

Определение формальной грамматики

и языка

Определение. Конечное множество символов, неделимых в данном рассмотрении, в теории формальных грамматик называется словарем, или алфавитом, а символы, входящие в множество, - буквами алфавита.

Например, алфавит $A = \{a, b, c, +, !\}$ содержит 5 букв, а алфавит $B = \{00, 01, 10, 11\}$ содержит 4 буквы, каждая из которых состоит из двух символов.

Определение. Последовательность букв алфавита называется словом или цепочкой в этом алфавите. Число букв, входящих в слово, называется его длиной.

Например, слово в алфавите A $a = ab++c$ имеет длину $l(a) = 5$, а слово в алфавите B $b = 00110010$ имеет длину $l(b) = 4$.

Если задан алфавит A , то обозначим A^* множество всевозможных цепочек, которые могут быть построены из букв алфавита A . При этом предполагается, что пустая цепочка, которую обозначим знаком ϵ , также входит в множество A^* .

Определение. Формальной порождающей грамматикой Γ называется следующая совокупность четырех объектов: $\Gamma = \{V_T, V_A, \langle I \rangle \in V_A, R\}$, где V_T - терминальный алфавит (словарь); буквы этого алфавита называются терминальными символами; из них строятся цепочки, порождаемые грамматикой;

V_A - нетерминальный, вспомогательный алфавит (словарь); буквы этого алфавита (нетерминальные символы) используются при построении цепочек; они могут входить в промежуточные цепочки, но не должны входить в результат порождения;

$\langle I \rangle$ - начальный символ грамматики $\langle I \rangle \in V_A$;

R - множество правил вывода, или порождающих правил вида $\alpha \rightarrow \beta$, где α и β - цепочки, построенные из букв алфавита $V_T \cup V_A$, который называют полным алфавитом (словарем) грамматики Γ .

В множество правил грамматики могут также входить правила с пустой правой частью вида $\langle E \rangle \rightarrow \epsilon$. Чтобы избежать неопределенности из-за отсутствия символа в правой части правила, условимся использовать символ пустой цепочки, записывая такое правило в виде $\langle E \rangle \rightarrow \epsilon$.

Чтобы установить правила построения цепочек, порождаемых грамматикой, введем следующие понятия.

Определение. Пусть $r = \alpha \rightarrow \gamma$ - правило грамматики Γ и $\alpha = x' t x''$ - цепочка символов, причем $x', x'' \in (V_T \cup V_A)^*$. Тогда цепочка $\beta = x' \gamma x''$ может быть получена из цепочки α путем применения правила r (т.е. заменой в цепочке α t на γ). В этом случае говорят, что цепочка β непосредственно выведена из цепочки α и обозначают $\alpha \Rightarrow \beta$.

Определение. Если задана совокупность цепочек $\Omega = (w_0, w_1, \dots, w_n)$, таких что существует последовательность непосредственных выводов:

$$w_0 \Rightarrow w_1, w_1 \Rightarrow w_2, \dots, w_{n-1} \Rightarrow w_n,$$

то такую последовательность называют **выводом** w_n из w_0 в грамматике Γ и обозначают

$$w_0 \Rightarrow^* w_n$$

Определение. Множество конечных цепочек терминального алфавита V_T грамматики Γ , выводимых из начального символа $\langle I \rangle$, называется языком, порождаемым грамматикой Γ и обозначается $L(\Gamma)$.

$$L(\Gamma) = \{ w \in V_T^* \mid \langle I \rangle \Rightarrow^* w \}$$

Определение. Если язык, порождаемый грамматикой Γ , не содержит ни одной конечной цепочки (конечного слова), то он называется пустым.

Утверждение. Для того, чтобы язык $L(\Gamma)$ не был пустым, в множестве R должно быть хотя бы одно правило вида $r = \chi \rightarrow \psi$, где $\psi \in V_T^*$ и должен существовать вывод

$$\langle I \rangle \Rightarrow^* \chi$$

Типы формальных языков и грамматик. Классификация по Хомскому

В теории формальных языков выделяются 4 типа грамматик, которым соответствуют 4 типа языков. Эти грамматики выделяются путем наложения усиливающих ограничений на правила грамматики

Граматики типа 0

Граматики типа 0, которые называют грамматиками общего вида, не имеют никаких ограничений на правила порождения. Эти грамматики порождают естественные языки.

