1A_2\dots A_m}$ в состояние $\underline{B_1B_2\dots B_k}$. (Если $k = 0$, то полагаем $\delta_1(\underline{A_1A_2\dots A_m}, t) = \emptyset$).
- Добавляем в K_1 новое состояние $\underline{B_1B_2\dots B_k}$.

Шаг 2 завершается, поскольку множество новых состояний K_1 конечно.

3. Заключительными состояниями построенного ДКА M_1 объявляются все состояния, содержащие в себе хотя бы одно заключительное состояние НКА M : $S_1 := \{\underline{A_1A_2\dots A_m} \mid \underline{A_1A_2\dots A_m} \in K_1, A_i \in S \text{ для некоторых } 1 \leq i \leq m\}$.

Замечание

Множество S_1 построенного ДКА может состоять более, чем из одного элемента. Не для всех регулярных языков существует ДКА с единственным заключительным состоянием (пример: язык всех цепочек в алфавите $\{a, b\}$, содержащих не более двух символов b). Однако для реализации

алгоритма детерминированного разбора заключительное состояние должно быть единственным. В таком случае изменяют входной язык, добавляя маркер \perp в конец каждой цепочки (на практике в роли маркера конца цепочки \perp может выступать признак конца файла, символ конца строки или другие разделители). Вводится новое состояние S , и для каждого состояния Q из множества S_1 добавляется переход по символу \perp : $\delta_1(Q, \perp) = S$. Состояния из S_1 больше не считаются заключительными, а S объявляется единственным заключительным состоянием. Теперь по такому ДКА можно построить автоматную грамматику, допускающую детерминированный разбор.

Проиллюстрируем работу алгоритма на примерах.

Пример 1. Задан НКА $M = \langle \{ H, A, B, S \}, \{ 0, 1 \}, \delta, \{ H \}, \{ S \} \rangle$, где

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

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

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

$L(M) = \{ 1(01)^n \mid n \geq 1 \}$. Диаграмма для M изображена на рис.7.

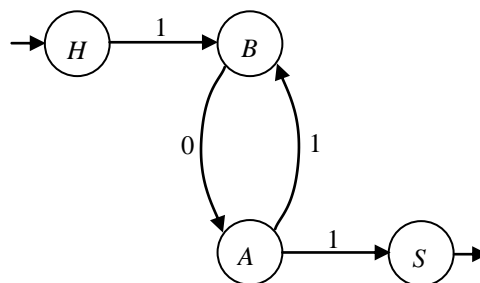


Рис. 7. ДС для автомата M .

Грамматика, соответствующая M :

$$S \rightarrow A1$$

$$A \rightarrow B0$$

$$B \rightarrow A1 \mid 1$$

Построим ДКА по НКА, пользуясь предложенным алгоритмом. Начальным состоянием будет \underline{H} .

$$\delta_1(\underline{H}, 1) = \underline{B}$$

$$\delta_1(\underline{B}, 0) = \underline{A}$$

$$\delta_1(\underline{A}, 1) = \underline{BS}$$

$$\delta_1(\underline{BS}, 0) = \underline{A}$$

Заключительным состоянием построенного ДКА является состояние \underline{BS} .

Таким образом, $M_1 = \langle \{ \underline{H}, \underline{B}, \underline{A}, \underline{BS} \}, \{ 0, 1 \}, \delta_1, \underline{H}, \underline{BS} \rangle$. Для удобства переименуем состояния в M_1 : \underline{BS} обозначается теперь как S_1 , а в однобуквенных именах состояний вместо подчеркивания используется индекс 1. Тогда $M_1 = \langle \{ H_1, B_1, A_1, S_1 \}, \{ 0, 1 \}, \{ \delta_1(H_1, 1) = B_1; \delta_1(B_1, 0) = A_1; \delta_1(A_1, 1) = S_1; \delta_1(S_1, 0) = A_1 \}, H_1, S_1 \rangle$.

Грамматика, соответствующая M_1 :

$$S_1 \rightarrow A_11$$

$$A_1 \rightarrow S_10 \mid B_10$$

$$B_1 \rightarrow 1$$

Построим диаграмму состояний (рис. 8).

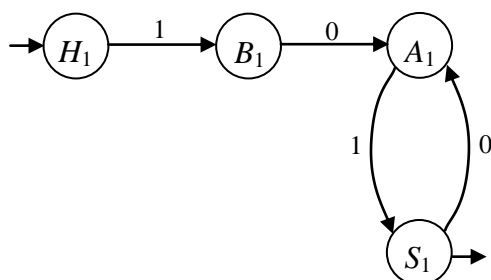


Рис. 8. Диаграмма состояний M_1 .

Пример 2. Дана регулярная грамматика $G = \langle T, N, P, S \rangle$

P :

$$S \rightarrow Sb \mid Aa \mid a$$

$$A \rightarrow Aa \mid Sb \mid b,$$

которой соответствует НКА. С помощью преобразования НКА в ДКА построить грамматику, по которой возможен детерминированный разбор.

Решение

Построим функцию переходов НКА:

$$\delta(H, a) = \{S\}$$

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

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

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

Процесс построения функции переходов ДКА удобно изобразить в виде таблицы, начав ее заполнение с начального состояния \underline{H} , добавляя строки для вновь появляющихся состояний:

символ \ состояние	a	b
\underline{H}	\underline{S}	\underline{A}
\underline{S}	\emptyset	\underline{AS}
\underline{A}	\underline{AS}	\emptyset
\underline{AS}	\underline{AS}	\underline{AS}

появились состояния \underline{S} и \underline{A} , они рассматриваются в следующих строках таблицы

Переходы ДКА:

$$\delta_1(\underline{H}, a) = \underline{S}$$

$$\delta_1(\underline{H}, b) = \underline{A}$$

$$\delta_1(\underline{S}, b) = \underline{AS}$$

$$\delta_1(\underline{A}, a) = \underline{AS}$$

$$\delta_1(\underline{AS}, a) = \underline{AS}$$

$$\delta_1(\underline{AS}, b) = \underline{AS}$$

Переобозначим состояния: $\underline{A} \Rightarrow A$, $\underline{H} \Rightarrow H$, $\underline{AS} \Rightarrow Y$, $\underline{S} \Rightarrow X$. Два заключительных состояния X и Y сведем в одно заключительное состояние S , используя маркер конца цепочки \perp . После такой модификации получаем ДКА с единственным заключительным состоянием S и функцией переходов:

$$\delta_1(H, a) = X$$

$$\delta_1(H, b) = A$$

$$\delta_1(X, b) = Y$$

$$\delta_1(X, \perp) = S$$

$$\delta_1(A, a) = Y$$

$$\delta_1(Y, a) = Y$$

$$\delta_1(Y, b) = Y$$

$$\delta_1(Y, \perp) = S$$

По ДКА строим грамматику G_1 , позволяющую воспользоваться алгоритмом детерминированного разбора:

G_1 :

$$S \rightarrow X\perp \mid Y\perp$$

$$Y \rightarrow Ya \mid Yb \mid Aa \mid Xb$$

$$X \rightarrow a$$

$$A \rightarrow b$$

В заключение раздела о регулярных языках приведем пример использования автоматов в решении теоретических задач.

Задача.

Доказать, что контекстно-свободный язык $L = \{ a^n b^n \mid n \geq 0 \}$ нерегулярен.

Доказательство (от противного).

Предположим, что язык L регулярен. Тогда существует конечный автомат (ДКА или НКА) $A = (K, \Sigma, \delta, I, F)$, допускающий язык $L : L(A) = L$. (Согласно утверждению 10, любой регулярный язык может быть задан конечным автоматом). Пусть число состояний автомата A равно k ($k > 0$). Рассмотрим цепочку $a^k b^k \in L$. Так как $L = L(A)$, $a^k b^k \in L(A)$. Это означает, что в автомате A есть успешный путь T с пометкой $a^k b^k$:

$$p_1 \xrightarrow{a} p_2 \xrightarrow{a} \dots \xrightarrow{a} p_k \xrightarrow{a} p_{k+1} \xrightarrow{b} p_{k+2} \xrightarrow{b} \dots \xrightarrow{b} p_{2k} \xrightarrow{b} p_{2k+1},$$

где $p_i \in K$ для $i = 1, \dots, 2k + 1$. Так как в автомате A всего k состояний, среди p_1, p_2, \dots, p_{k+1} найдутся хотя бы два одинаковых. Иными словами, существуют i, j , $1 \leq i < j \leq k$, такие что $p_i = p_j$. Таким образом, участок $p_i \xrightarrow{a} \dots \xrightarrow{a} p_j$ пути T является циклом. Пусть длина этого цикла равна m ($m > 0$, так как цикл — это непустой путь). Рассмотрим успешный путь T' , который отличается от T тем, что циклический участок $p_i \xrightarrow{a} \dots \xrightarrow{a} p_j$ присутствует в нем дважды:

$$p_1 \xrightarrow{a} \dots p_i \xrightarrow{a} \dots \xrightarrow{a} (p_i = p_j) \xrightarrow{a} \dots \xrightarrow{a} p_j \xrightarrow{a} \dots p_{k+1} \xrightarrow{b} \dots p_{2k+1} .$$

Пометкой пути T' является цепочка $a^{k+m}b^k$. Поскольку T' — успешный путь, $a^{k+m}b^k \in L(A)$. Так как $a^{k+m}b^k \notin L$, получаем, что $L \neq L(A)$. Это противоречит равенству $L = L(A)$. Следовательно, предположение о том, что L регулярен, неверно \square

Регулярные выражения

Кроме регулярных грамматик и конечных автоматов, существует еще один широко используемый в математических теориях и приложениях формализм, задающий регулярные языки. Это регулярные выражения. Они позволяют описать любой регулярный язык над заданным алфавитом, используя три вида операций: $+$ (объединение), \cdot (конкатенация), $*$ (итерация).¹⁵⁾

Определение. Пусть L, L_1, L_2 — языки над алфавитом Σ . Тогда будем называть язык $L_1 \cup L_2$ *объединением* языков L_1 и L_2 ; язык $L_1 \cdot L_2 = \{\varphi \cdot \psi \mid \varphi \in L_1, \psi \in L_2\}$ — *конкатенацией* (сцеплением) языков L_1 и L_2 (содержит всевозможные цепочки, полученные сцеплением цепочек из L_1 с цепочками из L_2); *i -ой степенью* языка L назовем язык $L^i = L^{i-1} \cdot L$ для $i > 0$, $L^0 = \{\varepsilon\}$. Язык $L^* = \bigcup_{n=0}^{\infty} L^n$ назовем *итерацией* языка L .

Определение. Пусть Σ — алфавит, не содержащий символов $*$, $+$, ε , \emptyset , $($, $)$. Определим рекурсивно *регулярное выражение* γ над алфавитом Σ и регулярный язык $L(\gamma)$, задаваемый этим выражением:

- 1) $a \in \Sigma \cup \{\varepsilon, \emptyset\}$ есть регулярное выражение; $L(a) = \{a\}$ для $a \in \Sigma \cup \{\varepsilon\}$; $L(\emptyset) = \emptyset$;
- 2) если α и β — регулярные выражения, то:
 - а) $(\alpha) + (\beta)$ — регулярное выражение; $L((\alpha) + (\beta)) = L(\alpha) \cup L(\beta)$;
 - б) $(\alpha)(\beta)$ — регулярное выражение; $L((\alpha)(\beta)) = L(\alpha)L(\beta)$;
 - в) $(\beta)^*$ — регулярное выражение; $L((\beta)^*) = (L(\beta))^*$;
- 3) все регулярные выражения конструируются только с помощью пунктов 1 и 2.

Если считать, что операция « $+$ » имеет самый низкий приоритет, а операция « $*$ » — наивысший, то в регулярных выражениях можно опускать лишние скобки.

Примеры регулярных выражений над алфавитом $\{a, b\}$:

$$a + b, \quad (a + b)^*, \quad (aa(ab)^*bb)^*.$$

Соответствующие языки:

$$L(a + b) = \{a\} \cup \{b\} = \{a, b\},$$

¹⁵⁾ Операцию итерации называют также звездочкой Клини, в честь ученого, предложившего регулярные выражения для описания регулярных языков. «Регулярность» в названии этого класса языков означает повторяемость некоторых фрагментов цепочек. На примере диаграмм конечных автоматов видно, что повторяемость фрагментов обусловлена наличием циклов в диаграмме.

$$L((a + b)^*) = \{a, b\}^*,$$

$$L((aa(ab)^*bb)^*) = \{\varepsilon\} \cup \{aa x_1 bb aa x_2 bb \dots x_k bb \mid k \geq 1, x_i \in \{(ab)^n \mid n \geq 0\} \text{ для } i = 1, \dots, k\}.$$

Существуют алгоритмы построения регулярного выражения по регулярной грамматике или конечному автомату и обратные алгоритмы, позволяющие преобразовать выражение в эквивалентную грамматику или автомат [3].

Регулярные выражения используются для описания лексем языков программирования, в задачах редактирования и обработки текстов; подходят для многих задач сравнения изображений и автоматического преобразования. Расширенные их спецификации (эквивалентные по описательной силе, но более удобные для практики) входят в промышленный стандарт открытых операционных систем *POSIX* и в состав базовых средств языка программирования *Perl*.

Задачи лексического анализа

Лексический анализ (ЛА) — это первый этап процесса компиляции. На этом этапе литеры, составляющие исходную программу, группируются в отдельные лексические элементы, называемые *лексемами*.

Лексический анализ важен для процесса компиляции по нескольким причинам:

- а) замена в программе идентификаторов, констант, ограничителей и служебных слов лексемами делает представление программы более удобным для дальнейшей обработки;
- б) лексический анализ уменьшает длину программы, устраняя из ее исходного представления несущественные пробельные символы и комментарии;
- в) если будет изменена кодировка в исходном представлении программы, то это отразится только на лексическом анализаторе.

Выбор конструкций, которые будут выделяться как отдельные лексемы, зависит от языка и от точки зрения разработчиков компилятора. Обычно принято выделять следующие типы лексем: идентификаторы, служебные слова, константы и ограничители. Каждой лексеме сопоставляется пара:

$$\langle \text{тип_лексемы}, \text{указатель_на_информацию_о_ней} \rangle.$$

Соглашение: эту пару тоже будем называть лексемой, если это не будет вызывать недоразумений.

Таким образом, лексический анализатор — это транслятор, входом которого служит цепочка символов, представляющих исходную программу, а выходом — последовательность лексем.

Очевидно, что лексемы перечисленных выше типов можно описать с помощью регулярных грамматик. Поскольку лексемы не пусты, для описания лексического состава любого языка программирования достаточно автоматных грамматик без ε -правил.

Например, идентификатор (*I*) описывается так:

$$I \rightarrow a \mid b \mid \dots \mid z \mid Ia \mid Ib \mid \dots \mid Iz \mid IO \mid IO \mid \dots \mid IO$$

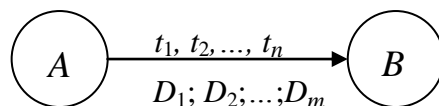
Целое без знака (N):

$$N \rightarrow 0 \mid 1 \mid \dots \mid 9 \mid N0 \mid N1 \mid \dots \mid N9 \quad \text{и т. д.}$$

Для грамматик этого класса, как мы уже видели, существует простой и эффективный алгоритм анализа того, принадлежит ли заданная цепочка языку, порождаемому этой грамматикой. Однако перед лексическим анализатором стоит более сложная задача:

- он должен сам выделить в исходном тексте цепочку символов, представляющую лексему, и проверить правильность ее записи;
- зафиксировать в специальных таблицах для хранения разных типов лексем факт появления соответствующих лексем в анализируемом тексте;
- преобразовать цепочку символов, представляющих лексему, в пару $\langle \text{тип_лексема}, \text{указатель_на_информацию_о_ней} \rangle$;
- удалить пробельные литеры и комментарии.

Для того чтобы решить эту задачу, опираясь на способ анализа с помощью диаграммы состояний, введем на дугах дополнительный вид пометок — пометки-действия. Теперь каждая дуга в ДС может выглядеть так:



Смысл t_i прежний: если в состоянии A очередной анализируемый символ совпадает с t_i для какого-либо $i = 1, 2, \dots, n$, то осуществляется переход в состояние B , при этом необходимо выполнить действия D_1, D_2, \dots, D_m .

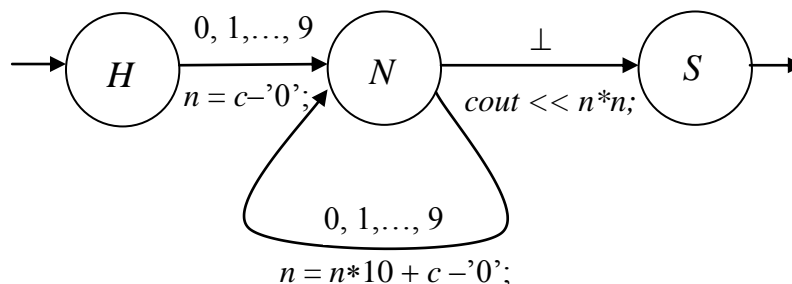
Задача. По заданной регулярной грамматке, описывающей целое число без знака, построить диаграмму с действиями по нахождению и печати квадрата данного числа.

$$S \rightarrow N \perp$$

$$N \rightarrow 0 \mid 1 \mid \dots \mid 9 \mid N0 \mid N1 \mid \dots \mid N9$$

Решение

Перевод числа во внутренне представление (переменная n типа int) будем осуществлять по схеме Горнера: распознав очередную цифру числа, умножаем текущее значение n на 10 и добавляем числовое значение текущей цифры (текущий символ хранится в переменной c типа $char$). Встретив маркер конца, печатаем квадрат числа n .



Лексический анализатор для М-языка

Описание модельного языка

$P \rightarrow \text{program } D_1; B \perp$
 $D_1 \rightarrow \text{var } D \{,D\}$
 $D \rightarrow I \{,I\}: [\text{int} \mid \text{bool}]$
 $B \rightarrow \text{begin } S \{;S\} \text{end}$
 $S \rightarrow I := E \mid \text{if } E \text{ then } S \text{ else } S \mid \text{while } E \text{ do } S \mid B \mid \text{read } (I) \mid \text{write } (E)$
 $E \rightarrow E_1 [= \mid < \mid > \mid !=] E_1 \mid E_1$
 $E_1 \rightarrow T \{ [+ \mid - \mid \text{or}] T \}$
 $T \rightarrow F \{ [* \mid / \mid \text{and}] F \}$
 $F \rightarrow I \mid N \mid L \mid \text{not } F \mid (E)$
 $L \rightarrow \text{true} \mid \text{false}$
 $I \rightarrow C \mid IC \mid IR$
 $N \rightarrow R \mid NR$
 $C \rightarrow a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z$
 $R \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

Замечания к грамматике модельного языка:

- а) запись вида $\{\alpha\}$ означает итерацию цепочки α , т. е. в порождаемой цепочке в этом месте может находиться либо ϵ , либо α , либо $\alpha\alpha$, либо $\alpha\alpha\alpha$ и т. д.;
- б) запись вида $[\alpha \mid \beta]$ означает, что в порождаемой цепочке в этом месте может находиться либо α , либо β ;
- в) P — цель грамматики; символ \perp — маркер конца текста программы.

Контекстные условия:

1. Любое имя, используемое в программе, должно быть описано и только один раз.
2. В операторе присваивания типы переменной и выражения должны совпадать.
3. В условном операторе и в операторе цикла в качестве условия возможно только логическое выражение.
4. Операнды операции отношения должны быть целочисленными.
5. Тип выражения и совместимость типов операндов в выражении определяются по обычным правилам; старшинство операций задано синтаксисом.

В любом месте программы, кроме идентификаторов, служебных слов и чисел, может находиться произвольное число пробелов и комментариев в фигурных скобках вида $\{ \langle \text{любые символы, кроме } \langle \rangle \text{ и } \langle \perp \rangle \rangle \}$.

true, false, read и *write* — служебные слова (в М-языке их нельзя переопределять в отличие от аналогичных стандартных идентификаторов в Паскале).

Разделители между идентификаторами, числами и служебными словами употребляются так же, как в Паскале.

Перейдем к построению лексического анализатора.

Вход лексического анализатора — символы исходной программы на М-языке; результат работы — исходная программа в виде последовательности лексем (их внутреннего представления). В нашем случае лексический анализатор будет вызываться, когда потребуется очередная лексема исходной программы, поэтому результатом работы лексического анализатора будет очередная лексема анализируемой программы на М-языке.

Лексический анализатор для модельного языка будем создавать поэтапно: сначала опишем на языке Си++ классы объектов, которые будут использоваться при ЛА, затем построим ДС с действиями для распознавания и формирования внутреннего представления лексем, а затем по ДС напишем программу ЛА. Отметим, что мы будем рассматривать один из возможных способов проектирования и реализации программы ЛА на Си++, быть может, не самый лучший.

Проектирование классовой структуры лексического анализатора

Представление лексем: выделим следующие типы лексем, введя следующий перечислимый тип данных:

```
enum type_of_lex
{
    LEX_NULL,                // 0
    LEX_AND, LEX_BEGIN, ... LEX_WRITE, // 18
    LEX_FIN,                // 19
    LEX_SEMICOLON, LEX_COMMA, ... LEX_GEQ, // 35
    LEX_NUM,                // 36
    LEX_ID,                 // 37
    POLIZ_LABEL,            // 38
    POLIZ_ADDRESS,         // 39
    POLIZ_GO,               // 40
    POLIZ_FGO               // 41
};
```

Содержательно внутреннее представление лексем — это пара (*тип лексемы, значение лексемы*). Значение лексемы — это номер строки в таблице лексем соответствующего класса, содержащей информацию о лексеме, или непосредственное значение, например, в случае целых чисел.

Соглашение об используемых таблицах лексем:

TW — таблица служебных слов М-языка;

TD — таблица ограничителей М-языка;

TID — таблица идентификаторов анализируемой программы;

Таблицы TW и TD заполняются заранее, т. к. их содержимое не зависит от исходной программы; TID формируется в процессе анализа.

Необходимые таблицы можно представить в виде объектов, конкретный вид которых мы рассмотрим чуть позже, или в виде массивов строк, например, для служебных слов.