Любое правило

$$r = \eta \rightarrow \psi$$

может быть построено с использованием произвольных цепочек η , $\psi \in (V_T \cup V_A)^*$. Например,

$$\langle T \rangle \langle W \rangle \rightarrow \langle W \rangle \langle T \rangle \text{ или } x \langle A \rangle b \langle C \rangle \langle D \rangle \rightarrow x \langle H \rangle \langle D \rangle$$

Граматики типа 1

Граматики типа 1, которые называют также контекстно-зависимыми грамматиками, не допускают использования любых правил. Правила вывода в таких грамматиках должны иметь вид:

$$\chi_1 \langle A \rangle \chi_2 \rightarrow \chi_1 \omega \chi_2,$$

где χ_1, χ_2 - цепочки, возможно пустые, из множества $(V_T \cup V_A)^*$, символ $\langle A \rangle \in V_A$ и цепочка $\omega \in (V_T \cup V_A)^*$. Цепочки χ_1 и χ_2 остаются неизменными при применении правила, поэтому их называют контекстом (соответственно левым и правым), а грамматику - контекстно-зависимой.

Граматики типа 1 значительно удобнее на практике, чем грамматики типа 0, поскольку в левой части правила заменяется всегда один нетерминальный символ, который можно связать с некоторым синтаксическим понятием, в то время как в грамматике типа 0 можно заменять сразу несколько символов, в том числе и терминальных.

Например, грамматика:

$\Gamma_{1.4}$:

$$V_T = \{a, b, c, d\}, V_A = \{\langle I \rangle, \langle A \rangle, \langle B \rangle\}$$

$$R = \{ \langle I \rangle \rightarrow aA\langle I \rangle,$$

$$A\langle I \rangle \rightarrow AA\langle I \rangle,$$

$$AAA \rightarrow A\langle B \rangle A,$$

$$A \rightarrow b,$$

$$b\langle B \rangle A \rightarrow bcdA,$$

$$b\langle I \rangle \rightarrow ba \}$$

является контекстно-зависимой, поскольку второе и шестое правила имеют непустой левый контекст, а третье и пятое правила содержат оба контекста. Вывод в такой грамматике может иметь вид:

$$\langle I \rangle \Rightarrow a\langle A \rangle \langle I \rangle \Rightarrow a\langle A \rangle \langle A \rangle \langle I \rangle \Rightarrow ab\langle A \rangle \langle I \rangle \Rightarrow abb\langle I \rangle \Rightarrow abba$$

Граматики типа 2

Граматики типа 2 называют контекстно-свободными и бесконтекстными грамматиками (КС-грамматики или Б-грамматики).

Правила вывода таких грамматик имеют вид:

$$\langle A \rangle \rightarrow \alpha,$$

где $\langle A \rangle \in V_A$ и $\alpha \in (V_T \cup V_A)^*$.

Очевидно, что эти правила получаются из правил грамматики типа 1 при условии $\chi_1 = \chi_2 = \$$. Поскольку контекстные условия отсутствуют, то правила КС-грамматик получаются проще, чем правила грамматик типа 1. Именно такие грамматики используются для описания языков программирования.

Примером КС-грамматики может служить следующая:

$\Gamma_{1.5}$:

$$V_T = \{a, b\}, V_A = \{\langle I \rangle\},$$

$$R = \{ \langle I \rangle \rightarrow a\langle I \rangle a,$$

$$\begin{aligned} \langle I \rangle &\rightarrow b\langle I \rangle b, \\ \langle I \rangle &\rightarrow aa, \\ \langle I \rangle &\rightarrow bb \} \end{aligned}$$

Эта грамматика порождает язык, который состоит из цепочек, каждая из которых в свою очередь состоит из двух частей, цепочки $\beta \in V_T^*$ и зеркального отображения этой цепочки β' .

$$L(\Gamma_5) = \{ bb' \mid b \in V_T^+ \},$$

где V_T^+ - это множество V_T^* без пустой цепочки. С помощью правил этой грамматики может быть построена, например, следующая цепочка:

$$\langle I \rangle \Rightarrow a\langle I \rangle a \Rightarrow ab\langle I \rangle ba \Rightarrow aba\langle I \rangle aba \Rightarrow ababbaba$$

Грамматики типа 3

Грамматики типа 3 называют автоматными грамматиками (А - грамматиками). Правила вывода в таких грамматиках имеют вид:

$$\langle A \rangle \rightarrow a \text{ или } \langle A \rangle \rightarrow a\langle B \rangle \text{ или } \langle A \rangle \rightarrow \langle B \rangle a,$$

где $a \in V_T$, $\langle A \rangle, \langle B \rangle \in V_A$, причем грамматика может иметь только правила вида $\langle A \rangle \rightarrow a\langle B \rangle$ - правосторонние правила, либо только вида $\langle A \rangle \rightarrow \langle B \rangle a$ - левосторонние правила. Примерами автоматных грамматик могут служить правосторонняя грамматика $\Gamma_{1.6}$ и левосторонняя грамматика $\Gamma_{1.7}$.