Класс *Lex*:

```

class Lex
{
    type_of_lex  t_lex;
    int          v_lex;
public:
    Lex ( type_of_lex t = LEX_NULL, int v = 0)
    {
        t_lex = t; v_lex = v;
    }
    type_of_lex  get_type () { return t_lex; }
    int  get_value () { return v_lex; }
    friend ostream& operator << ( ostream &s, Lex l )
    {
        s << '(' << l.t_lex << ',' << l.v_lex << "));";
        return s;
    }
};

```

Класс *Ident*:

```

class Ident
{
    char    *   name;
    bool    declare;
    type_of_lex  type;
    bool    assign;
    int     value;

public:
    Ident ()
    {
        declare = false;
        assign   = false;
    }

    char *get_name ()
    {
        return name;
    }
    void put_name (const char *n)
    {
        name = new char [ strlen(n) + 1 ];
        strcpy ( name, n );
    }
    bool get_declare ()
    {
        return declare;
    }

    void put_declare ()
    {
        declare = true;
    }
};

```

```
type_of_lex get_type ()
{
    return type;
}
void put_type ( kind_of_lex t )
{
    type = t;
}
bool get_assign ()
{
    return assign;
}
void put_assign ()
{
    assign = true;
}
int get_value ()
{
    return value;
}
void put_value (int v)
{
    value = v;
}
};
```

Класс *tabl_ident* :

```
class tabl_ident
{
    ident * p;
    int size;
    int top;

public:
    tabl_ident ( int max_size )
    {
        p = new ident[size=max_size];
        top = 1;
    }
    ~tabl_ident ()
    {
        delete []p;
    }
    ident& operator[] ( int k )
    {
        return p[k];
    }
    int put ( const char *buf );
};

int tabl_ident::put ( const char *buf )
{
```



```

    for ( int j=1; j<top; ++j )
        if ( !strcmp(buf, p[j].get_name() )
            return j;
    p[top].put_name(buf);
    ++top;
    return top-1;
}

```

Класс *Scanner* :

```

class Scanner
{
    enum state { H, IDENT, NUMB, COM, ALE, DELIM, NEQ };
    static char * TW[];
    static type_of_lex words[];
    static char * TD[];
    static type_of_lex d_lms[];
    state CS;
    FILE * fp;
    char c;
    char buf[80];
    int buf_top;

    void clear ()
    {
        buf_top = 0;
        for ( int j = 0; j < 80; ++j )
            buf[j] = '\0';
    }

    void add ()
    {
        buf [ buf_top ++ ] = c;
    }

    int look ( const char *buf, char **list )
    {
        int i = 0;
        while ( list[i] )
        {
            if ( !strcmp(buf, list[i]) )
                return i;
            ++i;
        }
        return 0;
    }

    void gc ()
    {
        c = fgetc (fp);
    }

public:

```

```

Lex get_lex ();

Scanner ( const char * program )
{
    fp = fopen ( program, "r" );
    CS = H;
    clear();
    gc();
}
};

```

Таблицы лексем М-языка можно описать на Си++, например, таким образом:

```

char * Scanner::TW[] =
{
    ""           // 0 позиция 0 не используется
    "and",       // 1
    "begin",     // 2
    "bool",      // 3
    "do",        // 4
    "else",      // 5
    "end",       // 6
    "if",        // 7
    "false",     // 8
    "int",       // 9
    "not",       // 10
    "or",        // 11
    "program",   // 12
    "read",      // 13
    "then",      // 14
    "true",      // 15
    "var",       // 16
    "while",     // 17
    "write",     // 18
    NULL
};

char * Scanner::TD[] =
{
    ""           // 0 позиция 0 не используется
    "@",         // 1
    ";",         // 2
    ",",         // 3
    ":",         // 4
    ":",         // 5
    "(",         // 6
    ")",         // 7
    "=",         // 8
    "<",         // 9
    ">",         // 10
    "+",         // 11
    "-",         // 12
    "*",         // 13
    "/",         // 14
    "<=",        // 15

```

```
"!=",          // 16
">=",         // 17
NULL
};

tabl_ident TID(100);

type_of_lex Scanner::words[] =
{
    LEX_NULL,
    LEX_AND,
    LEX_BEGIN,
    LEX_BOOL,
    LEX_DO,
    LEX_ELSE,
    LEX_END,
    LEX_IF,
    LEX_FALSE,
    LEX_INT,
    LEX_NOT,
    LEX_OR,
    LEX_PROGRAM,
    LEX_READ,
    LEX_THEN,
    LEX_TRUE,
    LEX_VAR,

    LEX_WHILE,
    LEX_WRITE,
    LEX_NULL
};

type_of_lex Scanner::dlms[] =
{
    LEX_NULL,
    LEX_FIN,
    LEX_SEMICOLON,
    LEX_COMMA,
    LEX_COLON,
    LEX_ASSIGN,
    LEX_LPAREN,
    LEX_RPAREN,
    LEX_EQ,
    LEX_LSS,
    LEX_GTR,
    LEX_PLUS,
    LEX_MINUS,
    LEX_TIMES,
    LEX_SLASH,
    LEX_LEQ,
    LEX_NEQ,
    LEX_GEQ,
    LEX_NULL
};
```

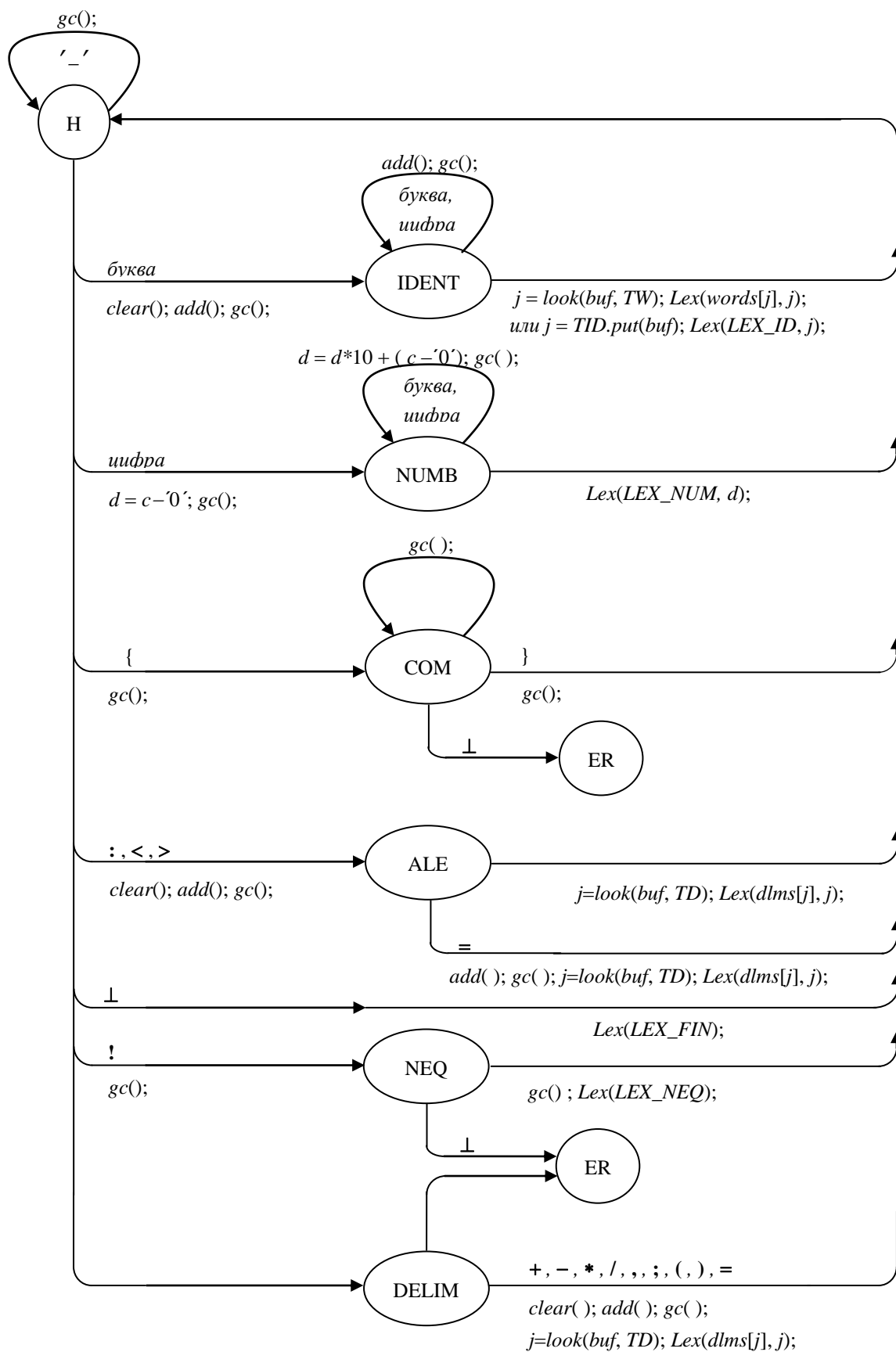


Рис. 9. ДС для модельного языка.

На рис. 9 приведена диаграмма состояний для модельного языка. Лексический анализ может проводиться независимо от последующих этапов трансляции. В этом случае исходный файл с текстом программы преобразуется в последовательность лексем во внутреннем представлении согласно построенной ДС; затем эта последовательность подается на вход синтаксическому анализатору.

Мы реализуем другой подход: лексический анализатор выдает очередную лексему по требованию синтаксического анализатора и затем «замирает», пока к нему не обратятся за следующей лексемой. При таком подходе действие $Lex(k, l)$; в ДС означает **return** $Lex(k, l)$; . Кроме того, вместо перехода в состояние ошибки ER генерируется исключение с указанием символа, который привел к ошибочной ситуации. Перехватываться и обрабатываться исключения будут не в лексическом анализаторе, а в основной программе, содержащей анализатор. Обработка исключения состоит в выводе сообщения об ошибке и корректном завершении программы.

Непосредственно реализует ЛА по ДС функция $get_lex()$:

```
Lex Scanner::get_lex ()
{
    int d, j;

    CS = H;
    do
    {
        switch ( CS )
        {
            case H:
                if ( c == ' ' || c == '\n' || c == '\r' || c == '\t' )
                    gc ();
                else if ( isalpha(c) )
                {
                    clear ();
                    add ();
                    gc ();
                    CS = IDENT;
                }
                else if ( isdigit(c) )
                {
                    d = c - '0';
                    gc ();
                    CS = NUMB;
                }
                else if ( c == '{' )
                {
                    gc ();
                    CS = COM;
                }
                else if ( c == ':' || c == '<' || c == '>' )
                {
                    clear ();
                    add ();
                    gc ();
                    CS = ALE;
                }
                else if ( c == '@' )
                    return Lex(LEX_FIN);
        }
    }
}
```

```
    else if ( c == '!' )
    {
        clear ();
        add ();
        gc ();
        CS = NEQ;
    }
    else
        CS = DELIM;
    break;
case IDENT:
    if ( isalpha(c) || isdigit(c) )
    {
        add ();
        gc ();
    }
    else
        if ( j = look (buf, TW) )
            return Lex (words[j], j);
        else
        {
            j = TID.put(buf);
            return Lex (LEX_ID, j);
        }
    break;
case NUMB:
    if ( isdigit(c) )
    {
        d = d * 10 + (c - '0');
        gc();
    }
    else
        return Lex ( LEX_NUM, d );
    break;
case COM:
    if ( c == '}' )
    {
        gc ();
        CS = H;
    }
    else if (c == '@' || c == '{' )
        throw c;
    else
        gc ();
    break;
case ALE:
    if ( c == '=' )
    {
        add ();
        gc ();
        j = look ( buf, TD );
        return Lex ( d\ms[j], j );
    }
    else
    {
```

```

        j = Look (buf, TD);
        return Lex ( d1ms[j], j );
    }
    break;
case NEQ:
    if ( c == '=' )
    {
        add ();
        gc ();
        j = look ( buf, TD );
        return Lex ( LEX_NEQ, j );
    }
    else
        throw '!';
    break;
case DELIM:
    clear ();
    add ();
    if ( j = look(buf, TD) )
    {
        gc ();
        return Lex ( d1ms[j], j );
    }
    else
        throw c;
    break;
} // end switch
}
while ( true );
}

```

Можно упростить реализацию метода *get_lex()*, вынеся обращение к *gc()* из тела **switch** и вставив его непосредственно перед оператором **switch**.

Синтаксический анализ

На этапе синтаксического анализа нужно:

- 1) установить, имеет ли цепочка лексем структуру, заданную синтаксисом языка, и
- 2) зафиксировать эту структуру.

Следовательно, снова надо решать задачу разбора: дана цепочка лексем, и надо определить, выводима ли она в грамматике, определяющей синтаксис языка. Если да, то построить вывод этой цепочки или дерево вывода. Однако структура таких конструкций как выражение, описание, оператор и т.п., более сложная, чем структура идентификаторов и чисел. Поэтому для описания синтаксиса языков программирования нужны более мощные грамматики, чем регулярные. Обычно для этого используют КС-грамматики, правила которых имеют вид $A \rightarrow \alpha$, где $A \in N$, $\alpha \in (T \cup N)^*$. Грамматики этого класса, с одной стороны, позволяют вполне адекватно описать синтаксис реальных языков программирования; с другой стороны, для разных подклассов КС-грамматик построены достаточно эффективные алгоритмы разбора.

Из теории синтаксического анализа известно, что существует алгоритм, который по любой данной КС-грамматике и данной цепочке выясняет, принадлежит ли цепочка языку, порождаемому этой грамматикой. Но время работы такого алгоритма (синтаксического анализа с возвратами) **экспоненциально** зависит от длины цепочки, что с практической точки зрения совершенно неприемлемо.

Существуют табличные методы анализа ([3]), применимые ко всему классу КС-грамматик и требующие для разбора цепочек длины n времени Cn^3 (алгоритм Кока-Янгера-Касами), где C — константа, либо Cn^2 (алгоритм Эрли). Их разумно применять только в том случае, если для интересующего нас языка не существует грамматики, по которой можно построить анализатор с линейной временной зависимостью от длины цепочки (такими языками могут быть, например, подмножества естественного языка).

При разработке языков программирования их синтаксис обычно стараются сделать таким, чтобы время на анализ было прямо пропорционально длине программы. Алгоритмы анализа, расходующие на обработку входной цепочки линейное время, применимы только к некоторым **подклассам** КС-грамматик.

Различные методы синтаксического анализа, или разбора, основываются на разных принципах, и используют различные техники построения дерева вывода. Каждый метод предполагает свой способ построения по грамматике программы-анализатора, которая будет осуществлять разбор цепочек. Корректный анализатор завершает свою работу для любой входной цепочки и выдает верный ответ о принадлежности цепочки языку. Анализатор *некорректен*, если:

- не распознает хотя бы одну цепочку, принадлежащую языку;
- распознает хотя бы одну цепочку, языку не принадлежащую;
- заикливается на какой-либо цепочке.

Говорят, что метод анализа *применим* к данной грамматике, если анализатор, построенный в соответствии с этим методом, корректен.

Рассмотрим один из фундаментальных методов разбора, применимый к некоторому подклассу КС-грамматик.

Метод рекурсивного спуска

Пример: пусть дана грамматика $G_1 = \langle \{a, b, c, d\}, \{S, A, B\}, P, S \rangle$, где

P :

$$\begin{aligned} S &\rightarrow ABd \\ A &\rightarrow a \mid cA \\ B &\rightarrow bA \end{aligned}$$

и надо определить, принадлежит ли цепочка $cabad$ языку $L(G_1)$. Построим левый вывод этой цепочки: $S \rightarrow ABd \rightarrow cABd \rightarrow caBd \rightarrow cabAd \rightarrow cabad$.

Следовательно, цепочка $cabad$ принадлежит языку $L(G_1)$.

Построение левого вывода эквивалентно построению дерева вывода методом сверху-вниз (нисходящим методом), при котором на очередном шаге раскрывается самый левый нетерминал в частично построенном дереве (рис. 10):

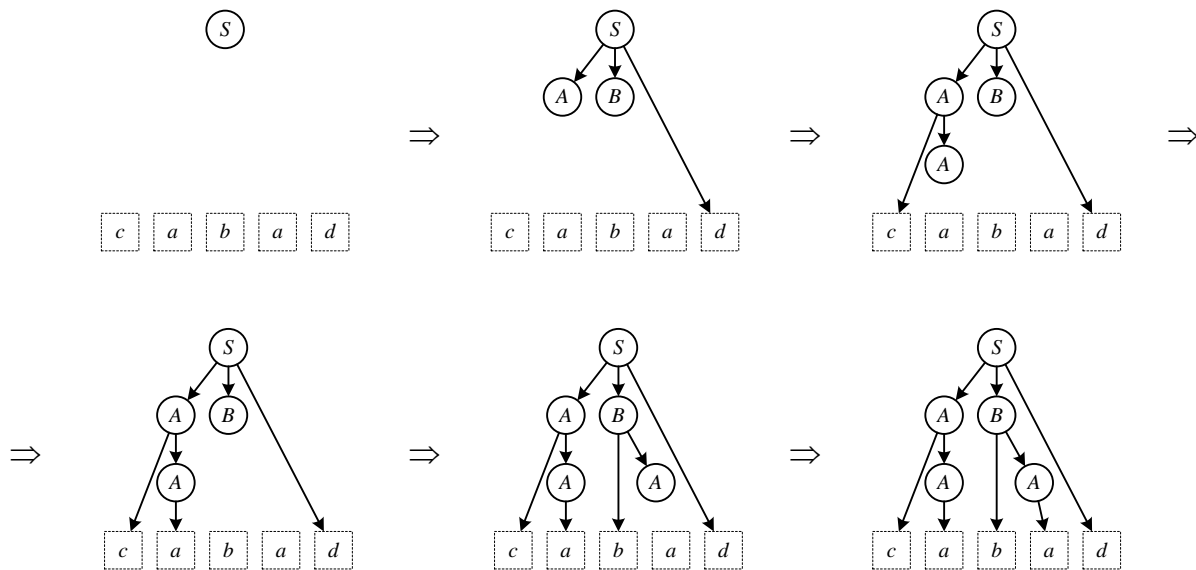


Рис. 10. Построение левого вывода.

Метод рекурсивного спуска (РС-метод) реализует разбор сверху-вниз и делает это с помощью системы рекурсивных процедур.

Для **каждого нетерминала** грамматики создается своя процедура, носящая **его имя**; ее задача — начиная с указанного места исходной цепочки найти подцепочку, которая выводится из этого нетерминала.

Если такую подцепочку найти не удастся, то процедура завершает свою работу, сигнализируя об ошибке. Это означает, что цепочка не принадлежит языку; разбор останавливается.

Если подцепочку удалось найти, то работа процедуры считается нормально завершённой и осуществляется возврат в точку вызова.

Тело каждой такой процедуры пишется непосредственно по правилам вывода (по альтернативам) соответствующего нетерминала: для правой части каждого правила осуществляется поиск подцепочки, выводимой из этой правой части. При этом терминалы из правой части распознаются самой процедурой, а нетерминалы соответствуют вызовам процедур, носящих их имена. После распознавания каждого терминала процедура считывает следующий символ из исходной цепочки, который становится текущим анализируемым символом. Выбор нужной альтернативы осуществляется процедурой по первому символу из еще нерассмотренной части исходной цепочки (т. е. по текущему символу).

Работа системы процедур начинается с главной функции *main()*. Она считывает первый символ исходной цепочки (заданной во входном потоке *stdin*) и вызывает процедуру *S()*, которая проверяет, выводится ли входная цепочка из начального символа *S* (в общем случае это делается с участием других процедур, которые, в свою очередь, рекурсивно могут вызывать и саму *S()* для анализа фрагмента исходной цепочки). Будем полагать, что в конце любой анализируемой цепочки всегда присутствует символ \perp (признак конца цепочки)¹⁶⁾, так что в задачу *main()* входит также распознавание символа \perp . Можно считать, что *main()* соответствует добавленному в грамматику правилу $M \rightarrow S\perp$, где *M* — новый начальный символ.

¹⁶⁾ На практике этим признаком может быть ситуация «конец файла» или «конец строки».

Пример. Совокупность процедур рекурсивного спуска для грамматики

G_1 :

$$\begin{aligned} S &\rightarrow ABd \\ A &\rightarrow a \mid cA \\ B &\rightarrow bA \end{aligned}$$

будет такой:

```
#include <iostream>
using namespace std;

int c;          // текущий анализируемый символ
void A ();
void B ();
void gc ()
{
    cin >> c; // считать очередной символ
}

void S ()
{
    cout << "S-->ABd, "; // применяемое правило вывода
    A();
    B();
    if ( c != 'd' )
        throw c;
    gc ();
}

void A ()
{
    if ( c == 'a' )
    {
        cout << "A-->a, ";
        gc ();
    }
    else if ( c == 'c' )
    {
        cout << "A-->cA, ";
        gc ();
        A ();
    }
    else
        throw c;
}

void B ()
{
    if ( c == 'b' )
    {
        cout << "B-->bA, ";
        gc ();
        A ();
    }
}
```

```

        else
            throw c;
    }

    int main ()
    {
        try
        {
            gc ();

            s ();
            if ( c != '⊥' )
                throw c;
            cout << "SUCCESS !!!" << endl;
            return 0;
        }
        catch ( int c )
        {
            cout << "ERROR on lexeme" << c << endl;
            return 1;
        }
    }

```

Для цепочки, выводимой из S , программа напечатает (помимо сообщения об успехе) последовательность правил, применяемых при нисходящем построении дерева вывода для данной цепочки (эта же последовательность годится для построения левого вывода). Вместо печати применяемых правил можно вставить действия по формированию дерева в динамической памяти в виде узлов, связанных указателями. Такое дерево может использоваться на последующих этапах трансляции.

Заметим, что даже если специально не фиксировать структуру анализируемой цепочки, система рекурсивных процедур все равно неявно обходит дерево вывода этой цепочки. Действительно, распознавание терминала b процедурой $B()$ соответствует в дереве вывода ветви B , а вызов процедуры $A()$ из процедуры $B()$ соответствует ветви $B \xrightarrow{A}$. Добавив в процедуры B анализ

дополнительные действия, можно наряду с проверкой синтаксиса определять смысл (семантику) входной цепочки. Например, смыслом арифметического выражения является его значение, и оно может быть вычислено в процессе неявного обхода дерева при разборе этого выражения системой рекурсивных процедур.

Выбор нужной альтернативы при анализе методом рекурсивного спуска легко осуществим, если все альтернативы начинаются с попарно различных терминальных символов. Сформулируем достаточное условие применимости метода рекурсивного спуска.

Достаточное условие применимости метода рекурсивного спуска

Для применимости метода рекурсивного спуска достаточно, чтобы каждое правило в грамматике удовлетворяло одному из двух видов:

(а) $X \rightarrow \alpha$,

где $\alpha \in (T \cup N)^*$ и это единственное правило вывода для этого нетерминала;

(б) $X \rightarrow a_1\alpha_1 \mid a_2\alpha_2 \mid \dots \mid a_n\alpha_n$,

где $a_i \in T$ для всех $i = 1, 2, \dots, n$; $a_i \neq a_j$ для $i \neq j$; $\alpha_i \in (T \cup N)^*$, т. е. если для нетерминала X правил вывода несколько, то они должны начинаться с терминалов, причем все эти терминалы должны быть попарно различными;

Это условие не является необходимым. Грамматику, удовлетворяющую данному условию, называют *s-грамматикой*.

Метод рекурсивного спуска является одной из возможных реализаций нисходящего анализа с прогнозируемым выбором альтернатив. Прогнозируемый выбор означает, что по грамматике можно заранее предсказать, какую альтернативу нужно будет выбрать на очередном шаге вывода в соответствии с текущим символом (т.е. первым символом из еще не прочитанной части входной цепочки). Далее мы подробно рассмотрим этот подход и сформулируем критерий его применимости.

Нисходящий анализ с прогнозируемым выбором альтернатив

В процессе построения левого вывода для произвольной цепочки в грамматике

G_1 :

$$\begin{aligned} S &\rightarrow ABd \\ A &\rightarrow a \mid cA \\ B &\rightarrow bA \end{aligned}$$

можно отметить следующее:

- (1) любой вывод начинается с применения правила $S \rightarrow ABd$;
- (2) если на очередном шаге сентенциальная форма имеет вид $wB\alpha$, где $w \in T^*$ — начало анализируемой цепочки, нетерминал B — самый левый в сентенциальной форме, то для продолжения вывода его нужно заменить на bA (других альтернатив нет);
- (3) если на очередном шаге сентенциальная форма имеет вид $wA\alpha$, где $w \in T^*$ — начало анализируемой цепочки, то выбор нужной альтернативы для замены A можно однозначно предсказать по тому, какой символ в анализируемой цепочке следует за начальной подцепочкой w : если символ a , то применяется альтернатива $A \rightarrow a$, если символ c , то альтернатива $A \rightarrow cA$; если какой-то иной символ — фиксируется ошибка: анализируемая цепочка не принадлежит языку $L(G_1)$;
- (4) если на каком-то шаге получилась сентенциальная форма вида $w\alpha$, отличная от (2) и (3), где w — максимально длинное начало, состоящее только из терминалов, то если α пуста и w совпадает с анализируемой цепочкой, процесс вывода успешно завершается, иначе фиксируется ошибка: анализируемая цепочка не принадлежит языку $L(G_1)$.

Отмеченные факты по поводу выбора нужной альтернативы на очередном шаге вывода в грамматике G_1 представим в виде так называемой таблицы прогнозов (или таблицы предсказаний):

	a	b	c	d
S	$S \rightarrow ABd$	$S \rightarrow ABd$	$S \rightarrow ABd$	$S \rightarrow ABd$
A	$A \rightarrow a$		$A \rightarrow cA$	
B		$B \rightarrow bA$		

Имея такую таблицу прогнозов (предсказаний) для КС-грамматики G , можно предложить следующий алгоритм нисходящего анализа (построение левого вывода):

1. Начать построение вывода с сентенциальной формы, состоящей из одного начального символа S .
2. Пока не будет получена цепочка, совпадающая с анализируемой, повторять следующие действия:

Пусть $wY\alpha$ — очередная сентенциальная форма, где $w \in T^*$. Если w не совпадает с началом анализируемой цепочки, то прекратить построение вывода и сообщить об ошибке: цепочка не принадлежит $L(G)$. В случае, когда w совпадает с началом, и следующим после w символом в анализируемой цепочке является символ z , заменить нетерминал Y на правую часть правила, которое находится в ячейке таблицы прогнозов на пересечении строки Y и столбца z . Если указанная ячейка пуста, прекратить построение вывода и сообщить об ошибке: цепочка не принадлежит $L(G)$.

Как мы видели выше, один из способов реализовать программу-анализатор для нисходящего анализа с прогнозируемым выбором альтернатив заключается в построении системы рекурсивных процедур¹⁷⁾. Это метод рекурсивного спуска.

Техника построения рекурсивных процедур уже была рассмотрена и продемонстрирована на примере, но остался открытым вопрос: как в общем случае «запрограммировать» процедуру на выбор нужной альтернативы по текущему символу. Ответ теперь известен — использовать таблицу прогнозов.

К сожалению, не для каждой КС-грамматики существует таблица с однозначными прогнозами, позволяющая безошибочно осуществить выбор альтернативы на каждом шаге вывода. В некоторых случаях заранее спрогнозировать выбор альтернативы невозможно: может оказаться, что подходящими в данной ситуации являются сразу несколько альтернатив (неоднозначный прогноз)¹⁸⁾.

Таким образом, нисходящий анализ с прогнозируемым выбором альтернатив пригоден лишь для некоторого подкласса КС-грамматик.

Метод рекурсивного спуска применим к грамматике, если для нее существует таблица однозначных прогнозов.

О применимости метода рекурсивного спуска

Очевидно, что метод рекурсивного спуска (без возвратов) неприменим к неоднозначным грамматикам. Например, по грамматике

$$G_2: \begin{array}{l} S \rightarrow aA / B \mid d \\ A \rightarrow d \mid aA \\ B \rightarrow aA \mid a \end{array}$$

¹⁷⁾ Другой способ, с явным использованием стека для хранения нетерминальной части сентенциальной формы, известен как $LL(1)$ -анализатор.

¹⁸⁾ Обобщение рассматриваемого метода — рекурсивный спуск с возвратами — «пробует» по очереди все возможные альтернативы в случае неоднозначного прогноза; о неудаче сообщается только в том случае, если ни одна из альтернатив не привела к успеху. Если грамматика неоднозначна, обобщенный метод построит все возможные деревья выводов. Мы не рассматриваем подробно рекурсивный спуск с возвратами, поскольку время, затрачиваемое на анализ при таком подходе, может экспоненциально зависеть от длины входной цепочки.

нельзя дать однозначный прогноз, что делать на первом шаге при анализе цепочки, начинающейся с символа a , т.е. по текущему символу a невозможно сделать однозначный выбор: $S \rightarrow aA$ или $S \rightarrow B$.

Как показывает следующий пример, грамматика может быть однозначной, однако однозначных прогнозов для нее также не существует:

$$G_3: \begin{array}{l} S \rightarrow A / B \\ A \rightarrow aA \mid d \\ B \rightarrow aB \mid b \end{array}$$

Действительно, каждая цепочка, выводимая в G_3 из S , оканчивается либо символом b , либо символом d , и имеет единственное дерево вывода. Но невозможно предсказать, с какой альтернативы ($S \rightarrow A$ или $S \rightarrow B$) начинать вывод, не просмотрев всю цепочку до конца и не увидев последний символ.

Наличие в грамматике нетерминала X с правилами вида $X \rightarrow \alpha$ и $X \rightarrow \beta$, из правых частей которых выводятся цепочки, начинающиеся одним и тем же терминалом a , т.е. $\alpha \Rightarrow a\alpha'$ и $\beta \Rightarrow a\beta'$, делает неоднозначным прогноз по символу a . Так что нисходящий анализ с прогнозируемым выбором альтернатив невозможен по такой грамматике, и метод рекурсивного спуска неприменим. Сформулируем это условие более формально.

Определение: множество $first(\alpha)$ в грамматике $G = \langle T, N, P, S \rangle$ — это множество терминальных символов, которыми начинаются цепочки, выводимые в G из цепочки $\alpha \in (T \cup N)^*$, т.е. $first(\alpha) = \{ a \in T \mid \alpha \Rightarrow a\alpha', \alpha' \in (T \cup N)^* \}$.

Например, для альтернатив правила $S \rightarrow A / B$ в грамматике G_3 имеем: $first(A) = \{ a, d \}$, $first(B) = \{ a, b \}$. Пересечение этих множеств непусто: $first(A) \cap first(B) = \{ a \} \neq \emptyset$, и поэтому метод рекурсивного спуска к G_3 неприменим.

Итак, наличие в грамматике правила с альтернативами $X \rightarrow \alpha \mid \beta$, такими что $first(\alpha) \cap first(\beta) \neq \emptyset$, делает метод рекурсивного спуска неприменимым.

Рассмотрим примеры грамматик с ϵ -правилами.

$$G_4: \begin{array}{l} S \rightarrow aA / BDc \\ A \rightarrow BAa \mid aB \mid b \\ B \rightarrow \epsilon \\ D \rightarrow B \mid b \end{array} \left| \begin{array}{l} first(aA) = \{ a \}, \quad first(BDc) = \{ b, c \}; \\ first(BAa) = \{ a, b \}, \quad first(aB) = \{ a \}, \quad first(b) = \{ b \}; \\ first(\epsilon) = \emptyset; \\ first(B) = \emptyset, \quad first(b) = \{ b \}. \end{array} \right.$$

Метод рекурсивного спуска неприменим к грамматике G_4 , так как $first(BAa) \cap first(aB) = \{ a \} \neq \emptyset$.

Следующий пример показывает еще одно свойство грамматик, наличие которого делает РС-метод неприменимым.

$$G_5: \begin{array}{l} S \rightarrow aA \\ A \rightarrow BC \mid B \\ C \rightarrow b \mid \epsilon \\ B \rightarrow \epsilon \end{array}$$

Пересечение множеств *first* пусто для любой пары альтернатив грамматики G_5 , однако наличие двух различных альтернатив, из которых выводится пустая цепочка, делает данную грамматику неоднозначной и, следовательно, метод рекурсивного спуска к ней неприменим. Действительно, $BC \Rightarrow \epsilon$ и $B \Rightarrow \epsilon$. Цепочка a имеет два различных дерева вывода :



Таким образом, если в грамматике для правила $X \rightarrow \alpha \mid \beta$ выполняются соотношения $\alpha \Rightarrow \epsilon$ и $\beta \Rightarrow \epsilon$, то метод рекурсивного спуска неприменим.

Осталось выяснить, как обстоят дела с применимостью метода, если для каждого нетерминала грамматики существует не более одной альтернативы, из которой выводится ϵ .

G_6 :

$S \rightarrow cAd \mid d$
 $A \rightarrow aA \mid \epsilon$

$first(cAd) = \{ c \}$, $first(d) = \{ d \}$;

Однозначные прогнозы для выбора альтернативы нетерминала S существуют, так как $first(cAd) \cap first(d) = \emptyset$.

Выбор альтернативы для A в данной грамматике также можно однозначно спрогнозировать: если текущим символом является a , применяется правило $A \rightarrow aA$, иначе правило $A \rightarrow \epsilon$. Это возможно благодаря тому, что за любой подцепочкой, выводимой из A , следует символ d , который сам в эту подцепочку не входит. Процедура $A()$ при выборе альтернативы $A \rightarrow \epsilon$ просто возвращает управление в точку вызова, не считывая следующий символ входной цепочки. Процедура $S()$, получив управление после вызова $A()$, проверяет, что текущим символом является d . Если это не так, фиксируется ошибка. Конечно, проверку символа d (без считывания следующего символа из входной цепочки) могла бы сделать и сама $A()$, но это излишне, так как $S()$ все равно будет проверять d , и если вместо d обнаружит другой символ, ошибка будет зафиксирована. Таблица прогнозов для G_6 :

	a	c	d
S		$S \rightarrow cAd$	$S \rightarrow d$
A	$A \rightarrow aA$	$A \rightarrow \epsilon$	$A \rightarrow \epsilon$

Итак, для грамматики G_6 , имеющей для каждого нетерминала не более одной альтернативы, из которой выводится пустая цепочка, метод рекурсивного спуска применим. Процедура $A()$ для нетерминала A , имеющего пустую альтернативу в грамматике G_6 , реализуется так:

```
void A ()
{
    if ( c == 'a' )
```

```

    {
        cout << "A->aA, ";
        gc ();
        A ();
    }
    else
    {
        cout << "A->epsilon, "; // след. символ не считывается
    }
}

```

Следующий пример показывает, что наличие альтернативы α , такой что $\alpha \Rightarrow \varepsilon$, все же может сделать метод рекурсивного спуска неприменимым.

G_7 :

$$\begin{aligned} S &\rightarrow Bd \\ B &\rightarrow cAa \mid a \\ A &\rightarrow aA \mid \varepsilon \end{aligned}$$

$$\text{first}(cAa) = \{ c \}, \text{first}(a) = \{ a \};$$

У нетерминала S правая часть единственна и проблема выбора альтернативы для S не стоит. Для выбора альтернативы нетерминала B существуют однозначные прогнозы, поскольку $\text{first}(cAa) \cap \text{first}(a) = \emptyset$.

Однако для нетерминала A прогноз по символу a неоднозначен. Дело в том, что любой вывод, содержащий A , имеет вид: $S \rightarrow Bd \rightarrow cAad \rightarrow \dots \rightarrow ca\dots aAad$. Поэтому альтернативу $A \rightarrow \varepsilon$ следует применять только тогда, когда текущим символом является a , а следующий за ним символ отличен от a (например, d). Если текущий — a и следующий за ним символ — тоже a , то выбирается альтернатива $A \rightarrow aA$. Но сделать однозначный выбор только по текущему символу в пользу какой-то одной из этих альтернатив невозможно, так как анализатор не умеет заглядывать вперед (в непрочитанную часть анализируемой цепочки).

Как видим, в G_7 существует сентенциальная форма, например $cAad$, в которой после нетерминала A , имеющего в грамматике пустую альтернативу, стоит символ a , с которого также начинается и непустая альтернатива для A . В таком случае процедура $A()$ не сможет правильно определить по текущему символу a , считать ли следующий символ и вызывать $A()$ (т. е. применять правило $A \rightarrow aA$) или возвращать управление без считывания символа (правило $A \rightarrow \varepsilon$). Опишем эту ситуацию более формально.

Определение: множество $\text{follow}(A)$ — это множество терминальных символов, которые могут появляться в сентенциальных формах грамматики $G = \langle T, N, P, S \rangle$ непосредственно справа от A (или от цепочек, выводимых из A), т.е.

$$\text{follow}(A) = \{ a \in T \mid S \Rightarrow \alpha A \beta, \beta \Rightarrow a \gamma, A \in N, \alpha, \beta, \gamma \in (T \cup N)^* \}.$$

Тогда, если в грамматике есть правило $X \rightarrow \alpha \mid \beta$, такое что $\beta \Rightarrow \varepsilon$, $\text{first}(\alpha) \cap \text{follow}(X) \neq \emptyset$, то метод рекурсивного спуска неприменим к данной грамматике.

Итак, на примерах мы рассмотрели все случаи, когда можно построить однозначные прогнозы по грамматике. Подытожив их, сформулируем **критерий применимости** метода рекурсивного спуска.

Утверждение 11. Пусть G — КС-грамматика. Метод рекурсивного спуска применим к G , если и только если для любой пары альтернатив $X \rightarrow \alpha \mid \beta$ выполняются следующие условия:

- (1) $first(\alpha) \cap first(\beta) = \emptyset$;
- (2) справедливо не более чем одно из двух соотношений: $\alpha \Rightarrow \varepsilon, \beta \Rightarrow \varepsilon$;
- (3) если $\beta \Rightarrow \varepsilon$, то $first(\alpha) \cap follow(X) = \emptyset$.

Рассмотрим грамматику

G_8 :

$S \rightarrow BDC$

$C \rightarrow Bd$

$D \rightarrow aB \mid d$

$B \rightarrow bB \mid \varepsilon$

$first(aB) = \{a\}, first(d) = \{d\}$;

$first(bB) = \{b\}$,

$follow(B) = \{a, b, d\}$, так как возможны следующие сентенциальные формы: $BdC, BaBbd$ ¹⁹⁾. Поскольку $first(bB) \cap follow(B) = \{b\} \neq \emptyset$, метод рекурсивного спуска неприменим к данной грамматике.

Естественно задаться вопросом: если грамматика не удовлетворяет критерию применимости метода рекурсивного спуска, то есть ли возможность построить эквивалентную грамматику, к которой данный метод применим.

Утверждение 12. Не существует алгоритма, определяющего для произвольной КС-грамматики, существует ли для нее эквивалентная грамматика, к которой метод рекурсивного спуска применим (т. е. это алгоритмически неразрешимая проблема²⁰⁾).

Построение таблицы прогнозов

Если критерий применимости метода рекурсивного спуска выполняется для грамматики G , то таблицу M однозначных прогнозов можно построить следующим образом:

1. Для каждого правила $X \rightarrow \alpha$ и для каждого терминала $a \in first(\alpha)$ помещаем $X \rightarrow \alpha$ в ячейку $M[X, a]$;
2. Для каждого правила $X \rightarrow \alpha$, такого что $\alpha \Rightarrow \varepsilon$, помещаем $X \rightarrow \alpha$ во все незаполненные на 1-м шаге ячейки строки X .
3. Для каждого правила $X \rightarrow Y\beta$, где $Y \in N$, $Y\beta$ — единственная альтернатива для X , помещаем $X \rightarrow Y\beta$ во все незаполненные на 1-м и 2-м шагах ячейки строки X .

На 2-м шаге могут быть заполнены ячейки для терминалов, не входящих в множество $follow(X)$, то есть соответствующих ошибочным ситуациям. Так как при анализе РС-

¹⁹⁾ В наших примерах мы вычисляем $first$ и $follow$ «интуитивно», опираясь на определения. Алгоритмы вычисления множеств $first$ и $follow$ можно найти в литературе (например, в [3]) или построить их самостоятельно.

²⁰⁾ Напомним, что алгоритмическая неразрешимость означает не то, что данную задачу нельзя решить для каждой конкретной грамматики, а то, что нет универсального способа решения, пригодного для всех грамматик.

методом считывания следующего символа в случае $X \Rightarrow \varepsilon$ не происходит, ошибка в текущей позиции обнаружится чуть позже, — той процедурой, которая анализирует текущий символ.

На 3-м шаге заполняются ячейки для терминалов, не входящих в $first(X)$, что также соответствует ошибочным ситуациям. Поскольку считывания следующего символа в случае единственной альтернативы $X \rightarrow Y\alpha$ в процедуре X не происходит, обнаружение ошибки произойдет позже, — в процедуре, анализирующей текущий символ.

Можно модифицировать построение таблицы прогнозов: третий шаг не выполнять вовсе (т.к. выбор единственной альтернативы уже осуществлен на шаге 1), на втором шаге каждое правило $X \rightarrow \alpha$, такое что $\alpha \Rightarrow \varepsilon$, помещать в ячейку $M[X, a]$ для каждого терминала $a \in follow(X)$ ²¹⁾; незаполненные на 1-м и 2-м шагах ячейки строки X оставить пустыми. Это позволит раньше обнаруживать ошибки в исходной цепочке, однако усложнит поведение самих процедур, так как им придется делать дополнительные проверки.

Пример. Построим таблицу прогнозов для грамматики

G_9 :

$S \rightarrow A / BS / cS$
 $B \rightarrow bB / d$
 $A \rightarrow aA \mid E / \varepsilon$
 $E \rightarrow e$

Вычислим необходимые для этого множества:

$first(A) = \{ a, e \}, first(BS) = \{ b, d \}, first(cS) = \{ c \}, follow(S) = \emptyset;$
 $first(bB) = \{ b \}, first(d) = \{ d \};$
 $first(aA) = \{ a \}, first(E) = \{ e \}, follow(A) = \emptyset;$
 $first(e) = \{ e \}.$

Как видно, множества $first$ для любой пары альтернатив не пересекаются, а также справедливы $first(A) \cap follow(A) = \emptyset$ и $first(S) \cap follow(S) = \emptyset$. Грамматика удовлетворяет критерию применимости метода рекурсивного спуска, и можно построить таблицу однозначных прогнозов:

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>
<i>S</i>	$S \rightarrow A$	$S \rightarrow BS$	$S \rightarrow cS$	$S \rightarrow BS$	$S \rightarrow A$
<i>A</i>	$A \rightarrow aA$	$A \rightarrow \varepsilon$	$A \rightarrow \varepsilon$	$A \rightarrow \varepsilon$	$A \rightarrow E$
<i>B</i>		$B \rightarrow bB$		$B \rightarrow d$	
<i>E</i>					$E \rightarrow e$

Построим для G_9 анализатор в виде системы рекурсивных процедур. Поведение процедур определяется полученной таблицей прогнозов. Заметим, что при выборе альтернативы, начинающейся с нетерминала, новый символ со входа не считывается, а сразу вызывается процедура, соответствующая этому нетерминалу.

```
#include <iostream>
using namespace std;
```

²¹⁾ Множество $follow(X)$ должно в этом случае содержать хотя бы один символ, — для этого предполагается, что в конце каждой входной цепочки языка есть маркер \perp .

```
int c;

void S ();
void A ();
void B ();
void E ();

void gc ()
{
    cin >> c; // считать очередной символ
}

void S ()
{
    if ( c =='a' || c =='e')
    {
        cout << "S-->A, "; // применяемое правило вывода
        // gc () не вызывается, текущий символ будет распознан процедурой A()
        A ();
    }
    else if ( c =='b' || c =='d')
    {
        cout << "S-->BS, ";
        // gc () не вызывается;
        B ();
        S ();
    }
    else if ( c =='c')
    {
        cout << "S-->cS, ";
        gc (); // символ 'c' распознан процедурой S(), считываем следующий
        S ();
    }
    else
        throw c; // недопустимый символ
}

void A ()
{
    if ( c =='a' )
    {
        cout << "A-->aA, ";
        gc ();
        A ();
    }
    else if ( c =='e' )
    {
        cout << "A-->E, ";
        // gc () не вызывается;
        E ();
    }
    else
    { // gc () не вызывается;
      cout << "A->epsilon, ";
    }
}

void B ()
{
    if ( c =='b' )
```

```
    {
        cout << "B-->bB, ";
        gc ();
        B ();
    }
    else if ( c == 'd' )
    {
        cout << "B-->d, ";
        gc ();
    }
    else
        throw c;
}

void E ()
{
    if ( c == 'e' )
    {
        cout << "E-->e, ";
        gc ();
    }
    else
        throw c;
}

int main ()
{
    try
    {
        gc ();
        S ();
        if ( c != '⊥' )
            throw c;
        cout << "SUCCESS !!!" << endl;
        return 0;
    }
    catch ( int c )
    {
        cout << "ERROR on lexeme" << c << endl;
        return 1;
    }
}
```

Рекурсивный спуск без построения прогнозов

Выделим подкласс грамматик, по которым можно строить систему рекурсивных процедур, минуя построение таблицы прогнозов.

Будем говорить, что правила КС-грамматики имеют *канонический* (для РС-метода) вид, если каждая группа правил с одинаковым нетерминалом в левой части относится к одному из перечисленных ниже видов и выполняются дополнительные условия:

- а) $X \rightarrow \alpha$,
где $\alpha \in (T \cup N)^*$ и это единственное правило вывода для этого нетерминала;
- б) $X \rightarrow a_1\alpha_1 \mid a_2\alpha_2 \mid \dots \mid a_n\alpha_n$,
где $a_i \in T$ для всех $i = 1, 2, \dots, n$; $a_i \neq a_j$ для $i \neq j$; $\alpha_i \in (T \cup N)^*$, т. е. если для нетерминала X правил вывода несколько, то они должны начинаться с терминалов, причем все эти терминалы попарно различны;
- в) $X \rightarrow a_1\alpha_1 \mid a_2\alpha_2 \mid \dots \mid a_n\alpha_n \mid \varepsilon$,
где $a_i \in T$ для всех $i = 1, 2, \dots, n$; $a_i \neq a_j$ для $i \neq j$; $\alpha_i \in (T \cup N)^*$, и $first(X) \cap follow(X) = \emptyset$.

Если правила вывода имеют такой вид, то рекурсивный спуск может быть реализован без промежуточного построения прогнозов: для правил с несколькими альтернативами выбирается альтернатива $a_i\alpha_i$, если текущий символ совпадает с a_i , иначе выбирается ε -альтернатива, если она присутствует; если нет совпадения текущего символа ни с одним из a_i и нет ε -альтернативы — фиксируется ошибка.

Канонический вид правил грамматики дает достаточное, но не необходимое условие применимости РС-метода.

Граматику с правилами канонического вида называют *q-грамматикой*. Рассмотренные выше *s-грамматики* являются *q-грамматиками*, обратное неверно.

Итерации в КС-грамматиках

При описании синтаксиса языков программирования часто встречаются правила, описывающие последовательность однотипных конструкций, отделенных друг от друга каким-либо знаком-разделителем (например, список идентификаторов при описании переменных, список параметров при вызове процедур и функций и т. п.). Такие правила обычно имеют вид: $L \rightarrow a \mid a, L$

Вместо обычных правил КС-грамматик для описания синтаксиса языков программирования нередко используют их модификации, например БНФ²²⁾. В БНФ, в частности, допускаются конструкции вида $\{\alpha\}$, где фигурные скобки — это спецсимволы для выделения фрагмента α , который может отсутствовать или повторяться произвольное число раз. Называют такую конструкцию *итерацией*.

С помощью итерации грамматику $L \rightarrow a \mid a, L$ можно переписать так: $L \rightarrow a \{, a\}$.

Наоборот, если в грамматике есть правило вида $X \rightarrow \alpha\{\beta\}\gamma$, содержащее итерацию $\{\beta\}$, то его можно заменить серией эквивалентных правил без итерации $\{\beta\}$:

$$\begin{aligned} X &\rightarrow \alpha Y \gamma \\ Y &\rightarrow \beta Y \mid \varepsilon, \end{aligned}$$

где Y — новый нетерминальный символ, добавляемый в алфавит нетерминалов грамматики.

Чтобы применить метод рекурсивного спуска для грамматики $L \rightarrow a \{, a\}$, преобразуем эту грамматику в эквивалентную без итерации:

$$\begin{aligned} L &\rightarrow a M \\ M &\rightarrow , a M \mid \varepsilon \end{aligned}$$

²²⁾ Бэкуса-Наура формула.

Метод применим к данной грамматике, так как $first(, a) \cap follow(M) = \emptyset$.

Можно построить систему рекурсивных процедур по преобразованной грамматике, но лучше, убедившись, что для преобразованной грамматики метод применим, вернуться к начальному варианту с итерацией.

Реализовать итерацию $\{\beta\}$ (где $\beta = b\beta'$, $b \in T$) удобно с помощью цикла: «пока текущий символ равен b , считать со входа следующий символ и проанализировать цепочку β' ».

Для правила с итерацией $L \rightarrow a \{, a\}$ процедура-анализатор реализуется так:

```

void L ()
{
    if ( c != 'a' )
        throw c;
    gc ();
    while ( c == ',' )
    {
        gc ();
        if ( c != 'a' )
            throw c;
        else
            gc ();
    }
}
    
```

Рассмотрим пример еще одной грамматики.

$G_{sequence}$:

$$S \rightarrow LB\perp$$

$$L \rightarrow a \{, a\}$$

$$B \rightarrow , b$$

Если для этой грамматики сразу написать анализатор, не убедившись в применимости метода к преобразованной (без итерации) грамматике, то цепочка a,a,a,b будет признана таким анализатором ошибочной, хотя в действительности a,a,a,b принадлежит порождаемому $G_{sequence}$ языку. Это произойдет потому, что процедура $L()$ захватит чужую запятую — ту, что выводится из B , и далее не обнаружив символ a , сообщит об ошибке.

Следует сначала проверить применимость метода, исключив из грамматики итерацию рассмотренным выше способом:

$$S \rightarrow LB\perp$$

$$L \rightarrow a M$$

$$M \rightarrow , a M \mid \varepsilon$$

$$B \rightarrow , b$$

Как видим, $first(, a) \cap follow(M) = \{, \} \neq \emptyset$ и поэтому метод рекурсивного спуска неприменим²³⁾. Можно попытаться поискать другую эквивалентную грамматику, к которой метод применим. Некоторые полезные для этого эквивалентные преобразования грамматик будут

²³⁾ Если бы в $G_{sequence}$ последовательность терминалов, перечисляемых через запятую, завершалась отличным от запятой символом (например, точкой с запятой), как это обычно и бывает в реальных языках программирования, то метод рекурсивного спуска был бы применим.

рассмотрены ниже. Однако, как следует из утверждения 12, успех в поиске эквивалентной грамматики, для которой метод применим, не гарантирован.

Преобразования грамматик

Если грамматика не удовлетворяет требованиям применимости метода рекурсивного спуска, то можно попытаться **преобразовать** ее, т. е. получить эквивалентную грамматику, пригодную для анализа этим методом.

(1) Если в грамматике есть нетерминалы, правила вывода которых **непосредственно леворекурсивны**, т. е. имеют вид:

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \dots \mid \beta_m,$$

где $\alpha_i \in (T \cup N)^+$ для $i = 1, 2, \dots, n$; $\beta_j \in (T \cup N)^*$ для $j = 1, 2, \dots, m$, то в таком случае применять метод рекурсивного спуска нельзя, поскольку $first(A\alpha_i) \cap first(A\alpha_k) \neq \emptyset$ для некоторых $i \neq k$, или $\beta_j = \varepsilon$ для некоторого j и $first(A\alpha_i) \cap follow(A) \neq \emptyset$ для $i = 1, 2, \dots, n$.