$$\begin{aligned} \Gamma_{1.6} \\ V_T &= \{a, b\}, V_A = \{\langle I \rangle, \langle A \rangle\}, \\ R &= \{ \langle I \rangle \rightarrow a\langle I \rangle, \\ &\langle I \rangle \rightarrow a\langle A \rangle, \\ &\langle A \rangle \rightarrow b\langle A \rangle, \\ &\langle A \rangle \rightarrow b\langle Z \rangle, \\ &\langle Z \rangle \rightarrow \$ \} \end{aligned}$$

$$\begin{aligned} \Gamma_{1.7} \\ V_T &= \{a, b\}, V_A = \{\langle I \rangle, \langle A \rangle\}, \\ R &= \{ \langle I \rangle \rightarrow \langle A \rangle b, \\ &\langle A \rangle \rightarrow \langle A \rangle b, \\ &\langle A \rangle \rightarrow \langle Z \rangle a, \\ &\langle Z \rangle \rightarrow \langle Z \rangle a, \\ &\langle Z \rangle \rightarrow \$ \} \end{aligned}$$

Эти грамматики являются эквивалентными и порождают язык

$$L(\Gamma_7) = \{ a...ab...b \mid n, m > 0 \}$$

Между множествами языков различных типов существует отношение включения:

$$\{ L_{\text{типа 3}} \} \subset \{ L_{\text{типа 2}} \} \subset \{ L_{\text{типа 1}} \} \subset \{ L_{\text{типа 0}} \}$$

Доказано, что существуют языки типа 0, не являющиеся языками типа 1, языки типа 2, не являющиеся языками типа 1, и языки типа 3, не являющиеся языками типа 2.

Учитывая, что наибольшее практическое применение находят грамматики типа 2 и типа 3, дальнейшее изложение посвящается рассмотрению именно этих типов грамматик.

Вывод в КС-грамматиках и правила построения дерева вывода

Формальные грамматики позволяют задавать языки, представляющие множества цепочек, построенных по определенным правилам. Используемый способ задания позволяет построить любую цепочку, принадлежащую языку. Чтобы сделать процесс построения, называемый выводом, наглядным, его изображают в виде графа, точнее, в виде дерева, которое называют синтаксическим деревом или деревом вывода. Учитывая, что вывод любой цепочки языка, принадлежащей языку, порождаемому заданной грамматикой, должен начинаться с начального символа, правила построения дерева можно сформулировать так:

1. В качестве начальной вершины или корня дерева возьмем вершину, которую обозначим начальным символом грамматики $\langle I \rangle$; эта вершина образует нулевой ярус дерева
2. Если при выводе цепочки на очередном шаге используется правило грамматики $\langle A \rangle \rightarrow a$ и вершина, помеченная нетерминалом $\langle A \rangle$, расположена на ярусе с номером $k-1$, то к построенному дереву нужно добавить столько вершин, сколько содержится символов в цепочке a , расположить эти вершины на ярусе k , обозначить их символами цепочки a и соединить эти вершины дугами с вершиной $\langle A \rangle$. Результатом вывода является множество конечных узлов - листьев, которые выписываются при обходе дерева слева - вниз - направо - вверх. Рассмотрим, например, грамматику $\Gamma_{1.8}$:

$\Gamma_{1.8}$:

$V_T = \{a, b\}, V_N = \{\langle I \rangle\},$

$R = \{\langle I \rangle \rightarrow a\langle I \rangle b,$
 $\langle I \rangle \rightarrow ab \},$

которая порождает язык $L(\Gamma_8) = \{aa...abb...b\}$, где a и b повторяются по n раз, $n=1,2,...$

Вывод цепочки с помощью правил этой грамматики имеет вид:

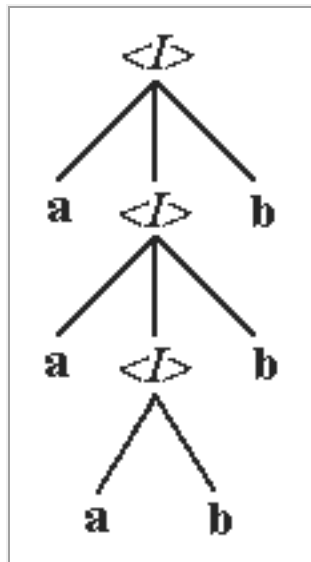


Рис. 1.

Синтаксический разбор

Вывод цепочки с помощью правил грамматики может быть задан не только в виде синтаксического дерева. Если пронумеровать правила грамматики, то последовательность номеров используемых правил также задает вывод.