Левую рекурсию всегда можно заменить правой:

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \dots \mid \beta_m A' \\ A' &\rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \mid \varepsilon \end{aligned}$$

Будет получена грамматика, эквивалентная данной, т.к. из нетерминала A по-прежнему выводятся цепочки вида $\beta_j \{\alpha_i\}$, где $i = 1, 2, \dots, n$; $j = 1, 2, \dots, m$.

(2) Если в грамматике есть нетерминал, у которого несколько правил вывода начинаются **одинаковыми терминальными символами**, т. е. имеют вид

$$A \rightarrow a\alpha_1 \mid a\alpha_2 \mid \dots \mid a\alpha_n \mid \beta_1 \mid \dots \mid \beta_m,$$

где $a \in T$; $\alpha_i, \beta_j \in (T \cup N)^*$, β_j не начинается с a , $i = 1, 2, \dots, n$, $j = 1, 2, \dots, m$, то непосредственно применять метод рекурсивного спуска нельзя, т.к. $first(a\alpha_i) \cap first(a\alpha_k) \neq \emptyset$ для $i \neq k$. Можно преобразовать правила вывода данного нетерминала, объединив правила с общими началами в одно правило:

$$\begin{aligned} A &\rightarrow aA' \mid \beta_1 \mid \dots \mid \beta_m \\ A' &\rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n \end{aligned}$$

Будет получена грамматика, эквивалентная данной.

(3) Если в грамматике есть нетерминал, у которого **несколько** правил вывода, и среди них есть правила, **начинающиеся нетерминальными символами**, т. е. имеют вид:

$$\begin{aligned} A &\rightarrow B_1\alpha_1 \mid \dots \mid B_n\alpha_n \mid a_1\beta_1 \mid \dots \mid a_m\beta_m \\ B_1 &\rightarrow \gamma_{11} \mid \dots \mid \gamma_{1k} \\ &\dots \\ B_n &\rightarrow \gamma_{n1} \mid \dots \mid \gamma_{np}, \end{aligned}$$

где $B_i \in N$; $a_j \in T$; $\alpha_i, \beta_j, \gamma_{ij} \in (T \cup N)^*$, то можно заменить вхождения нетерминалов B_i их правилами вывода в надежде, что правила нетерминала A станут удовлетворять условиям применимости метода рекурсивного спуска:

$$A \rightarrow \gamma_{11}\alpha_1 \mid \dots \mid \gamma_{1k}\alpha_1 \mid \dots \mid \gamma_{n1}\alpha_n \mid \dots \mid \gamma_{np}\alpha_n \mid a_1\beta_1 \mid \dots \mid a_m\beta_m$$

(4) Если есть правила с пустой альтернативой вида:

$$A \rightarrow \alpha_1 A \mid \dots \mid \alpha_n A \mid \beta_1 \mid \dots \mid \beta_m \mid \varepsilon$$

$$B \rightarrow \alpha A \beta$$

и $first(A) \cap follow(A) \neq \emptyset$ (из-за вхождения A в правила вывода для B), то можно построить такую грамматику:

$$B \rightarrow \alpha A'$$

$$A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \mid \beta_1 \beta \mid \dots \mid \beta_m \beta \mid \beta$$

Полученная грамматика будет эквивалентна исходной, т. к. из B по-прежнему выводятся цепочки вида $\alpha \{ \alpha_i \} \beta_j \beta$ либо $\alpha \{ \alpha_i \} \beta$.

Однако правило вывода для нетерминального символа A' будет иметь альтернативы, начинающиеся одинаковыми терминальными символами (т. к. $first(A) \cap follow(A) \neq \emptyset$); следовательно, потребуются дальнейшие преобразования, и успех не гарантирован.

Пример. Рассмотрим грамматику G_{origin} :

G_{origin}

$$S \rightarrow fASd \mid \varepsilon$$

$$A \rightarrow Aa \mid Ab \mid dB \mid f$$

$$B \rightarrow bcB \mid \varepsilon$$

$$first(Aa) = first(Ab) = \{ d, f \}$$

$$first(bcB) = \{ b \}, follow(B) = \{ a, b, d, f \}$$

Условия применимости метода рекурсивного спуска не выполняются для G_{origin} . С помощью преобразований приведем эту грамматику к каноническому виду для рекурсивного спуска. Будем подчеркивать не удовлетворяющие каноническому виду правила и при переходе к новой грамматике указывать номер примененного преобразования $\Rightarrow_{(i)}$, $1 \leq i \leq 4$:

G_{origin} :

$$S \rightarrow fASd \mid \varepsilon$$

$$A \rightarrow \underline{A}a \mid \underline{A}b \mid dB \mid f$$

$$B \rightarrow bcB \mid \varepsilon$$

$\Rightarrow_{(1)}$

$G_{transform_1}$:

$$S \rightarrow fASd \mid \varepsilon$$

$$A \rightarrow dBA' \mid fA'$$

$$A' \rightarrow aA' \mid bA' \mid \varepsilon$$

$$B \rightarrow bcB \mid \underline{\varepsilon}$$

$\Rightarrow_{(4)}$

$$first(S) = \{ f \}, follow(S) = \{ d \}, first(S) \cap follow(S) = \emptyset;$$

$$first(A') = \{ a, b \}, follow(A') = \{ f, d \}, first(A') \cap follow(A') = \emptyset;$$

$$first(B) = \{ b \}, follow(B) = \{ a, b, f, d \}, first(B) \cap follow(B) = \{ b \} \neq \emptyset.$$

$G_{transform_2}$:

$$S \rightarrow fASd \mid \varepsilon$$

$$A \rightarrow dB' \mid fA'$$

$$B' \rightarrow bcB' \mid \underline{A'}$$

$$A' \rightarrow aA' \mid bA' \mid \varepsilon$$

$\Rightarrow_{(4)}$

$\Rightarrow_{(3)}$

$G_{transform_3}$:

$$S \rightarrow fASd \mid \varepsilon$$

$$A \rightarrow dB' \mid fA'$$

$$B' \rightarrow \underline{bcB'} \mid aA' \mid \underline{bA'} \mid \varepsilon$$

$$A' \rightarrow aA' \mid bA' \mid \varepsilon$$

$$\begin{array}{l} G_{transform_4}: \\ \Rightarrow(2) \quad S \rightarrow fASd \mid \varepsilon \\ \quad \quad A \rightarrow dB' \mid fA' \\ \quad \quad B' \rightarrow bC \mid aA' \mid \varepsilon \\ \quad \quad C \rightarrow cB' \mid \underline{A'} \\ \quad \quad A' \rightarrow aA' \mid bA' \mid \varepsilon \end{array} \quad \Rightarrow(3) \quad \begin{array}{l} G_{object}: \\ S \rightarrow fASd \mid \varepsilon \\ A \rightarrow dB' \mid fA' \\ B' \rightarrow bC \mid aA' \mid \varepsilon \\ C \rightarrow cB' \mid aA' \mid bA' \mid \varepsilon \\ A' \rightarrow aA' \mid bA' \mid \varepsilon \end{array}$$

$$\begin{aligned} first(B') &= \{a, b\}, follow(B') = \{f, d\}; first(B') \cap follow(B') = \emptyset; \\ first(A') &= \{a, b\}, follow(A') = \{f, d\}; first(A') \cap follow(A') = \emptyset; \\ first(C) &= \{a, b, c\}, follow(C) = \{f, d\}; first(C) \cap follow(C) = \emptyset. \end{aligned}$$

Т. е. получили эквивалентную грамматику G_{object} , к которой применим метод рекурсивного спуска. G_{object} удобна для построения системы рекурсивных процедур, так как ее правила имеют канонический вид.

Задача разбора для неоднозначных грамматик

Для неоднозначных грамматик задача синтаксического анализа (задача разбора) может быть поставлена двумя основными способами.

(1) Даны КС-грамматика G и цепочка x . Требуется проверить: $x \in L(G)$? Если да, то построить все деревья вывода для x (или все левые выводы для x , или все правые выводы для x)²⁴⁾.

Для решения этой задачи можно обобщить метод рекурсивного спуска, чтобы он работал с возвратами, пробуя различные подходящие альтернативы.

(2) Даны КС-грамматика G и цепочка x . Требуется проверить: $x \in L(G)$? Если да, то построить одно дерево вывода для x (возможно, наиболее подходящее в некотором смысле).

При такой постановке для некоторых неоднозначных грамматик удастся модифицировать обычный РС-метод без возвратов так, что получаемый анализатор корректен, и строит наиболее подходящее в некотором смысле дерево.

Неприменимость метода рекурсивного спуска в «чистом» виде для неоднозначных грамматик обусловлена невозможностью однозначно спрогнозировать выбор альтернативы при разборе цепочки (прогноз может состоять из нескольких подходящих альтернатив). Модификация метода состоит в следующем: одна из альтернатив объявляется «наиболее подходящей», и процедура анализа всегда выбирает эту альтернативу, игнорируя другие.

Рассмотрим пример, иллюстрирующий ситуацию с условными (полным и сокращенным) операторами в языке Паскаль.

$$G_{if-then} = \langle \{if, then, else, a, b\}, \{S\}, P, S, S' \rangle,$$

где $P = \{ S \rightarrow \mathbf{if} \ b \ \mathbf{then} \ S \ S' \mid a ; S' \rightarrow \mathbf{else} \ S \mid \varepsilon \}$. В этой грамматике прогноз для S' по \mathbf{else} неоднозначен, так как $first(\mathbf{else} \ S) \cap follow(S') = \{\mathbf{else}\} \neq \emptyset$. Для цепочки $\mathbf{if} \ b \ \mathbf{then} \ \mathbf{if} \ b \ \mathbf{then} \ a \ \mathbf{else} \ a$ можно построить два различных дерева вывода, показанных на рис. 101 (а, б).

Если при построении анализатора отдать предпочтение непустой альтернативе для S' , то такой анализатор построит дерево, изображенное на рис. 101 (а), в котором \mathbf{else} соотносится с ближайшим (на его уровне вложенности) \mathbf{if} , что соответствует правилам, принятым в

²⁴⁾ Цепочка в неоднозначной грамматике может иметь и бесконечно много деревьев вывода. В таком случае можно ограничиться построением всех деревьев, высота которых не превосходит некоторой константы.

языке Паскаль при разрешении подобных неоднозначностей в комбинациях условных операторов.

Итак, мы модифицировали РС-метод для данного примера неоднозначной грамматики, отдав предпочтение одной из альтернатив (и получив тем самым подходящее для семантики языка Паскаль дерево разбора).

Нетрудно убедиться, что получаемый для грамматики $G_{if-then}$ анализатор корректен: он не заикливается, распознает правильные цепочки и отвергает неправильные.

Отметим, что подобная модификация РС-метода не всегда приводит к построению корректного анализатора. Корректность необходимо дополнительно проверять.

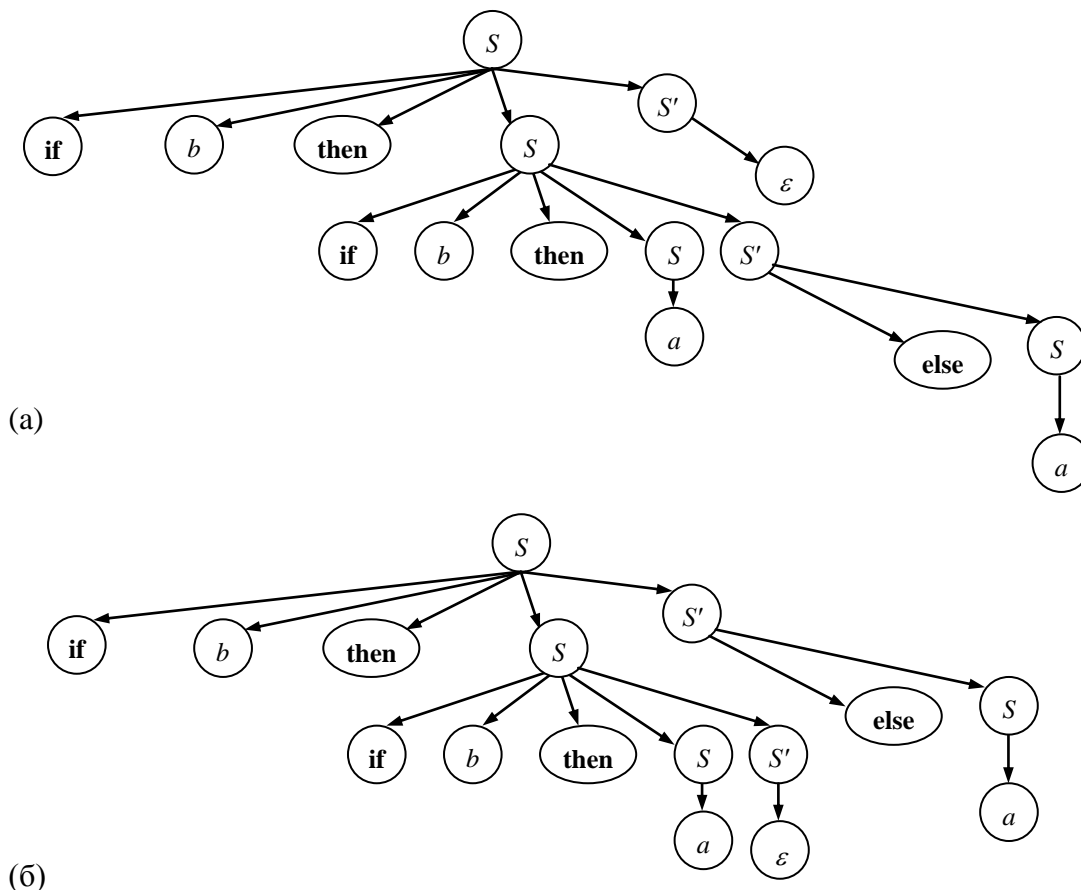


Рис. 10. Деревья вывода для цепочки *if b then if b then a else a*.

О других методах распознавания КС-языков

Метод рекурсивного спуска применим к достаточно узкому подклассу КС-грамматик. Известны более широкие подклассы КС-грамматик, для которых существуют эффективные анализаторы, обладающие тем же свойством, что и анализатор, построенный методом рекурсивного спуска, — входная цепочка считывается один раз слева направо и процесс разбора полностью детерминирован, в результате на обработку цепочки длины n расходуется время cn . К таким грамматикам относятся $LL(k)$ -грамматики, по которым, как правило, реализуется анализ сверху-вниз — нисходящий; $LR(k)$ -грамматики, грамматики предшествования, по которым, как правило, реализуется анализ снизу-вверх — восходящий; и некоторые другие (см., например, [2], [3]).

Анализатор для $LL(k)$ -грамматик просматривает входную цепочку слева направо и осуществляет детерминированный левый вывод, принимая во внимание k входных символов, расположенных справа от текущей позиции. Выбор альтернативы осуществляется на основе заранее составленной таблицы прогнозов.

Анализатор для $LR(k)$ -грамматик просматривает входную цепочку слева направо и осуществляет детерминированный правый вывод, принимая во внимание k входных символов, расположенных справа от текущей позиции. Вывод строится методом сверток, как при разборе по левосторонней автоматной грамматике. Предварительно по $LR(k)$ -грамматике строится таблица, которая на каждом шаге вывода позволяет анализатору однозначно выбрать нужную свертку.

Анализатор для грамматик предшествования просматривает входную цепочку слева направо и осуществляет детерминированный правый вывод, учитывая только некоторые отношения между парами смежных символов выводимой цепочки.

Часто одна и та же КС-грамматика может быть отнесена не к одному, а сразу к нескольким классам грамматик (например, любая LL -грамматика является LR -грамматикой, обратное неверно), допускающих построение линейных по временной сложности распознавателей. Но, на вопрос, какой лучше распознаватель выбрать, нисходящий или восходящий, нет однозначного ответа.

Нисходящий синтаксический анализ предпочтителен с точки зрения процесса трансляции, поскольку на его основе легче организовать процесс порождения цепочек результирующего языка. Восходящий синтаксический анализ привлекательнее тем, что часто для данного языка программирования легче построить LR -грамматику, поскольку ограничения на правила слабее, чем для LL -грамматик.

Конкретный выбор зависит от конкретного компилятора, от сложности грамматики входного языка программирования и от того, как будут использованы результаты работы распознавателя.

Синтаксический анализатор для М-языка

Будем считать, что синтаксический и лексический анализаторы взаимодействуют следующим образом: анализ исходной программы идет под управлением синтаксического анализатора; если для продолжения анализа ему нужна очередная лексема, то он запрашивает ее у лексического анализатора; тот выдает одну лексему и «замирает» до тех пор, пока синтаксический анализатор не запросит следующую лексему.

Соглашения:

- наш лексический анализатор — это функция-член `get_lex()` класса `Scanner`, которая в качестве результата выдает лексемы типа (*class*) `Lex`;
- в переменной `curr_lex` типа `Lex` будем хранить текущую лексему, выданную лексическим анализатором, а в переменной `c_val` — ее значение.

Анализатор методом рекурсивного спуска для М-языка реализуем в виде на Си++ в виде класса `Parser`. Кроме синтаксического анализа, процедуры данного класса осуществляют действия по семантическому контролю и переводу программы во внутреннее представление. Объяснение этих дополнительных действий и вспомогательных объектов будет дано в следующих разделах.

```

class Parser
{
    Lex          curr_lex; // текущая лексема
    type_of_lex  c_type;
    int          c_val;
    Scanner      scan;
    Stack < int, 100 > st_int;
    Stack < type_of_lex, 100 > st_lex;

    void P(); // процедуры РС-метода
    void D1();
    void D();
    void B();
    void S();
    void E();
    void E1();
    void T();
    void F();

    void dec ( type_of_lex type); // семантические действия
    void check_id ();
    void check_op ();
    void check_not ();
    void eq_type ();
    void eq_bool ();
    void check_id_in_read ();

    void gl () // получить очередную лексему
    {
        curr_lex = scan.get_lex();
        c_type = curr_lex.get_type();
        c_val = curr_lex.get_value();
    }

public:
    Poliz      prog; // внутреннее представление программы

    Parser ( const char *program) : scan (program), prog (1000) {}

    void      analyze(); // анализатор с действиями
};

void Parser::analyze ()
{
    gl ();
    P ();
    prog.print();
    cout << endl << "Yes!!!" << endl;
}

void Parser::P ()
{
    if ( c_type == LEX_PROGRAM )
        gl ();
}

```

```
    else
        throw curr_lex;
D1 ();
if ( c_type == LEX_SEMICOLON )
    gl ();
else
    throw curr_lex;
B ();
if ( c_type != LEX_FIN )
    throw curr_lex;
}

void Parser::D1 ()
{
    if ( c_type == LEX_VAR )
    {
        gl ();
        D ();
        while ( c_type == LEX_COMMA )
        { gl ();
          D ();
        }
    }
    else
        throw curr_lex;
}

void Parser::D ()
{
    st_int.reset();
    if ( c_type != LEX_ID )
        throw curr_lex;
    else
    {
        st_int.push ( c_val );
        gl ();
        while ( c_type == LEX_COMMA )
        {
            gl ();
            if ( c_type != LEX_ID )
                throw curr_lex;
            else {
                st_int.push ( c_val ); gl ();
            }
        }

        if ( c_type != LEX_COLON )
            throw curr_lex;
        else
        {
            gl ();
            if ( c_type == LEX_INT )
            {
                dec ( LEX_INT );
                gl ();
            }
            else

```

```
        if ( c_type == LEX_BOOL )
        {
            dec ( LEX_BOOL );
            gl();
        }
        else
            throw curr_lex;
    }
}

void Parser::B ()
{
    if ( c_type == LEX_BEGIN )
    {
        gl();
        S();
        while ( c_type == LEX_SEMICOLON )
        {
            gl();
            S();
        }
        if ( c_type == LEX_END )
            gl();
        else
            throw curr_lex;
    }
    else
        throw curr_lex;
}

void Parser::S ()
{
    int p10, p11, p12, p13;

    if ( c_type == LEX_IF )
    {
        gl();
        E();
        eq_bool();
        p12 = prog.get_free ();
        prog.blank();
        prog.put_lex (Lex(POLIZ_FGO));
        if ( c_type == LEX_THEN )
        {
            gl();

            S();
            p13 = prog.get_free();
            prog.blank();
            prog.put_lex (Lex(POLIZ_GO));
            prog.put_lex (Lex(POLIZ_LABEL, prog.get_free()),p12);
            if (c_type == LEX_ELSE)
            {
                gl();
                S();
            }
        }
    }
}
```

```
        prog.put_lex(Lex(POLIZ_LABEL,prog.get_free()),p13);
    }
    else
        throw curr_lex;
}
else
    throw curr_lex;
} //end if
else
if ( c_type == LEX_WHILE )
{
    p10 = prog.get_free();
    gl();
    E();
    eq_bool();
    p11 = prog.get_free();
    prog.blank();
    prog.put_lex (Lex(POLIZ_FGO));
    if ( c_type == LEX_DO )
    {
        gl();
        S();
        prog.put_lex (Lex (POLIZ_LABEL, p10));
        prog.put_lex (Lex ( POLIZ_GO));
        prog.put_lex (Lex(POLIZ_LABEL, prog.get_free()),p11);
    }
    else
        throw curr_lex;
} //end while
else
if ( c_type == LEX_READ )
{
    gl();
    if ( c_type == LEX_LPAREN )
    {
        gl();
        if ( c_type == LEX_ID )
        {
            check_id_in_read();
            prog.put_lex (Lex ( POLIZ_ADDRESS, c_val) );
            gl();
        }
        else
            throw curr_lex;
        if ( c_type == LEX_RPAREN )
        {
            gl();
            prog.put_lex (Lex (LEX_READ));
        }
        else
            throw curr_lex;
    }
    else
        throw curr_lex;
} //end read
else
```

```

        if ( c_type == LEX_WRITE )
        {
            gl();
            if ( c_type == LEX_LPAREN )
            {
                gl();
                E();
                if ( c_type == LEX_RPAREN )
                {
                    gl();
                    prog.put_lex (Lex(LEX_WRITE));
                }
                else
                    throw curr_lex;
            }
            else
                throw curr_lex;
        } //end write
    else
        if ( c_type == LEX_ID )
        {
            check_id ();
            prog.put_lex (Lex(POLIZ_ADDRESS,c_val));
            gl();
            if ( c_type == LEX_ASSIGN )
            {
                gl();
                E(); eq_type();
                prog.put_lex (Lex (LEX_ASSIGN) );
            }
            else
                throw curr_lex;
        } //assign-end
    else B();
}

void Parser::E ()
{
    E1();
    if ( c_type == LEX_EQ || c_type == LEX_LSS || c_type == LEX_GTR ||
        c_type == LEX_LEQ || c_type == LEX_GEQ || c_type == LEX_NEQ )
    {
        st_lex.push (c_type);
        gl();
        E1();
        check_op();
    }
}

void Parser::E1 ()
{
    T();
    while ( c_type==LEX_PLUS || c_type==LEX_MINUS || c_type==LEX_OR )
    {
        st_lex.push (c_type);
        gl();
    }
}

```



```
        T();
        check_op();
    }
}

void Parser::T ()
{
    F();
    while ( c_type==LEX_TIMES || c_type==LEX_SLASH || c_type==LEX_AND )
    {
        st_lex.push (c_type);
        gl();
        F();
        check_op();
    }
}

void Parser::F ()
{
    if ( c_type == LEX_ID )
    {
        check_id();
        prog.put_lex (Lex (LEX_ID, c_val));
        gl();
    }
    else
    if ( c_type == LEX_NUM )
    {
        st_lex.push ( LEX_INT );
        prog.put_lex ( curr_lex );
        gl();
    }
    else
    if ( c_type == LEX_TRUE )
    {
        st_lex.push ( LEX_BOOL );
        prog.put_lex (Lex (LEX_TRUE, 1) );
        gl();
    }
    else
    if ( c_type == LEX_FALSE )
    {
        st_lex.push ( LEX_BOOL );
        prog.put_lex (Lex (LEX_FALSE, 0) );
        gl();
    }
    else
    if ( c_type == LEX_NOT )
    {
        gl();

        F();
        check_not();
    }
    else
    if ( c_type == LEX_LPAREN )
    {
```

```

        gl();
        E();
        if ( c_type == LEX_RPAREN)
            gl();
        else
            throw curr_lex;
    }
    else
        throw curr_lex;
}

```

Семантический анализатор для М-языка

Контекстные условия, выполнение которых нам надо контролировать в программах на М-языке, таковы:

1. Любое имя, используемое в программе, должно быть описано и только один раз.
2. В операторе присваивания типы переменной и выражения должны совпадать.
3. В условном операторе и в операторе цикла в качестве условия возможно только логическое выражение.
4. Операнды операции отношения должны быть целочисленными.
5. Тип выражения и совместимость типов операндов в выражении определяются по обычным правилам (как в Паскале).