Определение. Последовательность номеров правил грамматики Γ , применение которых позволяет построить вывод рассматриваемой цепочки σ из начального символа грамматики, называется **синтаксическим разбором** σ .

Например, в следующей грамматике

$\Gamma_{1.9}$:

$V_T = \{ i, +, *, (,) \}, V_A = \{ \langle E \rangle, \langle T \rangle, \langle P \rangle \}$

$R = \{ (1) \langle E \rangle \rightarrow \langle E \rangle + \langle T \rangle,$

(2) $\langle E \rangle \rightarrow \langle T \rangle,$

(3) $\langle T \rangle \rightarrow \langle T \rangle * \langle P \rangle,$

(4) $\langle T \rangle \rightarrow \langle P \rangle,$

(5) $\langle P \rangle \rightarrow (E),$

(6) $\langle P \rangle \rightarrow i \},$

правила которой пронумерованы, вывод

$$\begin{aligned} \langle E \rangle &\Rightarrow \langle E \rangle + \langle T \rangle \Rightarrow \langle T \rangle + \langle T \rangle \Rightarrow \langle T \rangle * \langle P \rangle + \langle T \rangle \Rightarrow \\ \langle P \rangle * \langle P \rangle + \langle T \rangle &\Rightarrow i * \langle P \rangle + \langle T \rangle \Rightarrow i * i + \langle T \rangle \Rightarrow \\ i * i + \langle P \rangle &\Rightarrow i * i + i \end{aligned}$$

имеет синтаксический разбор [1,2,3,4,6,6,4,6].

Если в процессе построения вывода появляются промежуточные цепочки, содержащие несколько нетерминальных символов, то можно продолжать вывод, заменяя любой из них. Таким образом, одни и те же правила могут быть использованы при выводе цепочки в разном порядке.

Например, вывод цепочки $i + i$ в грамматике $\Gamma_{1.9}$ может быть получен десятью различными способами.

Левый и правый выводы

Среди всевозможных выводов наибольший интерес представляют следующие два типа выводов.

Определение. Если при построении вывода цепочки α при каждом применении правила заменяется самый левый нетерминальный символ, то такой вывод называется левым или левосторонним выводом α . Если при построении вывода α , всегда заменяется самый правый нетерминальный символ промежуточной цепочки, то вывод называется правым или правосторонним выводом α .

Например, приведенный выше вывод цепочки $i * i + i$ в грамматике $\Gamma_{1.9}$ является левосторонним выводом. Следует отметить, что различным выводам цепочки $i+i$ в грамматике $\Gamma_{1.9}$ соответствует одно и то же синтаксическое дерево. Аналогичная ситуация имеет место и при выводе цепочки $i * i + i$.

Неоднозначные и эквивалентные грамматики

Существуют грамматики, в которых одна и та же цепочка может быть получена с помощью различных выводов. Например, в грамматике $\Gamma_{1.10}$ цепочка abc может быть получена с помощью двух различных выводов, и ей соответствуют два различных синтаксических дерева.

$\Gamma_{1.10}$

$V_T = \{a, b, c, d\}, V_A = \{<I>, <A>, \},$

$R = \{ <I> \rightarrow <A>,$

$<A> \rightarrow ac,$

$ \rightarrow b,$

$ \rightarrow cb\}.$

Первый вывод этой цепочки имеет вид:

1) $<I> \Rightarrow <A> \Rightarrow <A>b \Rightarrow acb,$

а второй можно получить так:

2) $<I> \Rightarrow <A> \Rightarrow <A>cb \Rightarrow acb$

Этим выводам соответствуют разные синтаксические деревья и разборы:

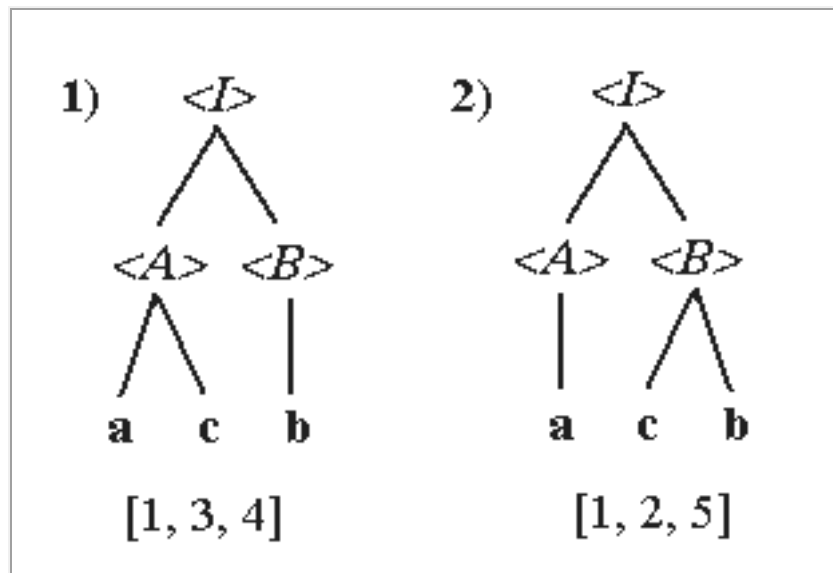


Рис. 1.