Проверку контекстных условий совместим с синтаксическим анализом. Для этого в синтаксические правила вставим вызовы процедур, осуществляющих необходимый контроль, а затем перенесем их в процедуры рекурсивного спуска.

Для реализации семантического анализа будем использовать следующий шаблонный класс *Stack*:

```

template <class T, int max_size > class Stack
{
    T s[max_size];
    int top;
public:
    Stack(){top = 0;}
    void reset() { top = 0; }
    void push(T i);
    T pop();
    bool is_empty(){ return top == 0; }
    bool is_full() { return top == max_size; }
};

template <class T, int max_size >

void Stack <T, max_size >::push(T i)
{
    if ( !is_full() )
    {
        s[top] = i;
        ++top;
    }
}

```

```

        else
            throw "stack_is_full";
    }

    template <class T, int max_size >
    T Stack <T, max_size >::pop()
    {
        if ( !is_empty() )
        {
            --top;
            return s[top];
        }
        else
            throw "stack_is_empty";
    }

```

Обработка описаний

Для контроля согласованности типов в выражениях и типов выражений в операторах, необходимо знать типы переменных, входящих в эти выражения. Кроме того, нужно проверять, нет ли повторных описаний идентификаторов. Эта информация становится известной в тот момент, когда синтаксический анализатор обрабатывает описания. Следовательно, в синтаксические правила для описаний нужно вставить действия, с помощью которых будем запоминать типы переменных и контролировать единственность их описания.

Лексический анализатор запомнил в таблице идентификаторов *TID* все идентификаторы-лексемы, которые были им обнаружены в тексте исходной программы. Информация о типе переменных и о наличии их описания заносится в ту же таблицу.

i-ая строка таблицы *TID* соответствует идентификатору-лексеме вида $\langle LEX_ID, i \rangle$.

Лексический анализатор заполнил поле *name*; значения полей *declare* и *type* будем заполнять на этапе семантического анализа.

Раздел описаний имеет вид

$$D \rightarrow I \{, I\}: [int \mid bool],$$

т. е. имени типа (*int* или *bool*) предшествует список идентификаторов. Эти идентификаторы (вернее, номера соответствующих им строк таблицы *TID*) надо запоминать (мы будем их запоминать в стеке целых чисел *Stack<int, 100> st_int*), а когда будет проанализировано имя типа, надо заполнить поля *declare* и *type* в этих строках.

Функция *void Parser::dec (type_of_lex type)* считывает из стека номера строк таблицы *TID*, заносит в них информацию о типе соответствующих переменных, о наличии их описаний и контролирует повторное описание переменных.

```

void Parser::dec ( type_of_lex type )
{
    int i;
    while ( !st_int.is_empty() )
    {
        i = st_int.pop();

        if ( TID[i].get_declare() )
            throw "twice";
    }
}

```

```

        else
        {
            TID[i].put_declare();
            TID[i].put_type(type);
        }
    }
}

```

С учетом имеющихся функций правило вывода с действиями для обработки описаний будет таким:

$$D \rightarrow \langle st_int.reset() \rangle I \langle st_int.push(c_val) \rangle \\ \{, I \langle st_int.push(c_val) \rangle \}: \\ [\mathit{int} \langle dec(LEX_INT) \rangle | \mathit{bool} \langle dec(LEX_BOOL) \rangle]$$

Контроль контекстных условий в выражении

Типы операндов и обозначения операций мы будем хранить в стеке *Stack* (*type_of_lex*, 100) *st_lex*.

Если в выражении встречается лексема *целое_число* или логические константы *true* или *false*, то соответствующий тип сразу заносится в стек.

Если операнд — лексема-переменная, то необходимо проверить, описана ли она; если описана, то ее тип надо занести в стек. Эти действия можно выполнить с помощью функции *check_id()*:

```

void parser::check_id()
{
    if ( TID[c_val].get_declare() )
        st_lex.push(TID[c_val].get_type());
    else
        throw "not declared";
}

```

Тогда для контроля контекстных условий каждой тройки — (операнд – операция – операнд) (для проверки соответствия типов операндов данной двуместной операции) будем использовать функцию *check_op()*:

```

void parser::check_op ()
{
    type_of_lex t1, t2, op, t = LEX_INT, r = LEX_BOOL;
    t2 = st_lex.pop();
    op = st_lex.pop();
    t1 = st_lex.pop();

    if ( op==LEX_PLUS || op==LEX_MINUS || op==LEX_TIMES || op==LEX_SLASH )
        r = LEX_INT;
    if ( op == LEX_OR || op == LEX_AND )
        t = LEX_BOOL;
    if ( t1 == t2 && t1 == t )
        st_lex.push(r);
    else
        throw "wrong types are in operation";
}

```

Для контроля за типом операнда одноместной операции **not** будем использовать функцию `check_not()`:

```
void Parser::check_not ()
{
    if (st_lex.pop() != LEX_BOOL)
        throw "wrong type is in not";
    else
    {
        st_lex.push (LEX_BOOL);
    }
}
```

Теперь главный вопрос: когда вызывать эти функции?

В грамматике модельного языка задано старшинство операций: наивысший приоритет имеет операция отрицания, затем в порядке убывания приоритета — группа операций умножения (*, /, and), группа операций сложения (+, -, or), операции отношения.

$$E \rightarrow E_1 \mid E_1 [= | < | >] E_1$$

$$E_1 \rightarrow T \{ [+ | - | \text{or}] T \}$$

$$T \rightarrow F \{ [* | / | \text{and}] F \}$$

$$F \rightarrow I \mid N \mid [\text{true} | \text{false}] \mid \text{not } F \mid (E)$$

Именно это свойство грамматики позволит провести синтаксически-управляемый контроль контекстных условий.

Упражнение. Сравните грамматики, описывающие выражения, состоящие из символов +, *, (,), i:

$$G_1: \quad E \rightarrow E + E \mid E * E \mid (E) \mid I$$

$$G_2: \quad E \rightarrow E + T \mid E * T \mid T \\ T \rightarrow i \mid (E)$$

$$G_3: \quad E \rightarrow T + E \mid T * E \mid T \\ T \rightarrow i \mid (E)$$

$$G_4: \quad E \rightarrow T \mid E + T \\ T \rightarrow F \mid T * F \\ F \rightarrow i \mid (E)$$

$$G_5: \quad E \rightarrow T \mid T + E \\ T \rightarrow F \mid F * T \\ F \rightarrow i \mid (E)$$

Оцените, насколько они удобны для трансляции выражений методом рекурсивного спуска (G_1 — неоднозначная, леворекурсивная; G_1, G_2, G_3 — не учитывается приоритет операций; G_4, G_5 — учитывается приоритет операций, G_2, G_4 — леворекурсивные, операции группируются слева направо, как принято в математике, G_3, G_5 — операции группируются справа налево).

Правила вывода выражений модельного языка с действиями для контроля контекстных условий:

$$\begin{aligned}
 E &\rightarrow E_1 \mid E_1 [= \mid < \mid >] \langle st_lex.push(c_type) \rangle E_1 \langle check_op() \rangle \\
 E_1 &\rightarrow T \{ [+ \mid - \mid or] \langle st_lex.push(c_type) \rangle T \langle check_op() \rangle \} \\
 T &\rightarrow F \{ [* \mid / \mid and] \langle st_lex.push(c_type) \rangle F \langle check_op() \rangle \} \\
 F &\rightarrow I \langle check_id() \rangle \mid N \langle st_lex(LEX_INT) \rangle \mid [true \mid false] \\
 &\quad \langle st_lex.push(LEX_BOOL) \rangle \mid not F \langle check_not() \rangle \mid (E)
 \end{aligned}$$

Контроль контекстных условий в операторах

$$S \rightarrow I := E \mid \text{if } E \text{ then } S \text{ else } S \mid \text{while } E \text{ do } S \mid B \mid \text{read } (I) \mid \text{write } (E)$$

1. Оператор присваивания

$$I := E$$

Контекстное условие: в операторе присваивания типы переменной I и выражения E должны совпадать.

В результате контроля контекстных условий выражения E в стеке останется тип этого выражения (как тип результата последней операции); если при анализе идентификатора I проверить, описан ли он, и занести его тип в тот же стек (для этого можно использовать функцию $check_id()$), то достаточно будет в нужный момент считать из стека два элемента и сравнить их:

```

void Parser::eq_type ()
{
    if ( st_lex.pop() != st_lex.pop() ) throw "wrong types are in :=";
}
    
```

Следовательно, правило для оператора присваивания:

$$I \langle check_id() \rangle := E \langle eq_type() \rangle$$

2. Условный оператор и оператор цикла

$$\text{if } E \text{ then } S \text{ else } S \mid \text{while } E \text{ do } S$$

Контекстные условия: в условном операторе и в операторе цикла в качестве условия возможны только логические выражения.

В результате контроля контекстных условий выражения E в стеке останется тип этого выражения (как тип результата последней операции); следовательно, достаточно извлечь его из стека и проверить:

```

void Parser::eq_bool ()
{
    if ( st_lex.pop() != LEX_BOOL )
        throw "expression is not boolean";
}
    
```

Тогда правила для условного оператора и оператора цикла будут такими:

$$\text{if } E \langle \text{eq_bool}() \rangle \text{ then } S \text{ else } S \mid \text{while } E \langle \text{eq_bool}() \rangle \text{ do } S$$

3. Для проверки операнда **оператора ввода** $\text{read}(I)$ можно использовать следующую функцию:

```
void Parser::check_id_in_read ()
{
    if ( !TID [c_val].get_declare() )
        throw "not declared";
}
```

Правило для оператора ввода будет таким:

$$\text{read} (I \langle \text{check_id_in_read}() \rangle) .$$

В итоге получаем процедуры для синтаксического анализа методом рекурсивного спуска с синтаксически-управляемым контролем контекстных условий, которые легко написать по правилам грамматики с действиями.

Класс *Parser* с дополнительными полями и методами для семантического анализа выглядит так:

```
class Parser
{
    Lex          curr_lex;
    type_of_lex  c_type;
    int          c_val;
    Scanner      scan;
    Stack < int, 100 > st_int;
    Stack < type_of_lex, 100 > st_lex;

    void P(); void D1(); void D(); void B(); void S();
    void E(); void E1(); void T(); void F();

    void dec ( type_of_lex type);
    void check_id ();
    void check_op ();
    void check_not ();
    void eq_type ();
    void eq_bool ();
    void check_id_in_read ();

    void gl ()
    {
        curr_lex = scan.get_lex();
        c_type = curr_lex.get_type();
        c_val = curr_lex.get_value();
    }

public:
    Parser ( const char *program) : scan (program) {}
    void analyze();
}
```

```
};  
  
void Parser::analyze ()  
{  
    gl ();  
    P ();  
    cout << endl << "Yes!!!" << endl;  
}
```

В качестве примера реализации приведем функцию с семантическими действиями для нетерминала *D* (раздел описаний):

```
void Parser::D ()  
{  
    st_int.reset();  
    if (c_type != LEX_ID)  
        throw curr_lex;  
    else  
    {  
        st_int.push ( c_val );  
        gl();  
        while (c_type == LEX_COMMA)  
        {  
            gl();  
            if (c_type != LEX_ID)  
                throw curr_lex;  
            else  
            {  
                st_int.push ( c_val );  
                gl();  
            }  
        }  
        if (c_type != LEX_COLON)  
            throw curr_lex;  
        else  
        {  
            gl();  
            if (c_type == LEX_INT)  
            {  
                dec ( LEX_INT );  
                gl();  
            }  
            else  
            if (c_type == LEX_BOOL)  
            {  
                dec ( LEX_BOOL );  
                gl();  
            }  
            else throw curr_lex;  
        }  
    }  
}
```


Генерация внутреннего представления программ

Результатом работы синтаксического анализатора должно быть некоторое внутреннее представление исходной цепочки лексем, которое отражает ее синтаксическую структуру. Программа в таком виде в дальнейшем может либо транслироваться в объектный код, либо интерпретироваться.

Язык внутреннего представления программы

Основные свойства языка внутреннего представления программ:

- а) он позволяет фиксировать синтаксическую структуру исходной программы;
- б) текст на нем можно автоматически генерировать во время синтаксического анализа;
- в) его конструкции должны относительно просто транслироваться в объектный код либо достаточно эффективно интерпретироваться.

Некоторые общепринятые способы внутреннего представления программ:

- а) постфиксная запись;
- б) префиксная запись;
- в) многоадресный код с явно именуемыми результатами;
- г) многоадресный код с неявно именуемыми результатами;
- д) связные списочные структуры, представляющие синтаксическое дерево.

В основе каждого из этих способов лежит некоторый метод представления синтаксического дерева.

Замечание

Чаще всего синтаксическим деревом называют дерево вывода исходной цепочки, в котором удалены вершины, соответствующие цепным правилам вида $A \rightarrow B$, где $A, B \in N$.

Выберем в качестве языка для представления промежуточной программы *постфиксную запись* (ее также называют **ПОЛИЗ** — польская инверсная запись).

ПОЛИЗ идеален для внутреннего представления интерпретируемых языков программирования, которые, как правило, удобно переводятся в ПОЛИЗ и легко интерпретируются.

В ПОЛИЗе операнды выписаны слева направо в порядке их следования в исходном тексте. Знаки операций стоят таким образом, что знаку операции непосредственно предшествуют ее операнды.

Простым будем называть выражение, состоящее из одной константы или имени переменной. Такие выражения в ПОЛИЗе остаются без изменений. При переводе в ПОЛИЗ сложных выражений важно правильно определять границы подвыражений, являющихся левыми и правыми операндами бинарных операций. Проблем не возникает, если сложные подвыражения явно ограничены скобками. Например, в выражении $(a + b) * c$ левым операндом операции $*$ является подвыражение $a + b$, а правым — простое выражение c . Когда скобки явно не расставлены, как в случаях $a + b * c$ и $a - b + c$, важно учитывать *приоритет* операций, а также *ассоциативность* операций одинакового приоритета. Умножение имеет больший приоритет, чем сложение, поэтому в выражении $a + b * c$ операнд b относится к опера-

ции умножения и эквивалентное выражение со скобками будет таким: $a + (b * c)$. В выражении $a - b + c$ операнд b относится к левой операции, т. е. к «минусу», а не к «плюсу» (в силу левой ассоциативности операций $+$ и $-$, имеющих одинаковый приоритет), и это выражение эквивалентно выражению $(a - b) + c$. Лево-ассоциативные операции группируются с помощью скобок слева направо: $a - b + c - d$ эквивалентно $((a - b) + c) - d$.

Теперь можем формально определить **постфиксную запись выражений** таким образом:

- 1) если E является простым выражением, то ПОЛИЗ выражения E — это само выражение E ;
- 2) ПОЛИЗом выражения $E_1 \theta E_2$, где θ — знак бинарной операции, E_1 и E_2 операнды для θ , является запись $E_1' E_2' \theta$, где E_1' и E_2' — ПОЛИЗ выражений E_1 и E_2 соответственно;
- 3) ПОЛИЗом выражения θE , где θ — знак унарной операции, а E — операнд θ , является запись $E' \theta$, где E' — ПОЛИЗ выражения E ;
- 4) ПОЛИЗом выражения (E) является ПОЛИЗ выражения E .

Пример. ПОЛИЗом выражения $a + b + c$ является $a b + c +$, но не $a b c + +$. Последняя запись является ПОЛИЗом для выражения $a + (b + c)$.

Алгоритм Дейкстры перевода в ПОЛИЗ выражений

Будем считать, что ПОЛИЗ выражения будет формироваться в массиве, содержащем лексемы — элементы ПОЛИЗа, и при переводе в ПОЛИЗ будет использоваться вспомогательный стек, также содержащий элементы ПОЛИЗа (операции, имена функций) и круглые скобки.

1. Выражение просматривается один раз слева направо.
2. Пока есть непрочитанные лексемы входного выражения, выполняем действия:
 - а) Читаем очередную лексему.
 - б) Если лексема является числом или переменной, добавляем ее в ПОЛИЗ-массив.
 - в) Если лексема является символом функции, помещаем ее в стек.
 - г) Если лексема является разделителем аргументов функции (например, запятая), то до тех пор, пока верхним элементом стека не станет открывающаяся скобка, выталкиваем элементы из стека в ПОЛИЗ-массив. Если открывающаяся скобка не встретилась, это означает, что в выражении либо неверно поставлен разделитель, либо не согласованы скобки.
 - д) Если лексема является операцией θ , тогда:
 - пока приоритет θ меньше либо равен приоритету операции, находящейся на вершине стека (для лево-ассоциативных операций), или приоритет θ строго меньше приоритета операции, находящейся на вершине стека (для право-ассоциативных операций), выталкиваем верхние элементы стека в ПОЛИЗ-массив;
 - помещаем операцию θ в стек.
 - е) Если лексема является открывающей скобкой, помещаем ее в стек.

ж) Если лексема является закрывающей скобкой, выталкиваем элементы из стека в ПОЛИЗ-массив до тех пор, пока на вершине стека не окажется открывающая скобка. При этом открывающая скобка удаляется из стека, но в ПОЛИЗ-массив не добавляется. Если после этого шага на вершине стека оказывается символ функции, выталкиваем его в ПОЛИЗ-массив. Если в процессе выталкивания открывающей скобки не нашлось и стек пуст, это означает, что в выражении не согласованы скобки.

3. Когда просмотр входного выражения завершен, выталкиваем все оставшиеся в стеке символы в ПОЛИЗ-массив. (В стеке должны были оставаться только символы операций; если это не так, значит в выражении не согласованы скобки.)

Например, обычной (инфиксной) записи выражения

$$a * (b + c) - (d - e) / f$$

соответствует такая постфиксная запись:

$$a b c + * d e - f / -$$

Замечание

Обратите внимание на то, что в ПОЛИЗе порядок операндов остался таким же, как и в инфиксной записи, учтено старшинство операций, а скобки исчезли.

Запись выражения в такой форме очень удобна для последующей интерпретации (т. е. вычисления значения этого выражения) с помощью стека.

Алгоритм вычисления выражений, записанных в ПОЛИЗе

- 1) Выражение просматривается один раз слева направо и для каждого элемента выполняются шаги (2) или (3);
- 2) Если очередной элемент ПОЛИЗа — операнд, то его значение заносится в стек;
- 3) Если очередной элемент ПОЛИЗа — операция, то на верхушке стека сейчас находятся ее операнды (это следует из определения ПОЛИЗа и предшествующих действий алгоритма); они извлекаются из стека, над ними выполняется операция, результат снова заносится в стек;
- 4) Когда выражение, записанное в ПОЛИЗе, прочитано, в стеке останется один элемент — это значение всего выражения, т. е. результат вычисления.

Замечание

1. Для интерпретации, кроме ПОЛИЗа выражения, необходима дополнительная информация об операндах, хранящаяся в таблицах.

2. Может оказаться так, что знак бинарной операции по написанию совпадает со знаком унарной; например, знак «-» в большинстве языков программирования означает и бинарную операцию вычитания, и унарную операцию изменения знака. В этом случае во время интерпретации операции «-» возникнет неоднозначность: сколько операндов надо извлекать из стека и какую операцию выполнять. Устранить неоднозначность можно, по крайней мере, двумя способами:

- заменить унарную операцию бинарной, т. е. считать, что $-a$ означает $0 - a$;
- либо ввести специальный знак для обозначения унарной операции; например, $-a$ заменить на $\&a$. Важно отметить, что это изменение касается только внутреннего представления программы и не требует изменения входного языка.

Оператор присваивания

$$I := E$$

в ПОЛИЗе будет записан как

$$\underline{I} E :=$$

где «:=» — это двухместная операция, а \underline{I} и E — ее операнды; подчеркнутое \underline{I} означает, что операндом операции «:=» является адрес переменной I , а не ее значение.

Оператор перехода в терминах ПОЛИЗа означает, что процесс интерпретации надо продолжить с того элемента ПОЛИЗа, который указан как операнд операции перехода.

Чтобы можно было ссылаться на элементы ПОЛИЗа, будем считать, что все они пронумерованы, начиная с 1 (допустим, занесены в последовательные элементы одномерного массива).

Пусть ПОЛИЗ оператора, помеченного меткой L , начинается с номера p , тогда оператор перехода $goto L$ в ПОЛИЗе можно записать как

$$p!$$

где «!» — операция выбора элемента ПОЛИЗа, номер которого равен p .

Немного сложнее окажется запись в ПОЛИЗе **условных операторов и операторов цикла**.

Например, если рассматривать оператор $if B then S_1 else S_2$ как обычную трехместную операцию с операндами B, S_1, S_2 , то ПОЛИЗ такого оператора должен выглядеть примерно так: $B' S_1' S_2' if$, где B' — ПОЛИЗ условия, $S_1' S_2'$ — ПОЛИЗ операторов S_1, S_2 , if — обозначение условной операции. Но тогда при интерпретации ПОЛИЗа обе ветви S_1, S_2 заранее вычисляются, независимо от условия B , что не соответствует семантике условного оператора. Для корректной реализации в ПОЛИЗе управляющих конструкций $if, while$ и т. п., их сначала заменяют эквивалентными фрагментами при помощи операторов перехода. Введем вспомогательную операцию — условный переход «по лжи» с семантикой

$$if (!B) goto L$$

Это двухместная операция с операндами B и L . Обозначим ее $!F$, тогда в ПОЛИЗе она будет записана как

$$B' p !F,$$

где p — номер элемента, с которого начинается ПОЛИЗ оператора, помеченного меткой L , B' — ПОЛИЗ логического выражения B .

Семантика условного оператора

$$if E then S_1 else S_2$$

с использованием введенной операции может быть описана так:

$$if (! E) goto L_2; S_1; goto L_3; L_2: S_2; L_3: \dots$$

Тогда ПОЛИЗ условного оператора будет таким (порядок операндов — прежний!):

$$E' p_2 !F S_1' p_3 ! S_2' \dots,$$

где p_i — номер элемента, с которого начинается ПОЛИЗ оператора, помеченного меткой L_i , $i = 2, 3$, E' — ПОЛИЗ логического выражения E . Заметим, что расположение внутренних

конструкций — условия E и операторов S_1, S_2 — относительно друг друга не изменилось. Это обязательное требование к ПОЛИЗу управляющих конструкций.

Семантика оператора цикла **while** E **do** S может быть описана так:

$$L_0: \text{if } (! E) \text{ goto } L_1; S; \text{goto } L_0; L_1: \dots$$

Тогда ПОЛИЗ оператора цикла **while** будет таким (порядок операндов — прежний!):

$$E' p_1 !F S' p_0 ! \dots,$$

где p_i — номер элемента, с которого начинается ПОЛИЗ оператора, помеченного меткой L_i , $i = 0, 1$, E' — ПОЛИЗ логического выражения E .

Операторы ввода и вывода М-языка являются одноместными операциями.

Оператор ввода **read** (I) в ПОЛИЗе будет записан как I read;

Оператор вывода **write** (E) в ПОЛИЗе будет записан как E' **write**, где E' — ПОЛИЗ выражения E .

Постфиксная польская запись операторов обладает всеми свойствами, характерными для постфиксной польской записи выражений, поэтому алгоритм интерпретации выражений пригоден для интерпретации всей программы, записанной в ПОЛИЗе (нужно только расширить набор операций; кроме того, выполнение некоторых из них не будет давать результата, записываемого в стек).

Постфиксная польская запись может использоваться не только для интерпретации промежуточной программы, но и для генерации по ней объектной программы. Для этого в алгоритме интерпретации вместо выполнения операции нужно генерировать соответствующие команды объектной программы.

Синтаксически управляемый перевод

На практике синтаксический, семантический анализ и генерация внутреннего представления программы часто осуществляются одновременно.

Существует несколько способов построения промежуточной программы. Один из них, называемый синтаксически управляемым переводом, особенно прост и эффективен.

В основе синтаксически управляемого перевода лежит уже известная нам грамматика с действиями (см. раздел о контроле контекстных условий). Теперь, параллельно с анализом исходной цепочки лексем, будем выполнять действия по генерации внутреннего представления программы. Для этого дополним грамматику вызовами соответствующих процедур генерации.

Содержательный пример — генерация внутреннего представления программы для М-языка — приведен ниже, а здесь в качестве иллюстрации рассмотрим более простой пример.