Следующая грамматика также допускает построение одной и той же цепочки с помощью двух выводов, имеющих разные синтаксические деревья.

$\Gamma_{1.11}$:

$V_T = \{0, +\}, V_A = \{<I>\},$

$R = \{ <I> \rightarrow 0,$

$<I> \rightarrow <I> + 0,$

$<I> \rightarrow 0 + <I> \}$

Два вывода этой грамматики, порождающие одинаковые цепочки, имеют вид:

1) $<I> \Rightarrow <I> + 0 \Rightarrow <I> + 0 + 0 \Rightarrow 0 + 0 + 0,$

2) $<I> \Rightarrow 0 + <I> \Rightarrow 0 + 0 + <I> \Rightarrow 0 + 0 + 0,$

а синтаксические деревья, соответствующие этим выводам, можно изобразить так:

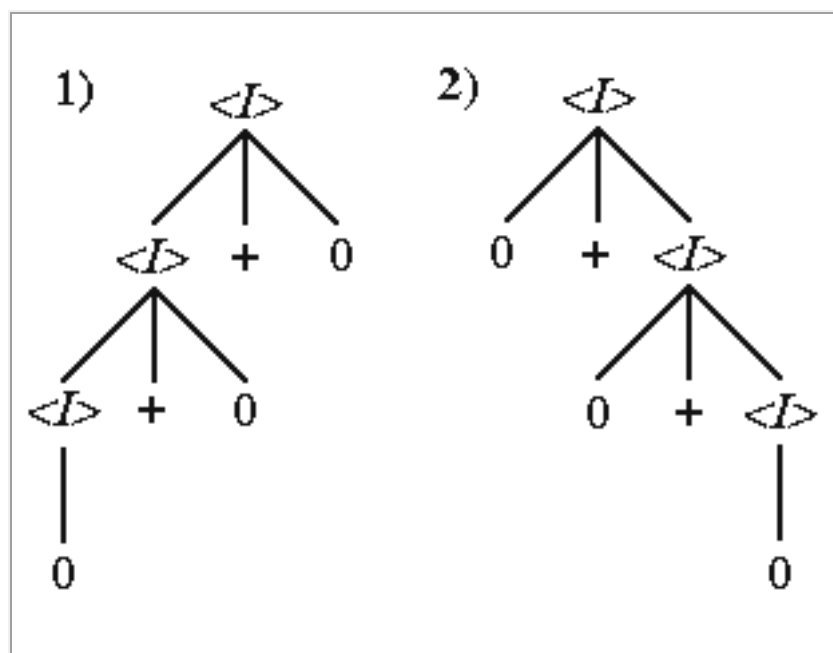


Рис. 2.

Рассмотренное свойство грамматик называется неоднозначностью. Оно может быть определено следующим образом.

Определение. Цепочка языка $L(\Gamma)$ называется неоднозначной, если для её вывода существует более чем одно синтаксическое дерево. Если грамматика Γ порождает неоднозначную цепочку, то она называется неоднозначной.

Свойство неоднозначности является крайне нежелательным для искусственных языков, поскольку оно не позволяет однозначным образом восстановить дерево вывода по заданной цепочке языка.

В общем случае можно сделать следующий вывод:

1. каждой цепочке, выводимой в грамматике, может соответствовать одно или несколько синтаксических деревьев,
2. каждому синтаксическому дереву могут соответствовать несколько выводов,
3. каждому синтаксическому дереву соответствует единственный правый и единственный левый выводы.

Кроме того, следует подчеркнуть, что один и тот же язык может быть получен с помощью различных грамматик

Определение. Две грамматики Γ_1 и Γ_2 называются эквивалентными, если они порождают один и тот же язык, т.е.

$$L(\Gamma_1) = L(\Gamma_2)$$

Способы задания схем грамматик.

Форма Наура-Бэкуса

Схема грамматики содержит правила вывода, определяющие синтаксис языка, или, другими словами, возможные компоненты и конструкции цепочек порождаемого языка. Для задания правил используются различные формы описания: символическая, форма Наура-Бэкуса, итерационная форма и синтаксические диаграммы.

В работах, связанных с рассмотрением общих свойств грамматик, обычно применяют символическую форму задания правил. Эта форма была рассмотрена в предыдущем параграфе. Она предусматривает использование в качестве элементов нетерминального словаря отдельных символов и стрелки в качестве разделителя правой и левой частей правила.