Пусть есть грамматика, описывающая простейшее арифметическое выражение.

$$G_{expr}: \\ E \rightarrow T \{+T\} \\ T \rightarrow F \{*F\} \\ F \rightarrow a | b | (E)$$

Тогда грамматика с действиями по переводу этого выражения в ПОЛИЗ будет такой.

$$\begin{aligned}
 G_{expr_polish}: \\
 E &\rightarrow T \{ +T \langle cout \ll '+'; \rangle \} \\
 T &\rightarrow F \{ *F \langle cout \ll '*'; \rangle \} \\
 F &\rightarrow a \langle cout \ll 'a'; \rangle \mid b \langle cout \ll 'b'; \rangle \mid (E)
 \end{aligned}$$

В процессе анализа методом рекурсивного спуска входной цепочки $a + b * c$ по грамматике G_{expr_polish} в выходной поток будет выведена цепочка $a b c * +$, являющаяся польской записью исходной цепочки. Таким образом, данная грамматика с действиями каждой цепочки языка арифметических выражений ставит в соответствие подходящую цепочку языка польских записей арифметических выражений. Иными словами, такая грамматика задает *перевод* из одного формального языка в другой.

Определение: Пусть T_1 и T_2 — алфавиты. *Формальный перевод* τ — это подмножество множества всевозможных пар цепочек в алфавитах T_1 и T_2 : $\tau \subseteq (T_1^* \times T_2^*)$.

Назовем *входным* языком перевода τ язык $L_{ex}(\tau) = \{ \alpha \mid \exists \beta: (\alpha, \beta) \in \tau \}$.

Назовем *целевым* (или *выходным*) языком перевода τ язык $L_y(\tau) = \{ \beta \mid \exists \alpha: (\alpha, \beta) \in \tau \}$.

Перевод τ *неоднозначен*, если для некоторых $\alpha \in T_1^*$, $\beta, \gamma \in T_2^*$, $\beta \neq \gamma$, $(\alpha, \beta) \in \tau$ и $(\alpha, \gamma) \in \tau$.

Рассмотренная выше грамматика G_{expr_polish} задает однозначный перевод: каждому выражению ставится в соответствие единственная польская запись. Неоднозначные переводы могут быть интересны при изучении моделей естественных языков; для трансляции языков программирования используются однозначные переводы.

Заметим, что для двух заданных языков L_1 и L_2 существует бесконечно много формальных переводов²⁵⁾. Чтобы задать перевод из L_1 в L_2 , важно точно указать *закон соответствия* между цепочками L_1 и L_2 .

Пример. Пусть $L_1 = \{ 0^n 1^m \mid n \geq 0, m > 0 \}$ — входной язык, $L_2 = \{ a^m b^n \mid n \geq 0, m > 0 \}$ — выходной язык и перевод τ определяется так: для любых $n \geq 0, m > 0$ цепочке $0^n 1^m \in L_1$ соответствует цепочка $a^m b^n \in L_2$. Можно записать τ с помощью теоретико-множественной формулы: $\tau = \{ (0^n 1^m, a^m b^n) \mid n \geq 0, m > 0 \}$, $L_{ex}(\tau) = L_1, L_y(\tau) = L_2$.

Задача.

Реализовать перевод $\tau = \{ (0^n 1^m, a^m b^n) \mid n \geq 0, m > 0 \}$ грамматикой с действиями.

Решение

Язык L_1 можно описать грамматикой:

$$\begin{aligned}
 S &\rightarrow 0S \mid 1A \\
 A &\rightarrow 1A \mid \varepsilon
 \end{aligned}$$

Вставим действия по переводу цепочек вида $0^n 1^m$ в соответствующие цепочки вида $a^m b^n$:

$$\begin{aligned}
 S &\rightarrow 0S \langle cout \ll 'b'; \rangle \mid 1 \langle cout \ll 'a'; \rangle A \\
 A &\rightarrow 1 \langle cout \ll 'a'; \rangle A \mid \varepsilon
 \end{aligned}$$

²⁵⁾ При условии, что хотя бы один из языков L_1, L_2 бесконечен. Действительно, пусть $L_1 = \{ a^n \mid n \geq 0 \}$, $L_2 = \{ b, c \}$. Существует бесконечно много переводов $\tau_i = \{ (a^i, b) \} \cup \{ (a^n, c) \mid n \neq i, i \geq 0 \}$, где $L_{ex}(\tau_i) = L_1, L_y(\tau_i) = L_2$.

Теперь при анализе методом рекурсивного спуска цепочек языка L_1 с помощью действий будут порождаться соответствующие цепочки языка L_2 .

Генератор внутреннего представления программы на М-языке

Каждый элемент в ПОЛИЗе — это лексема, т. е. пара вида $\langle \text{тип_лексемь, значение_лексемь} \rangle$. При генерации ПОЛИЗа будем использовать дополнительные типы лексем (для внутреннего представления операторов):

- POLIZ_GO — «!»;
- POLIZ_FGO — «!F»;
- POLIZ_LABEL — для ссылок на номера элементов ПОЛИЗа;
- POLIZ_ADDRESS — для обозначения операндов-адресов (например, в ПОЛИЗе оператора присваивания).

Будем считать, что генерируемая программа размещается в объекте *prog* класса *Poliz*, заданном описанием:

Poliz prog (1000);

```
class Poliz
{
    lex *p;
    int size;
    int free;
public:
    Poliz ( int max_size )
    {
        p = new Lex [size = max_size];
        free = 0;
    };
    ~Poliz() { delete []p; };
    void put_lex(Lex l) { p[free]=l; ++free; };
    void put_lex(Lex l, int place) { p[place]=l; };
    void blank() { ++free; };
    int get_free() { return free; };
    lex& operator[] ( int index )
    {
        if (index > size)
            throw "POLIZ:out of array";
        else
            if ( index > free )
                throw "POLIZ:indefinite element of array";
            else
                return p[index];
    };
    void print()
    {
        for ( int I = 0; i < free; ++i )
            cout << p[i];
    };
};
```

Объект *prog* для хранения внутреннего представления программы разместим в открытой части класса *Parser*:

```
class Parser
{
    Lex curr_lex;
    ...

public:
    Poliz      prog;
    Parser (const char *program ) : scan (program),prog (1000) {}
    void      analyze();
};
```

Генерация внутреннего представления программы будет проходить во время синтаксического анализа параллельно с контролем контекстных условий, поэтому для генерации можно использовать информацию, собранную синтаксическим и семантическим анализаторами; например, при генерации ПОЛИЗа выражений можно воспользоваться содержимым стека, с которым работает семантический анализатор. Кроме того, можно дополнить функции семантического анализа действиями по генерации.

Добавим в конец тела функции *check_op()* оператор *prog.put_lex (Lex (op))*; записывающий в ПОЛИЗ знак операции *op*, а в конец функции *check_not()* добавим оператор *prog.put_lex (Lex (LEX_NOT))*; записывающий лексему **not** (во внутреннем представлении) в ПОЛИЗ.

Тогда грамматика, содержащая действия по контролю контекстных условий и переводу выражений модельного языка в ПОЛИЗ, будет такой:

$$\begin{aligned}
 E &\rightarrow E_1 \mid E_1 [= | < | >] \langle st_lex.push (c_type) \rangle E_1 \langle check_op () \rangle \\
 E_1 &\rightarrow T \{ [+ | - | \mathbf{or}] \langle st_lex.push (c_type) \rangle T \langle check_op () \rangle \} \\
 T &\rightarrow F \{ [* | / | \mathbf{and}] \langle st_lex.push (c_type) \rangle F \langle check_op () \rangle \} \\
 F &\rightarrow I \langle check_id (); prog.put_lex (curr_lex); \rangle \mid \\
 &\quad N \langle st_lex.push (LEX_INT); prog.put_lex (curr_lex); \rangle \mid \\
 &\quad [\mathbf{true} \mid \mathbf{false}] \langle st_lex.push (LEX_BOOL); prog.put_lex (curr_lex); \rangle \mid \\
 &\quad \mathbf{not} F \langle check_not(); \rangle \mid (E)
 \end{aligned}$$

Действия, которыми нужно дополнить правило вывода оператора присваивания, также достаточно очевидны:

$$\begin{aligned}
 S &\rightarrow I \langle check_id (); prog.put_lex (Lex (POLIZ_ADDRESS, c_val)); \rangle := \\
 &\quad E \langle eqtype (); prog.put_lex (Lex (LEX_ASSIGN)); \rangle
 \end{aligned}$$

При генерации ПОЛИЗа выражений и оператора присваивания элементы объекта *prog* заполнялись последовательно. Семантика условного оператора *if E then S₁ else S₂* такова, что значения операндов для операций безусловного перехода и перехода «по лжи» в момент генерации операций еще неизвестны:

$$\mathbf{if} (!E) \mathbf{goto} l_2; S_1; \mathbf{goto} l_3; l_2: S_2; l_3: \dots$$

Поэтому придется оставлять соответствующие этим операндам элементы объекта *prog* незаполненными и запоминать их номера, а затем, когда станут известны их значения, заполнять пропущенное.

Тогда грамматика, содержащая действия по контролю контекстных условий и переводу условного оператора модельного языка в ПОЛИЗ, будет такой:

$$\begin{aligned}
 S \rightarrow & \text{if } E \langle \text{eqbool}(); pl_2 = \text{prog.get_free}(); \text{prog.blank}(); \\
 & \text{prog.put_lex} (\text{Lex} (\text{POLIZ_FGO})); \rangle \\
 & \text{then } S_1 \langle pl_3 = \text{prog.get_free}(); \text{prog.blank}(); \\
 & \text{prog.put_lex} (\text{Lex} (\text{POLIZ_GO})); \\
 & \text{prog.put_lex} (\text{Lex} (\text{POLIZ_LABEL}, \text{prog.get_free}()), pl_2); \rangle \\
 & \text{else } S_2 \langle \text{prog.put_lex} (\text{Lex} (\text{POLIZ_LABEL}, \text{prog.get_free}()), pl_3); \rangle
 \end{aligned}$$

Семантика оператора цикла *while E do S* описывается так:

$$L_0: \text{if } (!E) \text{goto } l_1; S; \text{goto } l_0; l_1: \dots,$$

а грамматика с действиями по контролю контекстных условий и переводу оператора цикла в ПОЛИЗ будет такой:

$$\begin{aligned}
 S \rightarrow & \text{while } \langle pl_0 = \text{prog.get_free} (); \rangle E \langle \text{eqbool} (); \\
 & pl_1 = \text{prog.get_free} (); \text{prog.blank} (); \\
 & \text{prog.put_lex} (\text{Lex} (\text{POLIZ_FGO})); \rangle \\
 & \text{do } S \langle \text{prog.put_lex} (\text{Lex} (\text{POLIZ_LABEL}, pl_0); \\
 & \text{prog.put_lex} (\text{Lex} (\text{POLIZ_GO})); \\
 & \text{prog.put_lex} (\text{Lex} (\text{POLIZ_LABEL}, \text{prog.get_free}()), pl_1); \rangle
 \end{aligned}$$

Замечание

Переменные pl_i ($i = 0, 1, 2, 3$) должны быть локализованы в процедуре *S*, иначе возникнет ошибка при обработке вложенных операторов.

Наконец, запишем соответствующие действия для операторов ввода и вывода:

$$\begin{aligned}
 S \rightarrow & \text{read} (I \langle \text{check_id_in_read}(); \\
 & \text{prog.put_lex} (\text{Lex} (\text{POLIZ_ADDRESS}, c_val)); \rangle) \\
 & \langle \text{prog.put_lex} (\text{Lex} (\text{LEX_READ})); \rangle
 \end{aligned}$$

$$S \rightarrow \text{write} (E) \langle \text{prog.put_lex} (\text{Lex} (\text{LEX_WRITE})); \rangle$$

Интерпретатор ПОЛИЗа для модельного языка

Польская инверсная запись была выбрана нами в качестве языка внутреннего представления программы, в частности, потому, что записанная таким образом программа может быть легко проинтерпретирована.

Идея алгоритма очень проста: просматриваем ПОЛИЗ слева направо; если встречаем операнд, то записываем его в стек; если встретили знак операции, то извлекаем из стека нужное количество операндов и выполняем операцию, результат (если он есть) заносим в стек и т. д.

Итак, записанная в ПОЛИЗе программа хранится в виде последовательности лексем в объекте *prog* класса *Poliz*. Лексемы могут быть следующие: лексемы-константы (числа, *true*,

false), лексемы-метки ПОЛИЗа, лексемы-операции (включая дополнительно введенные в ПОЛИЗе) и лексемы-переменные (их значения или адреса — номера строк в таблице *TID*).

В программе-интерпретаторе будем использовать некоторые переменные и функции, введенные нами ранее.

```
class Executer
{
    Lex pc_el;
public:
    void execute ( Poliz& prog );
};
void Executer::execute ( Poliz& prog )
{
    Stack < int, 100 > args;
    int i, j, index = 0, size = prog.get_free();

    while ( index < size )
    {
        pc_el = prog [ index ];

        switch ( pc_el.get_type () )
        {
            case LEX_TRUE:
            case LEX_FALSE:
            case LEX_NUM:
            case POLIZ_ADDRESS:
            case POLIZ_LABEL:
                args.push ( pc_el.get_value () );
                break;
            case LEX_ID:
                i = pc_el.get_value ();
                if ( TID[i].get_assign () )
                {
                    args.push ( TID[i].get_value () );
                    break;
                }
                else
                    throw "POLIZ: indefinite identifier";
            case LEX_NOT:
                args.push( !args.pop() );
                break;
            case LEX_OR:
                i = args.pop();
                args.push ( args.pop() || i );
                break;
            case LEX_AND:
                i = args.pop();
                args.push ( args.pop() && i );
                break;
            case POLIZ_GO:
                index = args.pop() - 1;
                break;
            case POLIZ_FGO:
                i = args.pop();

                if ( !args.pop() )
```

```

        index = i - 1;
    break;
case LEX_WRITE:
    cout << args.pop () << endl;
    break;
case LEX_READ:
    {
        int k;
        i = args.pop ();
        if ( TID[i].get_type () == LEX_INT )
        {
            cout << "Input int value for";
            cout << TID[i].get_name () << endl;
            cin >> k;
        }
        else
        {
            char j[20];
            rep:
            cout << "Input boolean value;
            cout << (true or false) for";

            cout << TID[i].get_name() << endl;
            cin >> j;
            if ( !strcmp(j, "true") )
                k = 1;
            else if ( !strcmp(j, "false") )
                k = 0;
            else
            {
                cout << "Error in input:true/false";
                cout << endl;
                goto rep;
            }
        }
        TID[i].put_value (k);
        TID[i].put_assign ();
        break;}
case LEX_PLUS:
    args.push ( args.pop() + args.pop() );
    break;
case LEX_TIMES:
    args.push ( args.pop() * args.pop() );
    break;
case LEX_MINUS:
    i = args.pop();
    args.push ( args.pop() - i );
    break;
case LEX_SLASH:
    i = args.pop();
    if ( !i )
    {
        args.push ( args.pop() / i );
        break;
    }
    else throw "POLIZ:divide by zero";
case LEX_EQ:

```

```

        args.push ( args.pop() == args.pop() );

        break;
    case LEX_LSS:
        i = args.pop();
        args.push ( args.pop() < i );
        break;
    case LEX_GTR:
        i = args.pop();
        args.push ( args.pop() > i );
        break;
    case LEX_LEQ:
        i = args.pop();
        args.push ( args.pop() <= i );
        break;
    case LEX_GEQ:
        i = args.pop();
        args.push ( args.pop() >= i );
        break;
    case LEX_NEQ:
        i = args.pop();

        args.push ( args.pop() != i );
        break;
    case LEX_ASSIGN:
        i = args.pop();
        j = args.pop();
        TID[j].put_value(i);
        TID[j].put_assign();
        break;
    default:
        throw "POLIZ: unexpected elem";
    }
    // end of switch
    ++index;
};
//end of while
cout << "Finish of executing!!!" << endl;
}

class Interpretator
{
    Parser pars;
    Executer E;
public:
    Interpretator ( char* program ): pars (program) {};
    void interpretation ();
};

void Interpretator::interpretation ()
{
    pars.analyze ();
    E.execute ( pars.prog );
}

int main ()
{
    try
    {
        Interpretator I ( "program.txt" );
    }
}

```

```

        I.interpretation ();
        return 0;
    }

    catch ( char c )
    {
        cout << "unexpected symbol " << c << endl;
        return 1;
    }
    catch ( Lex l )
    {
        cout << "unexpected lexeme";
        cout << l;
        return 1;
    }
    catch ( const char *source )
    {
        cout << source << endl;
        return 1;
    }
}

```

Задачи

I. Грамматики и языки. Классификация по Хомскому

1. Даны грамматика и цепочка. Построить вывод заданной цепочки.

$$\text{a) } S \rightarrow T \mid T+S \mid T-S$$

$$T \rightarrow F \mid F*T$$

$$F \rightarrow a \mid b$$

$$\text{Цепочка: } a-b*a+b$$

$$\text{b) } S \rightarrow aSBC \mid abC$$

$$CB \rightarrow BC$$

$$bB \rightarrow bb$$

$$bC \rightarrow bc$$

$$cC \rightarrow cc$$

$$\text{Цепочка: } aaabbbccc$$

2. Построить все сентенциальные формы для грамматики с правилами:

$$S \rightarrow A+B \mid B+A$$

$$A \rightarrow a$$

$$B \rightarrow b$$

3. К какому типу по Хомскому относится данная грамматика? Какой язык она порождает? Каков тип языка? Указать максимально возможный номер типа грамматики и языка.

$$\text{a) } S \rightarrow APA$$

$$P \rightarrow + \mid -$$

$$A \rightarrow a \mid b$$

$$\text{b) } S \rightarrow A \mid SA \mid SB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

- c) $S \rightarrow 1B$
 $B \rightarrow B0 \mid 1$
- e) $S \rightarrow a \mid Ba$
 $B \rightarrow Bb \mid b$
- g) $S \rightarrow 0A1 \mid 01$
 $0A \rightarrow 00A1$
 $A \rightarrow 01$
- i) $S \rightarrow A \mid B$
 $A \rightarrow aAb \mid 0$
 $B \rightarrow aBbb \mid 1$
- k) $S \rightarrow 0S \mid S0 \mid D$
 $D \rightarrow DD \mid 1A \mid \varepsilon$
 $A \rightarrow 0B \mid \varepsilon$
 $B \rightarrow 0A \mid 0$
- m) $S \rightarrow SS \mid A$
 $A \rightarrow a \mid bb$
- o) $S \rightarrow aBA \mid \varepsilon$
 $B \rightarrow bSA$
 $AA \rightarrow c$
- r) 1. $S \rightarrow KF$
 2. $K \rightarrow KB \mid CB$
 3. $C \rightarrow CA \mid DA$
 4. $DA \rightarrow aAD$
 5. $Aa \rightarrow aA$
 6. $DB \rightarrow bBD$
- s) 1. $S \rightarrow DC$
 2. $D \rightarrow aDA \mid bDB \mid aA \mid bB$
 3. $AC \rightarrow aC$
 4. $BC \rightarrow bC$
 5. $Aa \rightarrow aA$
- t) $S \rightarrow aAc$
 $aA \rightarrow aaBbC \mid ab$
 $Bb \rightarrow bb \mid abbbc \mid aDbbbcc$
 $C \rightarrow c$
 $D \rightarrow ab$
- d) $S \rightarrow aQb \mid \varepsilon$
 $Q \rightarrow cSc$
- f) $S \rightarrow Ab$
 $A \rightarrow Aa \mid ba$
- h) $S \rightarrow AB$
 $AB \rightarrow BA$
 $A \rightarrow a$
 $B \rightarrow b$
- j) $S \rightarrow 0A \mid 1S$
 $A \rightarrow 0A \mid 1B$
 $B \rightarrow 0B \mid 1B \mid \perp$
- l) $S \rightarrow 0A \mid 1S \mid \varepsilon$
 $A \rightarrow 1A \mid 0B$
 $B \rightarrow 0S \mid 1B$
- n) $S \rightarrow AB\perp$
 $A \rightarrow a \mid cA$
 $B \rightarrow bA$
- p) $S \rightarrow Ab \mid c$
 $A \rightarrow Ba$
 $B \rightarrow cS$
7. $Bb \rightarrow bB$
8. $Ab \rightarrow bA$
9. $DF \rightarrow E$
10. $BE \rightarrow Eb$
11. $AE \rightarrow Ea$
12. $bE \rightarrow b$
6. $Ba \rightarrow aB$
7. $Ab \rightarrow bA$
8. $Bb \rightarrow bB$
9. $C \rightarrow \varepsilon$
- u) $S \rightarrow 0A1$
 $0A \rightarrow 0B1 \mid 0$
 $B1 \rightarrow 0C11 \mid 01$
 $C \rightarrow 0D \mid 00D1 \mid 0$
 $D \rightarrow 01$

4. Построить грамматику, порождающую заданный язык L . Каков тип построенной грамматики? Каков тип языка?

a) $L = \{ a^n b^m \mid n, m \geq 1 \}$;

b) $L = \{ \alpha c \beta c \gamma c \mid \alpha, \beta, \gamma \in \{a, b\}^* \}$ (т.е. α, β, γ — любые цепочки из a и b);

c) $L = \{ a_1 a_2 \dots a_n a_n \dots a_2 a_1 \mid a_i \in \{0, 1\}, n \geq 1 \}$;

d) $L = \{ 0^n 1^{\lfloor n/2 \rfloor} \mid n \geq 1 \}$;

e) $L = \{ c a^n c b^m c^n \mid n \geq 0, m \geq 0 \}$;

f) $L = \{ a^n b^m \mid n \neq m; n, m \geq 0 \}$;

g) $L = \{ \omega \mid \omega \in \{0, 1\}^*, |\omega|_0 \neq |\omega|_1 \}$ (т.е. все цепочки из 0 и 1 с неравным числом 0 и 1);

h) $L = \{ \alpha \alpha \mid \alpha \in \{a, b\}^+ \}$;

i) $L = \{ \omega \mid \omega \in \{0, 1\}^+, |\omega|_0 = |\omega|_1, \forall x, y \in \{0, 1\}^* : \omega = xy \Rightarrow |x|_1 \leq |x|_0 \}$
(т.е. все непустые цепочки, содержащие равное количество 0 и 1, причем любая подцепочка, взятая с левого конца, содержит единиц не больше, чем нулей);

j) $L = \{ \omega_1 \omega_2 \dots \omega_n \mid n \geq 0, \omega_i \in \{ a^{2^m} b^m \mid m \geq 1 \}$ для $1 \leq i \leq n \}$;

k) $L = \{ a^{3^n+1} \mid n \geq 1 \}$;

l) $L = \{ a^{n^2} \mid n \geq 1 \}$;

m) $L = \{ a^{n^3+1} \mid n \geq 1 \}$.

5. К какому типу по Хомскому относится данная грамматика (указать максимально возможный номер)? Какой язык она порождает? Каков тип языка? Выписать подтверждающую ответ грамматику, в состав которой входит только один нетерминал — цель грамматики.

a) $S \rightarrow AB \mid ASB$ $A \rightarrow a$ $aB \rightarrow b$ $bB \rightarrow bb$	б) $S \rightarrow 1A0$ $1A \rightarrow 11A0 \mid 01$
--	---

6. Эквивалентны ли грамматики с правилами:

а) $S \rightarrow AB$ $A \rightarrow a \mid Aa$ $B \rightarrow b \mid Bb$	и	$S \rightarrow AS \mid SB \mid AB$ $A \rightarrow a$ $B \rightarrow b$
б) $S \rightarrow aSL \mid aL$ $L \rightarrow Kc$ $cK \rightarrow Kc$ $K \rightarrow b$	и	$S \rightarrow aSBc \mid abc$ $cB \rightarrow Bc$ $bB \rightarrow bb$

7. Построить КС-грамматику, эквивалентную грамматике с правилами:

- | | |
|-----------------------------|--------------------------------|
| a) $S \rightarrow aAb$ | b) $S \rightarrow AB \mid ABS$ |
| $aA \rightarrow aaAb$ | $AB \rightarrow BA$ |
| $A \rightarrow \varepsilon$ | $BA \rightarrow AB$ |
| | $A \rightarrow a$ |
| | $B \rightarrow b$ |

8. Построить регулярную грамматику, эквивалентную грамматике с правилами:

- | | |
|------------------------------|---------------------------|
| a) $S \rightarrow A \mid AS$ | b) $S \rightarrow A.A$ |
| $A \rightarrow a \mid bb$ | $A \rightarrow B \mid BA$ |
| | $B \rightarrow 0 \mid 1$ |

9. Привести пример грамматики, каждое правило которой относится к одному из трех видов $A \rightarrow Bt$, либо $A \rightarrow tB$, либо $A \rightarrow t$, для которой не существует эквивалентной регулярной грамматики.

10. Доказать, что грамматика с правилами:

- $$S \rightarrow aSBC \mid abC$$
- $$CB \rightarrow BC$$
- $$bB \rightarrow bb$$
- $$bC \rightarrow bc$$
- $$cC \rightarrow cc$$

порождает язык $L = \{a^n b^n c^n \mid n \geq 1\}$.

Для этого показать, что в данной грамматике:

- I) выводится любая цепочка вида $a^n b^n c^n$ ($n \geq 1$) и
 II) не выводятся никакие другие цепочки.

11. Дана грамматика с правилами:

- | | |
|---|---|
| a) $S \rightarrow S0 \mid S1 \mid D0 \mid D1$ | b) $S \rightarrow \mathbf{if} \ B \ \mathbf{then} \ S \mid B = E$ |
| $D \rightarrow H$. | $E \rightarrow B \mid B + E$ |
| $H \rightarrow 0 \mid 1 \mid H0 \mid H1$ | $B \rightarrow a \mid b$ |

Построить восходящим и нисходящим методами дерево вывода для цепочки:

- | | |
|------------|--|
| a) 10.1001 | b) $\mathbf{if} \ a \ \mathbf{then} \ b = a+b+b$ |
|------------|--|

12. Определить тип грамматики. Описать язык, порождаемый этой грамматикой. Построить для этого языка КС-грамматику.

- $$S \rightarrow P\perp$$
- $$P \rightarrow 1P1 \mid 0P0 \mid T$$
- $$T \rightarrow 021 \mid 120R$$
- $$R1 \rightarrow 0R$$
- $$R0 \rightarrow 1R$$
- $$R\perp \rightarrow 1\perp$$

13. Построить регулярную грамматику, порождающую цепочки в алфавите $\{a, b\}$, в которых символ a не встречается два раза подряд.

14. Построить КС-грамматику для языка L , построить дерево вывода и левосторонний вывод для цепочки $aabbbcccc$.

$$L = \{a^{2n}b^m c^{2k} \mid m = n + k, m > 1\}.$$

15. Построить грамматику, порождающую сбалансированные относительно круглых скобок цепочки в алфавите $\{a, (,), \perp\}$. Сбалансированную цепочку α определим рекурсивно: цепочка α сбалансирована, если

- а) α не содержит скобок,
- б) $\alpha = (\alpha_1)$ или $\alpha = \alpha_1\alpha_2$, где α_1 и α_2 сбалансированы.

16. Построить КС-грамматику, порождающую язык L , и вывод для цепочки $aacbhbca$ в этой грамматике.

$$L = \{a^n cb^m ca^n \mid n, m > 0\}.$$

17. Построить КС-грамматику, порождающую язык L , и вывод для цепочки 110000111 в этой грамматике.

$$L = \{1^n 0^m 1^p \mid n + p > m; n, p, m > 0\}.$$

18. Дан язык $L = \{1^{3n+2} 0^n \mid n \geq 0\}$. Определить его тип, построить грамматику, порождающую L . Построить левосторонний и правосторонний выводы, дерево разбора для цепочки 1111111100 .

19. Найти общие алгоритмы построения по данным КС-грамматикам G_1 и G_2 , порождающим языки L_1 и L_2 , КС-грамматики для

- а) $L_1 \cup L_2$;
- б) $L_1 \cdot L_2$;
- в) L_1^* .

Замечание

$L = L_1 \cdot L_2$ — это конкатенация языков L_1 и L_2 , т.е. $L = \{\alpha\beta \mid \alpha \in L_1, \beta \in L_2\}$; $L = L_1^*$ — это итерация языка L_1 , т.е. объединение $\{\varepsilon\} \cup L_1 \cup L_1 \cdot L_1 \cup L_1 \cdot L_1 \cdot L_1 \cup \dots$

20. Построить КС-грамматику для $L = \{\omega_i 2 \omega_{i+1}^R \mid i \in N, \omega_i = (i)_2 \text{ — двоичное представление числа } i \text{ (без незначащих нулей), } \omega^R \text{ — обращение цепочки } \omega\}$. Построить КС-грамматику для языка L^* (см. замечание к задаче 19).

21. Показать, что грамматика

$$E \rightarrow E+E \mid E^*E \mid (E) \mid i$$

неоднозначна. Как описать этот же язык с помощью однозначной грамматики?

22. Показать, что наличие в приведенной КС-грамматике правил вида

- a) $A \rightarrow AA \mid \alpha$
- b) $A \rightarrow A\alpha A \mid \beta$
- c) $A \rightarrow \alpha A \mid A\beta \mid \gamma,$

где $\alpha, \beta, \gamma \in (T \cup N)^*$, $A \in N$, делает ее неоднозначной. Можно ли преобразовать эти правила таким образом, чтобы полученная эквивалентная грамматика была однозначной?

23. Показать, что грамматика G неоднозначна. Какой язык она порождает? Является ли этот язык однозначным?

$$G: S \rightarrow aAc \mid aB$$

$$B \rightarrow bc$$

$$A \rightarrow b$$

24. Дана КС-грамматика $G = \langle T, N, P, S \rangle$. Предложить алгоритм построения множества

$$X = \{A \in N \mid A \Rightarrow \varepsilon\}.$$

25. Для произвольной КС-грамматики G предложить алгоритм, определяющий, пуст ли язык $L(G)$.

26. Построить приведенную грамматику, эквивалентную данной КС-грамматике.

a) $S \rightarrow aABS \mid bCACd$ $A \rightarrow bAB \mid cSA \mid cCC$ $B \rightarrow bAB \mid cSB$ $C \rightarrow cS \mid c$	b) $S \rightarrow aAB \mid E$ $A \rightarrow dDA \mid \varepsilon$ $B \rightarrow bE \mid f$ $C \rightarrow cAB \mid dSD \mid a$ $D \rightarrow eA$ $E \rightarrow fA \mid g$
--	--

27. Построить неукорачивающую КС-грамматику, эквивалентную данной, применив алгоритм устранения правил с пустой правой частью.

$$S \rightarrow aS \mid Sa \mid C$$

$$C \rightarrow CC \mid bA \mid \varepsilon$$

$$A \rightarrow aB \mid \varepsilon$$

$$B \rightarrow aA \mid a$$

28. Язык называется распознаваемым, если существует алгоритм, который за конечное число шагов позволяет получить ответ о принадлежности любой цепочки этому языку. Если число шагов зависит от длины цепочки и может быть оценено до выполнения алгоритма, язык называется легко распознаваемым. Доказать, что язык, порождаемый неукорачивающей грамматикой, легко распознаваем.

29. Доказать, что любой конечный язык является регулярным языком.

30. Доказать, что нециклическая КС-грамматика порождает конечный язык.

Замечание

Нетерминальный символ $A \in N$ — циклический, если в грамматике существует вывод $A \Rightarrow \xi_1 A \xi_2$. КС-грамматика называется циклической, если в ней имеется хотя бы один циклический символ.

31. Показать, что условие цикличности грамматики (см. задачу 29) не является достаточным условием бесконечности порождаемого ею языка.

32. Доказать, что язык, порождаемый циклической приведенной КС-грамматикой, содержащей хотя бы один эффективный циклический символ, бесконечен.

Замечание

Циклический символ называется эффективным, если $A \Rightarrow \alpha A \beta$, где $|\alpha A \beta| > 1$; иначе циклический символ называется фиктивным.

II. Регулярные грамматики, конечные автоматы, разбор по ДС

1. Дана регулярная грамматика с правилами:

$$S \rightarrow S0 \mid S1 \mid P0 \mid P1$$

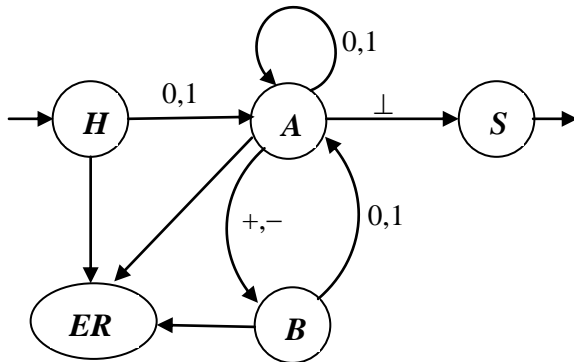
$$P \rightarrow N.$$

$$N \rightarrow 0 \mid 1 \mid N0 \mid N1$$

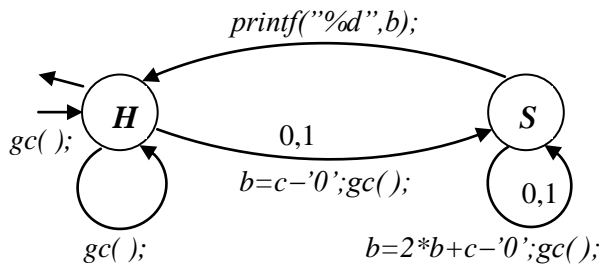
Построить по ней диаграмму состояний (ДС) и использовать ДС для разбора цепочек: 11.010, 0.1, 01., 100. Какой язык порождает эта грамматика?

2. Дана ДС.

- Осуществить разбор цепочек $1011\perp$, $10+011\perp$ и $0-101+1\perp$.
- Восстановить регулярную грамматику, по которой была построена данная ДС.
- Какой язык порождает полученная грамматика?



3. Пусть имеется переменная c и функция $gc()$, считывающая в c очередной символ анализируемой цепочки. Дана ДС с действиями:



- Определить, что будет выдано на печать при разборе цепочки $1+101//p11+++1000/5\perp$.

б) Написать на Си++ анализатор по этой ДС.

4. Построить регулярную (автоматную) левостороннюю грамматику для заданного языка, по ней построить ДС, а по ДС — написать программу анализатора:

а) $L = \{ x\alpha y \perp \mid \alpha \in \{x, y\}^* \}$;

б) $L = \{ (xy^3)^n \perp \mid n \geq 1 \}$;

в) $L = \{ (abb)^k \perp \mid k \geq 1 \}$;

г) $L = \{ \omega \perp \mid \omega \in \{0,1\}^* \text{, где за каждой 1 непосредственно следует 0} \}$;

д) $L = \{ 1\omega 1 \perp \mid \omega \in \{0,1\}^+ \text{, где все подряд идущие 0 — подцепочки нечетной длины} \}$.

5. Дана регулярная грамматика:

$$S \rightarrow A \perp$$

$$A \rightarrow Ab \mid Bb \mid b$$

$$B \rightarrow Aa$$

Определить язык, который она порождает; построить ДС; написать на Си++ анализатор.

6. Построить ДС, по которой в заданном тексте, оканчивающемся на \perp , выявляются все парные комбинации $\langle \rangle$, $\langle =$ и $\rangle =$ и подсчитывается их общее количество.

7. Написать на Си++ анализатор, выделяющий из текста вещественные числа без знака (они определены как в Паскале) и преобразующий их из символьного представления в числовое.

8. Написать на Си++ анализатор, выделяющий из текста вещественные числа (они определены как в Си) и преобразующий их из символьного представления в числовое.

9. Даны две грамматики G_1 и G_2 .

$$\begin{aligned} G_1: \quad S &\rightarrow 0C \mid 1B \\ B &\rightarrow 0B \mid 1C \mid \varepsilon \\ C &\rightarrow 0C \mid 1C \end{aligned}$$

$$\begin{aligned} G_2: \quad S &\rightarrow 0D \mid 1B \\ B &\rightarrow 0C \mid 1C \\ C &\rightarrow 0D \mid 1D \mid \varepsilon \\ D &\rightarrow 0D \mid 1D \end{aligned}$$

Пусть $L_1 = L(G_1)$, $L_2 = L(G_2)$. Построить регулярную (автоматную) грамматику для:

а) $L_1 \cup L_2$

б) $L_1 \cap L_2$

в) $L_1^* \setminus \{\varepsilon\}$

г) $L_2^* \setminus \{\varepsilon\}$

д) $L_1 \cdot L_2$

Если разбор по ней оказался недетерминированным, построить эквивалентную ей грамматику, допускающую детерминированный разбор.

10. Построить левостороннюю регулярную грамматику, эквивалентную данной правосторонней, допускающую детерминированный разбор.

- | | |
|--|--|
| a) $S \rightarrow 0S \mid 0B$
$B \rightarrow 1B \mid 1C$
$C \rightarrow 1C \mid \perp$ | b) $S \rightarrow 0B$
$B \rightarrow 1C \mid 1S$
$C \rightarrow \perp$ |
| c) $S \rightarrow aB$
$B \rightarrow aC \mid aD \mid dB$
$C \rightarrow aB$
$D \rightarrow \perp$ | |

11. Для данной грамматики

- a) определить ее тип;
- b) определить порождаемый ею язык;
- c) построить эквивалентную автоматную грамматику;
- d) построить ДС и анализатор на Си++.

$$\begin{aligned}
 S &\rightarrow 0S \mid S0 \mid D \\
 D &\rightarrow DD \mid 1A \mid \varepsilon \\
 A &\rightarrow 0B \mid \varepsilon \\
 B &\rightarrow 0A \mid 0
 \end{aligned}$$

12. Построить ДС, соответствующую заданной левосторонней автоматной грамматике. Если ДС задает НКА, то построить эквивалентный ему ДКА, используя алгоритм. Для полученного ДКА построить анализатор. Построить соответствующую ДКА грамматику.

- | | |
|--|--|
| a) $S \rightarrow Sa \mid Ab \mid b$
$A \rightarrow Ab \mid Sa \mid a$ | b) $S \rightarrow Sb \mid Aa \mid a$
$A \rightarrow Sb \mid a \mid b$ |
| c) $S \rightarrow C\perp$
$C \rightarrow A1 \mid B1 \mid 1$
$A \rightarrow A1 \mid C0 \mid 0$
$B \rightarrow C0 \mid 0$ | d) $S \rightarrow A\perp$
$A \rightarrow Bb \mid a$
$B \rightarrow Bb \mid b$ |
| e) $S \rightarrow B0 \mid C0$
$B \rightarrow B0 \mid 0$
$C \rightarrow C1 \mid A1$
$A \rightarrow 0$ | f) $S \rightarrow Bb \mid Cc$
$B \rightarrow Bb \mid Ab$
$C \rightarrow Cc \mid Ab$
$A \rightarrow a$ |
| g) $S \rightarrow S1 \mid A0$
$A \rightarrow B1 \mid C1$
$B \rightarrow A0$
$C \rightarrow C0 \mid 0$ | h) $S \rightarrow Sa \mid Cc \mid a$
$C \rightarrow Bb$
$B \rightarrow Sa \mid a$ |

- i) $S \rightarrow C\perp$
 $C \rightarrow A1 \mid B1 \mid 1$
 $A \rightarrow A1 \mid C0 \mid 0$
 $B \rightarrow C0 \mid 0$
- j) $S \rightarrow A\perp$
 $A \rightarrow Bb \mid a$
 $B \rightarrow Bb \mid b$
- k) $S \rightarrow C\perp$
 $B \rightarrow B1 \mid 0 \mid D0$
 $C \rightarrow B1 \mid C1$
 $D \rightarrow D0 \mid 0$
- l) $S \rightarrow C\perp$
 $C \rightarrow B1$
 $B \rightarrow 0 \mid D0$
 $D \rightarrow B1$
- m) $S \rightarrow A0$
 $A \rightarrow A0 \mid S1 \mid 0$
- n) $S \rightarrow B0 \mid 0$
 $B \rightarrow B0 \mid C1 \mid 0 \mid 1$
 $C \rightarrow B0$
- o) $S \rightarrow A0 \mid A1 \mid B1 \mid 0 \mid 1$
 $A \rightarrow A1 \mid B1 \mid 1$
 $B \rightarrow A0$
- p) $S \rightarrow S0 \mid A1 \mid 0 \mid 1$
 $A \rightarrow A1 \mid B0 \mid 0 \mid 1$
 $B \rightarrow A0$
- r) $S \rightarrow Sb \mid Aa \mid a \mid b$
 $A \rightarrow Aa \mid Sb \mid a$

13. Грамматика G определяет язык $L=L_1 \cup L_2$, причем $L_1 \cap L_2 = \emptyset$. Построить регулярную (автоматную) грамматику G_1 , которая порождает язык $L_1 \cdot L_2$ (см. замечание к задаче 19 раздела I). Для нее построить ДС и анализатор.

$$\begin{aligned} S &\rightarrow A\perp \\ A &\rightarrow A0 \mid A1 \mid B1 \\ B &\rightarrow B0 \mid C0 \mid 0 \\ C &\rightarrow C1 \mid 1 \end{aligned}$$

14. Даны две грамматики G_1 и G_2 , порождающие языки L_1 и L_2 .

$$\begin{aligned} G_1: \quad S &\rightarrow S1 \mid A0 \\ \quad A &\rightarrow A1 \mid 0 \end{aligned} \qquad \begin{aligned} G_2: \quad S &\rightarrow A1 \mid B0 \mid E1 \\ \quad A &\rightarrow S1 \\ \quad B &\rightarrow C1 \mid D1 \\ \quad C &\rightarrow 0 \\ \quad D &\rightarrow B1 \\ \quad E &\rightarrow E0 \mid 1 \end{aligned}$$

Построить регулярные (автоматные) грамматики для языков

- a) $L_1 \cup L_2$;
 b) $L_1 \cap L_2$;
 c) $L_1 \cdot L_2$ (см. замечание к задаче 19 раздела I).

Для полученной грамматики построить ДС и анализатор.

15. По данной грамматике G_1 построить регулярную грамматику G_2 для языка $L_1 \cdot L_1$ (см. замечание к задаче 19 раздела I), где $L_1 = L(G_1)$; по грамматике G_2 построить ДС и анализатор.

$$G_1: \quad S \rightarrow S1 \mid A1 \\ A \rightarrow A0 \mid 0$$

16. Построить лексический блок (преобразователь) для кода Морзе. Входом служит последовательность "точек", "тире" и "пауз" (например, ..— . — ...— ⊥). Выходом являются соответствующие буквы, цифры и знаки пунктуации. Особое внимание обратить на организацию таблицы.

III. Метод рекурсивного спуска. КС-грамматики с действиями

1. Определить, применим ли РС-метод к грамматике. Ответ обосновать.

a) $S \rightarrow cA \mid B \mid d$
 $A \rightarrow abA \mid c \mid \varepsilon$
 $B \rightarrow bSc \mid aAb$

b) $S \rightarrow aAbc \mid A$
 $A \rightarrow bB \mid cBc$
 $B \rightarrow bcB \mid a \mid \varepsilon$

c) $S \rightarrow aSB \mid bAf \mid \varepsilon$
 $A \rightarrow bAc \mid cS$
 $B \rightarrow cB \mid d$

d) $S \rightarrow aSB \mid bA$
 $A \rightarrow aS \mid cA \mid \varepsilon$
 $B \rightarrow bB \mid d$

e) $S \rightarrow bABCb \mid d$
 $A \rightarrow aA \mid cB \mid \varepsilon$
 $B \rightarrow Sc$
 $C \rightarrow a \{ bb \}$

f) $S \rightarrow aAb \mid cC$
 $A \rightarrow a \{ bab \}$
 $B \rightarrow cAc \mid aB \mid \varepsilon$
 $C \rightarrow Bb$

g) $S \rightarrow aA \{ xx \}$
 $A \rightarrow bA \mid cBx \mid \varepsilon$
 $B \rightarrow bSc$

h) $S \rightarrow aSc \mid bA \mid \varepsilon$
 $A \rightarrow cS \{ da \} bA \mid d$

i) $S \rightarrow bS \mid aAB$
 $A \rightarrow bcA \mid ccA \mid \varepsilon$
 $B \rightarrow cbB \mid \varepsilon$

j) $S \rightarrow aASb \mid cfAd$
 $A \rightarrow bA \mid c \mid \varepsilon$

k) $S \rightarrow A \mid B$
 $A \rightarrow bA \mid \varepsilon$
 $B \rightarrow cB \mid b \mid \varepsilon$

l) $S \rightarrow AS \mid B$
 $A \rightarrow b \mid c$
 $B \rightarrow dB \mid a \mid \varepsilon$

2. Пусть имеется анализатор в виде системы рекурсивных процедур, построенных по некоторой грамматике в соответствии с методом рекурсивного спуска (S — начальный символ грамматики).

```

#include <iostream>
using namespace std;

int c; // текущий символ

void S(); // объявления процедур, соответствующих нетерминалам грамматики
void A();
...

void gc() {cin >> c;} // считать очередной символ

void S() { ... } // реализация процедур РС-метода
void A() { ... }
...

int main() {

    try {
        gc(); S(); if ( c != '\1' ) throw c;
        cout << "SUCCESS !!!" << endl;
        return 0;
    }
    catch (int c) {
        cout << "ERROR on lexeme " << c << endl;
        return 1;
    }
}

```

Восстановить грамматику по функциям, реализующим синтаксический анализ методом рекурсивного спуска. Удовлетворяет ли полученная грамматика критерию применимости метода рекурсивного спуска?

a)

```

void S() { A(); if ( c != 'd' ) throw c; }

void A() { B(); while ( c == 'a' ) { gc(); B(); } }

void B() { if ( c == 'b' ) gc(); }

```

b)

```

void S() { if ( c == 'a' ) { gc(); A(); }
           else if ( c == 'b' ) { gc(); B(); } else throw c;
}

void A() { if ( c == 'c' ) { gc(); S(); } }

void B() { while ( c == ',' ) { gc(); if ( c != 'b' ) throw c; gc(); }
}

```


c)

```
void s () { if (c == 'a') { gc(); s(); if (c == 'b') gc();
                else throw c;}
            else A();
        }
void A () { if (c == 'b') gc(); else throw c;
            while (c == 'b') gc();
        }
```

d)

```
void s () { A(); if (c != 'd') throw c; }
void A () { B(); while (c == 'a') { gc(); B(); } B(); }
void B () { if (c == 'b') gc(); }
```

e)

```
void s () { if (c == 'a' || c == 'b') { A(); s(); }
            else if (c == 'd') B();
        }
void A () { if (c == 'a') gc();
            else if (c == 'b') { gc(); B(); }
        }
void B () { while (c == 'c') { gc(); B(); } }
```

3. Задана КС-грамматика $G = \langle T, N, P, S \rangle$. По ней написать синтаксический анализатор, реализующий РС-метод, предварительно преобразовав заданную грамматику, если это требуется для применимости РС-метода и если это возможно.

a) $S \rightarrow bS \mid aAB$
 $A \rightarrow bcA \mid ccA \mid \varepsilon$
 $B \rightarrow cbB \mid \varepsilon$

b) $S \rightarrow aASb \mid cfAd$
 $A \rightarrow bA \mid c \mid \varepsilon$

c) $S \rightarrow Sa \mid Sbb \mid fAc$
 $A \rightarrow aB \mid d$
 $B \rightarrow abB \mid Sb$

d) $S \rightarrow cAd$
 $A \rightarrow Aa \mid bB$
 $B \rightarrow abB \mid \varepsilon$

e) $S \rightarrow E \perp$
 $E \rightarrow () \mid (E \{ , E \}) \mid A$
 $A \rightarrow a \mid b$

f) $S \rightarrow P := E \mid \text{if } E \text{ then } S \mid$
 $\quad \text{if } E \text{ then } S \text{ else } S$
 $P \rightarrow I \mid I(E)$
 $E \rightarrow T \{ +T \}$
 $T \rightarrow F \{ *F \}$
 $F \rightarrow P \mid (E)$
 $I \rightarrow a \mid b$

- g) $F \rightarrow \mathbf{function} I(I) S; I:=E \mathbf{end}$ h) $S \rightarrow SaAb \mid Sb \mid bABa$
 $S \rightarrow ; I:=E S \mid \varepsilon$ $A \rightarrow acAb \mid cA \mid \varepsilon$
 $E \rightarrow E*I \mid E+I \mid I$ $B \rightarrow bB \mid \varepsilon$
- i) $S \rightarrow Ac \mid dBea$ j) $S \rightarrow fASd \mid \varepsilon$
 $A \rightarrow Aa \mid Ab \mid daBc$ $A \rightarrow Aa \mid Ab \mid dB \mid f$
 $B \rightarrow cB \mid \varepsilon$ $B \rightarrow bcB \mid \varepsilon$

4. Какой язык задает данная грамматика с действиями? Провести анализ цепочки $(a,(b,a),(a,(b)),b)\perp$.

$$S \rightarrow \langle k = 0 \rangle E \perp$$

$$E \rightarrow A \mid (\langle k = k+1; \mathbf{if} (k == 3) \mathbf{ERROR}(); \rangle E \{,E\} \langle k = k-1 \rangle$$

$$A \rightarrow a \mid b$$

Замечание

Функция *ERROR()* сообщает об ошибке и завершает работу программы.

5. Есть грамматика, описывающая цепочки в алфавите $\{0, 1, 2, \perp\}$:

$$S \rightarrow A\perp$$

$$A \rightarrow 0A \mid 1A \mid 2A \mid \varepsilon$$

Дополнить эту грамматику действиями, исключаящими из языка все цепочки, содержащие подцепочки 002.