При описании синтаксиса конкретных языков программирования приходится вводить большое число нетерминальных символов, и символическая форма записи теряет свою наглядность. В этом случае применяют форму Наура-Бэкуса (ФНБ), которая предполагает использование в качестве нетерминальных символов комбинаций слов естественного языка, заключенных в угловые скобки, а в качестве разделителя - специального знака, состоящего из двух двоеточий и равенства. Например, если правила

$\langle L \rangle \rightarrow \langle L \rangle$ и $\langle L \rangle \rightarrow \langle E \rangle$ записаны в символической форме, и символ $\langle L \rangle$ соответствует синтаксическому понятию "список", а символ $\langle E \rangle$ - "элемент списка", то их можно представить в форме Наура-Бэкуса так:

$\langle \text{список} \rangle ::= \langle \text{элемент списка} \rangle \langle \text{список} \rangle,$

$\langle \text{список} \rangle ::= \langle \text{элемент списка} \rangle$

Чтобы сократить описание схемы грамматики, в ФБН разрешается объединять правила с одинаковой левой частью в одно правило, правая часть которого должна включать правые части объединяемых правил, разделенные вертикальной чертой. Используя объединение правил, для рассматриваемого примера получаем:

$\langle \text{список} \rangle ::= \langle \text{элемент списка} \rangle \langle \text{список} \rangle | \langle \text{элемент списка} \rangle$

Итерационная форма

Для получения более компактных описаний синтаксиса применяют итерационную форму описания. Такая форма предполагает введение специальной операции, которая называется итерацией и обозначается парой фигурных скобок со звездочкой. Итерация вида $\{a\}^*$ определяется как множество, включающее цепочки всевозможной длины, построенные с использованием символа a , и пустую цепочку.

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

Используя итерацию для описания множества цепочек, задаваемых символическими правилами, для списка получаем:

$\langle L \rangle \rightarrow \langle E \rangle \{ \langle E \rangle \}^*$

Например, описание множества цепочек, каждая из которых должна начинаться знаком $\#$ и может состоять из произвольного числа букв x и y , может быть представлено в итерационной форме так:

$\langle I \rangle \rightarrow \# \{x | y\}^*$

В итерационных формах описания наряду с итерационными скобками часто применяют квадратные скобки для указания того, что цепочка, заключенная в них, может быть опущена. С помощью таких скобок правила:

$\langle A \rangle \rightarrow x \langle A \rangle y \langle B \rangle z$ и $\langle A \rangle \rightarrow x \langle B \rangle z$

могут быть записаны так:

$\langle A \rangle \rightarrow x [\langle A \rangle y] \langle B \rangle z.$

Синтаксические диаграммы

Для того, чтобы улучшить зрительное восприятие и облегчить понимание сложных синтаксических описаний, применяют представление правил грамматики в виде синтаксических диаграмм. Правила построения таких диаграмм можно сформулировать в следующем виде:

1. Каждому правилу вида $\langle A \rangle \rightarrow a_1 \mid a_2 \mid \dots \mid a_k$ ставится в соответствие диаграмма, структура которой определяется правой частью правила.
2. Каждое появление терминального символа x в цепочке a_i изображается на диаграмме дугой, помеченной этим символом x , заключенным в кружок.

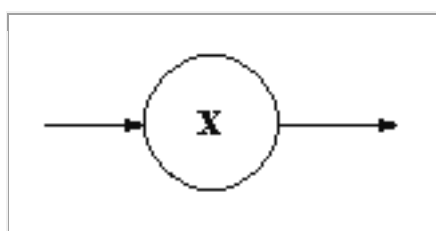


Рис. 1.

3. Каждому появлению нетерминального символа $\langle A \rangle$ в цепочке a_i ставится в соответствие на диаграмме дуга, помеченная символом, заключённым в квадрат.

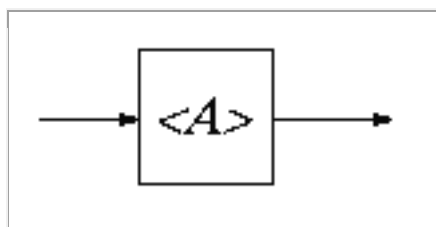


Рис. 2.