6. Дана грамматика, описывающая цепочки в алфавите $\{a, b, c, \perp\}$:

$$S \rightarrow A\perp$$

$$A \rightarrow aA \mid bA \mid cA \mid \varepsilon$$

Дополнить эту грамматику действиями, исключаящими из языка все цепочки, в которых не выполняется хотя бы одно из условий:

- в цепочку должно входить не менее трех букв *c* ;
- если встречаются подряд две буквы *a*, то непосредственно за ними обязательно должна идти буква *b*.

7. Есть грамматика, описывающая цепочки в алфавите $\{0, 1\}$:

$$S \rightarrow 0S \mid 1S \mid \varepsilon$$

Дополнить эту грамматику действиями, исключаящими из языка любые цепочки, содержащие подцепочку 101.

8. Построить КС-грамматику с действиями для порождения

$$L = \{a^m b^n c^k \mid m + k = n \text{ либо } m - k = n\}.$$

9. Построить КС-грамматику с действиями для порождения

$$L = \{1^n 0^m 1^p \mid n+p > m, m \geq 0\}.$$

10. Дана грамматика с семантическими действиями:

$$S \rightarrow \langle A = 0; B = 0 \rangle L \{L\} \langle \mathbf{if} (A > 5) \mathbf{ERROR}() \rangle \perp$$

$$L \rightarrow a \langle A = A + 1 \rangle \mid b \langle B = B + 1; \mathbf{if} (B > 2) \mathbf{ERROR}() \rangle \mid$$

$$c \langle \mathbf{if} (B == 1) \mathbf{ERROR}() \rangle$$

Какой язык описывает эта грамматика? (см. замечание к задаче 4)

11. Дана грамматика:

$$S \rightarrow E \perp$$

$$E \rightarrow () \mid (E \{, E\}) \mid A$$

$$A \rightarrow a \mid b$$

Вставить в заданную грамматику действия, контролирующие соблюдение следующих условий:

- уровень вложенности скобок не больше четырех;
- на каждом уровне вложенности происходит чередование скобочных и бесскобочных элементов.

12. Включить в правила вывода действия, проверяющие выполнение следующих контекстных условий:

а) Пусть в языке L есть переменные и константы целого, вещественного и логического типов, а также есть оператор цикла

$$S \rightarrow \mathbf{for} I = E \mathbf{step} E \mathbf{to} E \mathbf{do} S$$

Включить в это правило вывода действия, проверяющие выполнение следующих ограничений:

- тип I и всех вхождений E должен быть одинаковым;
- переменная логического типа недопустима в качестве параметра цикла.

Для каждой используемой процедуры привести ее текст на Си++.

б) Дан фрагмент грамматики

$$P \rightarrow \mathbf{program} D; \mathbf{begin} S \{; S \} \mathbf{end}$$

$$D \rightarrow \dots \mid \mathbf{label} L \{, L\} \mid \dots$$

$$S \rightarrow L \{ , L \} : S' \mid S'$$

$$S' \rightarrow \dots \mid \mathbf{goto} L \mid \dots$$

$$L \rightarrow I$$

где I — идентификатор.

Вставить в грамматику действия, контролирующие выполнение следующих условий:

- каждая метка, используемая в программе, должна быть описана и только один раз;
- не должно быть одинаковых меток у различных операторов;
- если метка используется в операторе **goto**, то обязательно должен быть оператор, помеченный такой меткой.

Для каждой используемой процедуры привести ее текст на Си++.

IV. Синтаксически управляемый перевод

1. Построить грамматику арифметического выражения, использующего операции $+$, $-$, $*$, $/$ и круглые скобки (приоритет операций стандартный), простые аргументы операций — переменные a и b , например: $a+(b-a)*b/a+b$. Предполагая, что анализ грамматики будет производиться РС-методом, вставить в нее действия вида $cout \ll$ 'символ' по переводу таких выражений в ПОЛИЗ.

2. Дана грамматика языка L_1 , в которую вставлены действия по переводу цепочек языка L_1 в цепочки языка L_2 . Определить языки L_1 и L_2 .

$$\begin{aligned} \text{a)} \quad S &\rightarrow a \langle a = 1; b = 0; \rangle A \perp \\ A &\rightarrow a \langle \text{if} (a) \{ cout \ll 'a'; a = 0; \} \text{ else } a++; \rangle A \mid \\ &\quad bA \langle \text{if} (b) \{ cout \ll 'b'; b = 0; \} \text{ else } b++; \rangle \mid \varepsilon \end{aligned}$$

$$\begin{aligned} \text{b)} \quad S &\rightarrow \langle a = 0; \rangle E \perp \langle cout \ll '\perp'; \rangle \\ E &\rightarrow a \langle a = 1; \rangle E \langle cout \ll 'a'; \rangle \mid \\ &\quad b \langle \text{if} (a == 0) cout \ll 'b'; \rangle E \langle cout \ll 'b'; \rangle \mid \varepsilon \end{aligned}$$

3. Построить грамматику для языка L_1 . Вставить в нее действия по генерации цепочек языка L_2 в процессе анализа **методом рекурсивного спуска**. Соответствие между цепочками задается формальным переводом τ . В качестве действий можно использовать только операторы вида $cout \ll$ 'символ'.

$$\begin{aligned} \text{a)} \quad L_1 &= \{ a^n c^m b^n \mid n \geq 0, m \geq 1 \}, \quad L_2 = \{ 0^i 1^k \mid i \geq 0, k > i \}, \\ \tau &= \{ (a^n c^m b^n, 0^n 1^{n+m}) \mid n \geq 0, m \geq 1 \}; \end{aligned}$$

$$\begin{aligned} \text{b)} \quad L_1 &= \{ \alpha c^n \mid \alpha \in \{a, b\}^*, n \geq 1 \}, \quad L_2 = \{ a^n c^m \mid n \geq 1, m \geq 0 \}, \\ \tau &= \{ (\alpha c^n, a^n c^m) \mid \alpha \in \{a, b\}^*, n \geq 1, m = |\alpha|_a \text{ (т.е. } m \text{ — количество символов } a \text{ в} \\ &\quad \text{цепочке } \alpha \text{)} \}; \end{aligned}$$

$$\begin{aligned} \text{c)} \quad L_1 &= \{ a, b \}^+, \quad L_2 = \{ 1^i 0^j \mid i \geq 1, j \geq 0; i \geq j \}, \\ \tau &= \{ (\omega, 1^{n+m} 0^n) \mid \omega \in \{a, b\}^+, n = |\omega|_a, m = |\omega|_b \}; \end{aligned}$$

$$\begin{aligned} \text{d)} \quad L_1 &= \{ a, b \}^+, \quad L_2 = \{ 2^n \alpha \mid n \geq 0, \alpha \in \{a, b\}^+ \}, \\ \tau &= \{ (\omega, 2^n \omega^R) \mid \omega \in \{a, b\}^+, n = |\omega|_a \text{ (здесь } \omega^R \text{ — реверс цепочки } \omega \text{)} \}; \end{aligned}$$

$$\begin{aligned} \text{e)} \quad L_1 &= \{ 1^n 0^m 1^m 0^n \mid m, n > 0 \}, \quad L_2 = \{ 1^i 0^k \mid k > i > 0 \}, \\ \tau &= \{ (1^n 0^m 1^m 0^n, 1^m 0^{n+m}) \mid m, n > 0 \}; \end{aligned}$$

$$\begin{aligned} \text{f)} \quad L_1 &= \{ \alpha_1 \alpha_2 \dots \alpha_n \perp \mid n \geq 1, \alpha_i \in \{ab, ba\} \text{ для } 1 \leq i \leq n \}, \quad L_2 = \{ \omega \perp \mid \omega \in \{a, b\}^+ \}, \\ \tau &= \{ (\alpha_1 \alpha_2 \dots \alpha_n \perp, \beta_1 \beta_2 \dots \beta_n \perp) \mid n \geq 1; \text{ для } 1 \leq i \leq n \quad \alpha_i \in \{ab, ba\}, \beta_i = b, \\ &\quad \text{если } \alpha_i = ab, \beta_i = a, \text{ если } \alpha_i = ba \}; \end{aligned}$$

4. Построить грамматику для языка L_1 . Вставить в нее действия по генерации цепочек языка L_2 в процессе анализа **методом рекурсивного спуска**. Соответствие между цепочками задается формальным переводом τ . В качестве действий можно использовать любые операторы.

$$a) \quad L_1 = \{ 1^m 0^n \mid n, m > 0 \}, \quad L_2 = \{ 1^k \mid k > 0 \} \cup \{ 0^i \mid i > 0 \} \cup \{ \varepsilon \},$$

$$\tau = \{ (1^m 0^n, 1^{m-n}) \mid m > n > 0 \} \cup \{ (1^m 0^n, 0^{n-m}) \mid n > m > 0 \} \cup \{ (1^m 0^n, \varepsilon) \mid m = n \};$$

$$b) \quad L_1 = \{ \omega_i \mid i \geq 0, \omega_i = (i)_2 \text{ (т.е. } \omega_i \text{ — это двоичное представление числа } i \in N \text{ без незначащих ведущих нулей)} \},$$

$$L_2 = \{ (\omega_i)^R \mid i \geq 1, \omega_i = (i)_2 \text{ (} \omega^R \text{ — обращение цепочки } \omega) \},$$

$$\tau = \{ (\omega_i, (\omega_{i+1})^R) \mid i \geq 0, \omega_i = (i)_2, \omega_{i+1} = (i+1)_2 \};$$

$$c) \quad L_1 = \{ \alpha \perp \mid \alpha \in \{a, b\}^* \}, \quad L_2 = \{ b^n \beta \perp \mid n \geq 0, \beta \in \{a, b\}^* \},$$

$$\tau = \{ (\alpha \perp, b^n \alpha^R \perp) \mid \alpha \in \{a, b\}^*, n = |\alpha|_a \text{ (} n \text{ — количество символов } a \text{ в цепочке } \alpha, \alpha^R \text{ — реверс цепочки } \alpha) \};$$

$$d) \quad L_1 = \{ \omega \perp \mid \omega \in \{a, b\}^+ \}, \quad L_2 = \{ a^i b^k \mid i, k \geq 0, i+k > 0 \},$$

$$\tau = \{ (\omega \perp, a^{\lfloor n/2 \rfloor} b^{\lfloor m/2 \rfloor}) \mid \omega \in \{a, b\}^+, n = |\omega|_a, m = |\omega|_b \};$$

$$e) \quad L_1 = \{ \omega \perp \mid \omega \in \{a, b\}^+ \}, \quad L_2 = \{ a^i b^k \mid i, k \geq 0, i+k > 0 \},$$

$$\tau = \{ (\omega \perp, a^{\lfloor (n+1)/3 \rfloor} b^{\lfloor (m-n)/2 \rfloor}) \mid \omega \in \{a, b\}^+, n = |\omega|_a, m = |\omega|_b \};$$

$$f) \quad L_1 = \{ \omega \perp \mid \omega \in \{0, 1\}^+, \omega = (i)_2^R, i \geq 0 \text{ (т. е. } \omega \text{ — реверс двоичной записи числа } i) \},$$

$$L_2 = \{ |^n \mid n \geq 0 \}, \quad \tau = \{ (\omega \perp, |^i) \mid \omega \in \{0, 1\}^+, \omega = (i)_2^R, i \geq 0 \}.$$

5. Построить грамматику, описывающую целые двоичные числа (количество 0 и 1 четно, допускаются незначащие нули). Вставить в нее действия по переводу этих целых чисел в четверичную систему счисления.

6. Построить грамматику для выражений, содержащих переменные, знаки операций $+$, $-$, $*$, $/$, $**$ и скобки $()$ с обычным приоритетом операций и скобок. Включить в эту грамматику действия по переводу выражений в префиксную запись (операции предшествуют операндам). Предложить интерпретатор префиксной записи выражений.

7. В грамматику, описывающую выражения, включить действия по переводу выражений из инфиксной формы (операция между операндами) в функциональную запись.

Например,
$$a+b \quad \Rightarrow \quad +(a, b),$$

$$a+b*c \quad \Rightarrow \quad +(a, *(b, c)).$$

V. ПОЛИЗ, перевод в ПОЛИЗ

1. Представить в ПОЛИЗе следующие выражения:

$$a) \quad a+b-c$$

$$b) \quad a*b+c/a$$

- | | |
|-------------------------------------|---|
| c) $a/(b+c)*a$ | d) $(a+b)/(c+a*b)$ |
| e) $a \text{ and } b \text{ or } c$ | f) $\text{not } a \text{ or } b \text{ and } a$ |
| g) $x+y=x/y$ | h) $(x*x+y*y<1) \text{ and } (x>0)$ |

2. Для следующих выражений в ПОЛИЗе дать обычную инфиксную запись:

- | | | |
|----------------|---------------------------------------|------------------------------|
| a) $ab*c+$ | b) $abc*/$ | c) $ab+c*$ |
| d) $ab+bc-/a+$ | e) $a \text{ not } b \text{ and not}$ | f) $abca \text{ and or and}$ |
| g) $2x+2x*<$ | | |

3. Используя стек, вычислить следующие выражения в ПОЛИЗе:

- | | |
|--|----------------------------|
| a) $x \ y*x \ y / +$ | при $x = 8, y = 2;$ |
| b) $a \ 2+b / b \ 4*+$ | при $a = 4, b = 3;$ |
| c) $a \ b \ \text{not and } a \ \text{or not}$ | при $a = b = \text{true};$ |
| d) $x \ y*0 > y \ 2 \ x - < \text{and}$ | при $x = y = 1.$ |

4. Перевести в ПОЛИЗ фрагмент программы на Си:

- | |
|--|
| a) $S=0; i=1; \text{while} (i < 10) \{ S = S*(i+i); i++; \}$ |
| b) $\text{if} ((x+1) > (2*y)) x = y; \text{else } y = (x+y)*2;$ |
| c) $i = 1; S = 0; \text{while} (i < 10 \ \&\& \ S < 40) \{ S = S + f(i); ++i; \}$ |
| d) $\text{if} (z<x*y+5) a = x<y, z = (x+6)/(a-y); \text{else } z = y << 2;$ |
| e) $a = x + y < z*(t + x) ? - (a + b)/(c-d)*2 : ++x+5;$ |
| f) $S = x + y; i = 1; \text{for} (j = 0; j < n; j++) \{ S = S + i*j*S; i = i*x; \}$ |
| g) $i = 1; S = 0; \text{while} (i < 10 \ \&\& \ S < 40) \{ S = S + f(i); ++i; \}$ |
| h) $\text{if} (z<x*y+5) a = x<y, z = (x+6)/(a-y); \text{else } z = y << 2;$ |
| i) $\text{do} \{ x = y; y = 2*y; \} \text{while} (x < k);$ |
| j) $S = 0; \text{for} (i = 1; i <= k; i = i + 1) S += i*i;$ |
| k) $\text{switch} (k) \{$
<div style="margin-left: 2em;"> $\text{case } 1: a = \text{not } a; \text{break};$
 $\text{case } 2: b = a \ \text{or not } b;$
 $\text{case } 3: a = b ;$ </div> $\}$ |

Замечание

ПОЛИЗ управляющих операторов языка Си составляется аналогично ПОЛИЗу операторов М-языка и Паскаля. При переводе в ПОЛИЗ сложных операторов порядок внутренних конструкций (операторов и выражений) относительно друг друга сохраняется. Например, в ПОЛИЗе для оператора цикла вида $\text{for} (expr1; expr2; expr3) \{body\}$ ПОЛИЗ $expr3$ должен предшествовать ПОЛИЗу $body$.

В языке Си присваивание является операцией, а не оператором, поэтому при интерпретации ПОЛИЗа присвоенное значение сохраняется в стеке (как результат операции). Чтобы удалить ненужное значение с вершины стека (такие значения остаются в стеке, например, после испол-

нения оператора-выражения), в ПОЛИЗе языка Си используется операция ";".

Чтобы отличить префиксные операции ++ и -- от постфиксных, префиксные обозначают в ПОЛИЗе как #+ и #-, а постфиксные как #+ и #- соответственно.

ПОЛИЗ вызова функции представляет собой последовательность ее аргументов в ПОЛИЗе, за которыми следует имя функции. Для функций с переменным числом параметров перед именем функции в ПОЛИЗ вставляется дополнительный аргумент — количество фактических параметров в данном вызове функции. При интерпретации сначала из стека извлекается этот дополнительный аргумент, а затем — соответствующее ему число фактических параметров.

5. По заданному ПОЛИЗу выражения записать его в инфиксной форме (на Си).

а) $\underline{x} \underline{y} \underline{z} a x 5 y / + * z 6 + 8 * - = = = ;$ б) $\underline{x} a x z y / + * z 6 a - * + =$

6. Является ли запись

ПОЛИЗ	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	
	\underline{y}	15	=	;	\underline{x}	x	a	b	c	2	/	1	+	*	-	*	a	
	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35
	-	=	;	\underline{y}	y	2	-	=	;	y	10	<=	5	!F				

правильной записью в ПОЛИЗе следующего фрагмента программы на Си (считаем, что элементы ПОЛИЗа нумеруются с 1): $y = 15; \mathbf{do} \{x = x*(a-b*(c/2+1))-a; y = y-2;\} \mathbf{while} (y > 10);$ Если не является, объясните почему, и предложите свой вариант ПОЛИЗа для этого фрагмента программы.

7. Является ли запись

ПОЛИЗ	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17		
	\underline{i}	1	=	;	i	n	<	36	!F	\underline{a}	a	b	+	1	-	x	y		
	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36
	y	2	+	/	-	*	5	+	=	;	\underline{i}	i	2	+	=	;	5	!	

правильной записью в ПОЛИЗе следующего фрагмента программы на Си:

$$\mathbf{for} (i = 1; i < n; i = i + 2) a = (a+b-1)*(x-y/(y+2))+5;$$

Если не является, объясните почему и предложите свой вариант ПОЛИЗа для этого фрагмента программы.

8. а) перевести в ПОЛИЗ фрагмент программы на Си:

$$i = 1; s = 0; \mathbf{while} (i < 10 \ \&\& \ s < 40) \{ s = s + f(i); ++i; \};$$

с) выражение в ПОЛИЗе записать в инфиксной форме (на Си).

$$\underline{x} \underline{y} \underline{z} a x 5 y / + * z 6 + 8 * - = = = ;$$

9. а) перевести в ПОЛИЗ фрагмент программы на Си:

$$\mathbf{if} (z < x*y+5) a = x < y, z = (x+6)/(a-y); \mathbf{else} z = y << 2;$$

б) выражение в ПОЛИЗе записать в инфиксной форме (на Си).

$$\underline{x} a x z y / + * z 6 a - * + = ;$$

10. Предложить ПОЛИЗ для следующих операторов:

a) *if* (E) S_1 ; S_2 ; S_3

Семантика этого оператора такова: вычисляется значение выражения E ; если его значение меньше 0, то выполняется оператор S_1 ; если равно 0 — оператор S_2 , иначе — оператор S_3 .

b) *choice* (S_1 ; S_2 ; S_3), E

Семантика: вычисляется значение выражения E ; если его значение равно i , то выполняется оператор S_i для $i = 1, 2, 3$; иначе оператор *choice* эквивалентен пустому оператору.

c) *cycle* (E_1 ; E_2 ; E_3), S

Семантика этого оператора отличается от семантики оператора *for* в языке Си только тем, что оператор S выполняется, по крайней мере, один раз (т.е. после вычисления выражения E_1 сразу выполняется оператор S , затем вычисляется значение E_3 , потом — значение E_2 , которое используется для контроля за количеством повторений цикла также, как и в цикле *for*).

11. Построить грамматику для выражений, содержащих переменные, знаки операций $+$, $-$, $*$, $/$ и скобки $()$. Грамматика должна отражать **одинаковый** приоритет и лево-ассоциативность всех четырех операций. Определить действия по переводу таких выражений в ПОЛИЗ.

12. Изменить приоритет операций отношения в М-языке — сделать его наивысшим. Построить соответствующую грамматику, отражающую этот приоритет. Написать синтаксический анализатор, обеспечить контроль типов, задать перевод в ПОЛИЗ.

13. Построить КС-грамматику, аналогичную данной,

$$E \rightarrow T \{+T\}$$

$$T \rightarrow F \{*F\}$$

$$F \rightarrow (E) | i$$

с той лишь разницей, что в новом языке перед идентификатором будет допускаться унарный минус, имеющий наивысший приоритет (например, $a*-b+-c$ допускается и означает $a*(-b)+(-c)$). В созданную грамматику вставить действия по переводу такого выражения в ПОЛИЗ. Для каждой используемой процедуры привести ее текст на Си++.

14. Дана грамматика, описывающая выражения:

$$E \rightarrow TE'$$

$$T \rightarrow FT'$$

$$F \rightarrow PF'$$

$$P \rightarrow (E) | i$$

$$E' \rightarrow + TE' | \varepsilon$$

$$T' \rightarrow * FT' | \varepsilon$$

$$F' \rightarrow ^ PF' | \varepsilon$$

Включить в эту грамматику действия по переводу этих выражений в ПОЛИЗ. Для каждой используемой процедуры привести ее текст на Си++.

Литература

- [1]. Д. Грис. Конструирование компиляторов для цифровых вычислительных машин. — М.: Мир, 1975.
- [2]. Ф. Льюис, Д. Розенкранц, Р. Стирнз. Теоретические основы проектирования компиляторов. — М.: Мир, 1979.
- [3]. А. Ахо, Дж. Ульман. Теория синтаксического анализа, перевода и компиляции. — Т. 1,2. — М.: Мир, 1979.
- [4]. Ф. Вайнгартен. Трансляция языков программирования. — М.: Мир, 1977.
- [5]. И. Л. Братчиков. Синтаксис языков программирования. — М.: Наука, 1975.
- [6]. С. Гинзбург. Математическая теория контекстно-свободных языков. — М.: Мир, 1970.
- [7]. Дж. Фостер. Автоматический синтаксический анализ. — М.: Мир, 1975.
- [8]. В. Н. Лебедев. Введение в системы программирования. — М.: Статистика, 1975.
- [9]. Б. Ф. Мельников. Подклассы класса контекстно-свободных языков. — М.: МГУ, 1995.
- [10]. В. Н. Пильщиков, В. Г. Абрамов, А. А. Вылиток, И. В. Горячая. Машины Тьюринга и алгоритмы Маркова. Решение задач. — М.: МГУ, 2006.
- [11]. А. Ахо., Р. Сети, Дж. Ульман. Компиляторы: принципы, технологии, инструменты. — М.: «Вильямс», 2001.
- [12]. А. Ахо, М. Лам, Р. Сети, Дж. Ульман. Компиляторы: принципы, технологии и инструментарий. — М.: «Вильямс», 2008.

Содержание

<i>МГУ им. М. В. Ломоносова</i>	2
Рецензенты:.....	2
Элементы теории формальных языков и грамматик.....	3
Введение.....	3
Основные понятия и определения.....	3
Классификация грамматик и языков по Хомскому.....	7
Граматики с ограничениями на вид правил вывода.....	7
Иерархия Хомского.....	9
Примеры грамматик и языков.....	11
Регулярные.....	12
Контекстно-свободные.....	12
Неукорачивающие и контекстно-зависимые.....	13
Без ограничений на вид правил (тип 0).....	13
Замечание о связи между языком и грамматикой.....	14
Разбор цепочек.....	15
Преобразования грамматик.....	19
Алгоритм удаления недостижимых символов.....	19
Алгоритм удаления бесплодных символов.....	19
Алгоритм приведения грамматики.....	20
Алгоритм устранения правил с пустой правой частью.....	20
Элементы теории трансляции.....	21
Введение.....	21
Разбор по регулярным грамматикам.....	22
Алгоритм разбора по диаграмме состояний.....	24
Пример разбора цепочки.....	27
О недетерминированном разборе.....	28
Регулярные выражения.....	34
Задачи лексического анализа.....	35
Лексический анализатор для М-языка.....	37
Синтаксический анализ.....	47
Метод рекурсивного спуска.....	48
Нисходящий анализ с прогнозируемым выбором альтернатив.....	52
О применимости метода рекурсивного спуска.....	53
Задача разбора для неоднозначных грамматик.....	65
О других методах распознавания КС-языков.....	66
Синтаксический анализатор для М-языка.....	67
Семантический анализатор для М-языка.....	74
Генерация внутреннего представления программ.....	81
Язык внутреннего представления программы.....	81
Синтаксически управляемый перевод.....	85

Генератор внутреннего представления программы на М-языке	87
Интерпретатор ПОЛИЗа для модельного языка	89
Задачи	93
I. Грамматики и языки. Классификация по Хомскому.....	93
II. Регулярные грамматики, конечные автоматы, разбор по ДС	99
III. Метод рекурсивного спуска. КС-грамматики с действиями	103
IV. Синтаксически управляемый перевод	108
V. ПОЛИЗ, перевод в ПОЛИЗ	109
Литература	113