4. Порождающее правило, имеющее вид:

$$\langle A \rangle \rightarrow a_1 a_2 \dots a_n$$

изображается на диаграмме следующим образом:

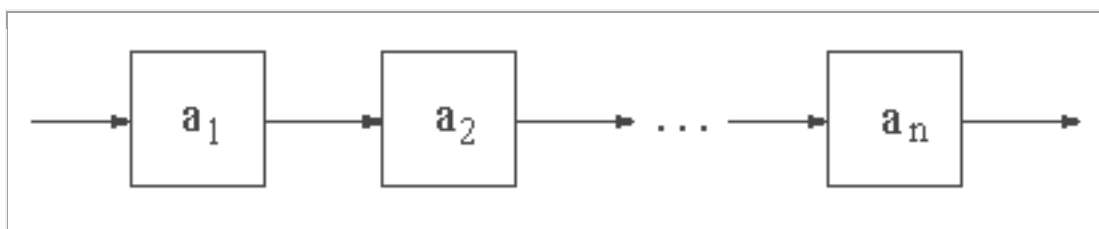


Рис. 3.

5. Порождающее правило, имеющее вид:

$$\langle A \rangle \rightarrow a_1 \mid a_2 \mid \dots \mid a_n$$

изображается на диаграмме так:

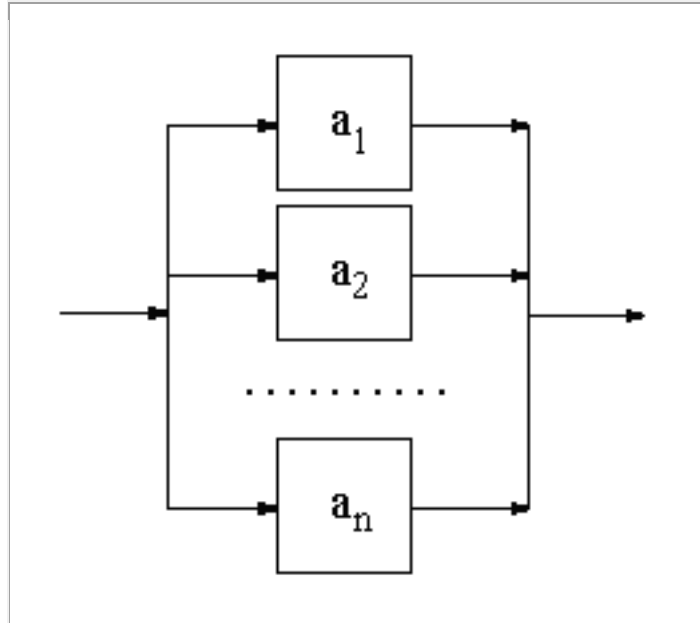


Рис. 4.

6. Если порождающее правило задано в виде итерации:

$$\langle A \rangle \rightarrow \{a\}^*,$$

то ему соответствует диаграмма:

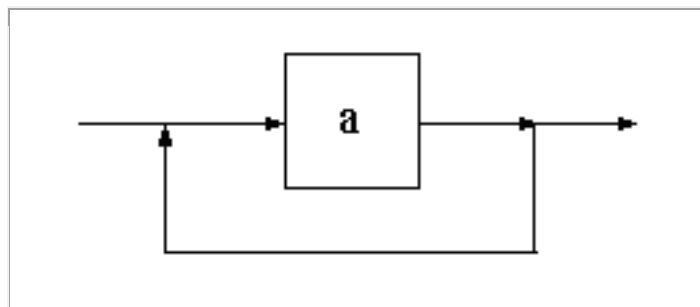


Рис. 5.

Число синтаксических диаграмм, которые можно построить для заданной схемы грамматики, определяется числом правил. Чтобы сократить число диаграмм, их объединяют, заменяя нетерминальные символы, входящие в диаграмму, построенными для них диаграммами.

Правила 3-6 предусматривают, что в качестве цепочки a_1 на объединенной диаграмме могут быть использованы диаграммы построенные для этих цепочек. В качестве примера рассмотрим следующую грамматику с начальным символом $\langle A \rangle$:

$\Gamma_{1.14}$:

$$V_T = \{x, +, (,)\}, V_A = \{\langle A \rangle, \langle B \rangle, \langle C \rangle\},$$

$$R = \{\langle A \rangle \rightarrow x \mid (\langle B \rangle),$$

$$\langle B \rangle \rightarrow \langle A \rangle \langle C \rangle,$$

$$\langle C \rangle \rightarrow \{+\langle A \rangle\}^*\}$$

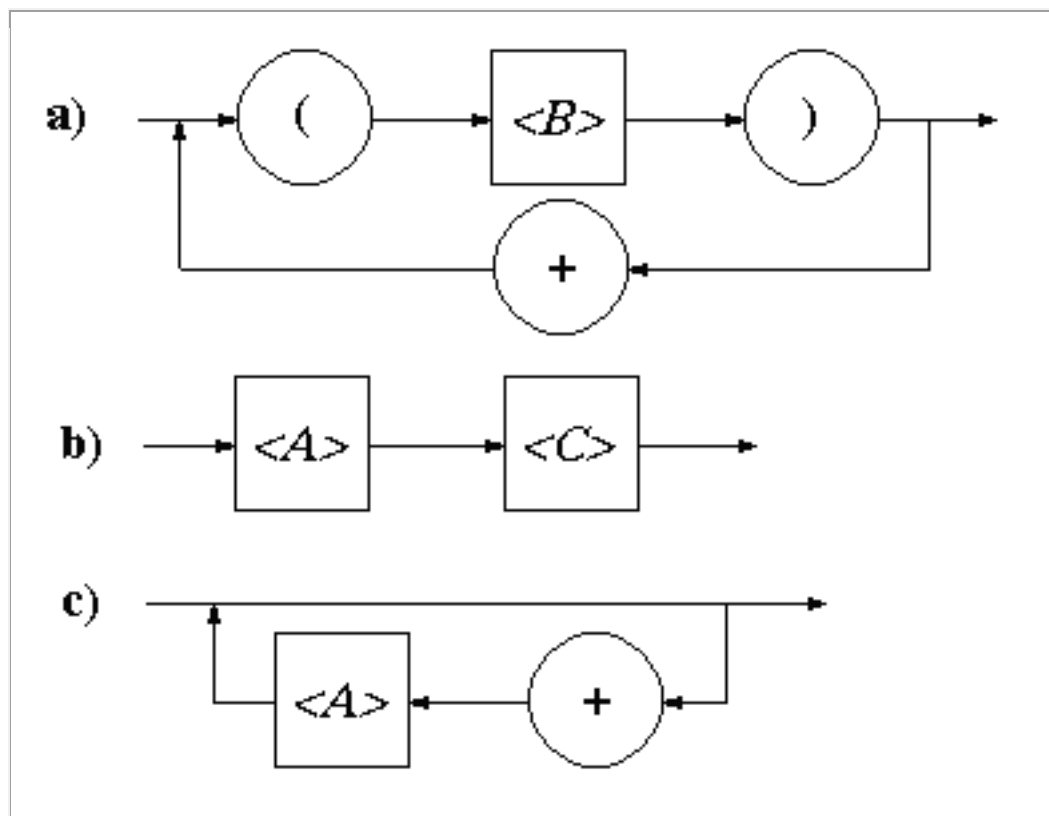


Рис. 6.

Заменяя нетерминальные символы, соответствующими диаграммами, получаем объединенную диаграмму в виде:

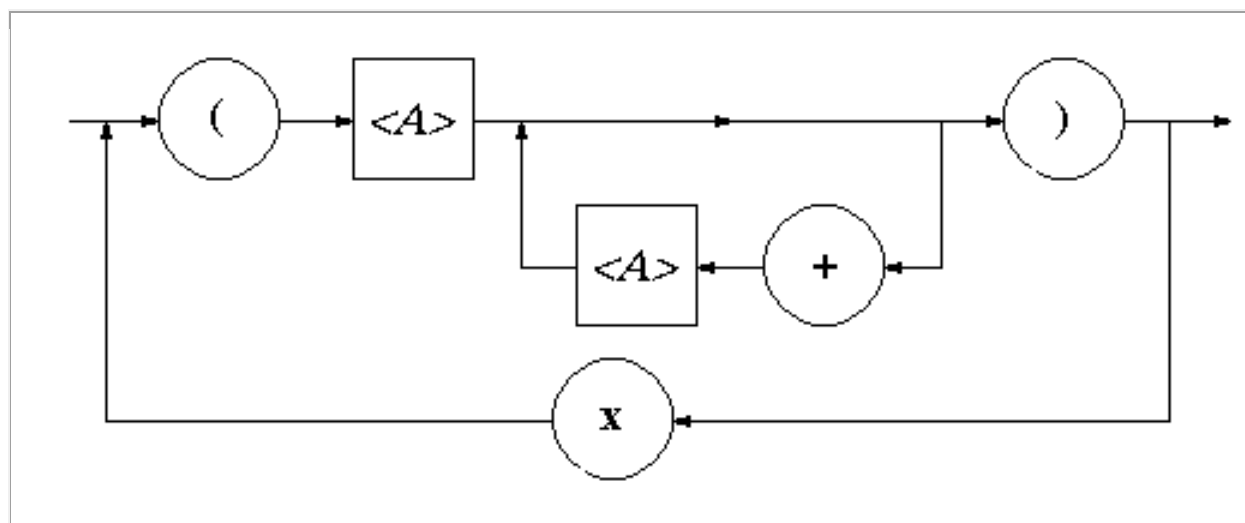


Рис. 7.