

ЛЕКЦИЯ 2

“ЭКВИВАЛЕНТНЫЕ ПРЕОБРАЗОВАНИЯ ГРАММАТИК. СИНТАКСИЧЕСКИЕ ДЕРЕВЬЯ И ИХ СВОЙСТВА.”

ПЛАН

1. Эквивалентные преобразования грамматик.
2. Получение ε -свободной грамматики.
3. Преобразование грамматик в нормальную форму Хомского.
4. Преобразование грамматик в нормальную форму Грейбах.
5. Синтаксические деревья и их свойства.
6. Неоднозначность грамматик: установление неоднозначности и коррекция её при помощи синтаксических деревьев.
7. Двоичные деревья и правила их построения.
8. Скелеты деревьев трансляции.

ЭКВИВАЛЕНТНЫЕ ГРАММАТИКИ

Пример 1.

$G_1[S]$

$S \rightarrow A \mid B$

$A \rightarrow a A \mid a$

$B \rightarrow b B \mid b$

$G_2[S]$

$S \rightarrow a A \mid b B \mid b \mid a$

$A \rightarrow a A \mid a$

$B \rightarrow b B \mid b$

Язык: $L = \{a^n, b^n; n \geq 1\}$.

Пример 2.

$G_1[S]$

$S \rightarrow a A \mid b B$

$A \rightarrow a$

$B \rightarrow b$

$G_2[S]$

$S \rightarrow A b$

$A \rightarrow a$

Язык $L = \{ab\}$.

ПОЛУЧЕНИЕ ε -СВОБОДНОЙ ГРАММАТИКИ

Правило вида $\mathbf{B} \rightarrow \varepsilon$ называется *ε -продукцией*. Грамматика, *не имеющая* ε -продукций, называется *ε -свободной* (e-free).

Теорема. Для любого контекстно-свободного языка L существует ε -свободная КС-грамматика $G[S]$, такая, что $L\{G[S]\} = L \setminus \{\varepsilon\}$.

Цель преобразования формально записывается так

$$S \xRightarrow[G]{*} x \quad \Leftrightarrow \quad S \xRightarrow[G']{*} x, \quad \text{для} \quad \forall x \in V_T^+, \quad (1)$$

$$x \in L\{G[S]\} \setminus \{\varepsilon\} \Leftrightarrow x \in L\{G'[S]\}. \quad (2)$$

АЛГОРИТМ ПРЕОБРАЗОВАНИЯ

Пусть $G[S] = (V_T, V_N, S, R)$ — КС-грамматика, имеющая ε -продукции.

1. Объединим *все* ε -продукции (непосредственного вывода) из исходного множества правил R в множество R_ε .

2. Отыщем все нетерминальные символы $N_i \in V_N$, такие что $N_i \Rightarrow^+ \varepsilon$ или $N_i \Rightarrow^* \varepsilon$. Указанные символы будем называть *ε -порождающими*.

3. Каждой продукции $p \in R$, в *правой* части которой находятся ε -порождающие символы, поставим в соответствие продукцию, в правой части которой, по сравнению с исходной, опущены один или более ε -порождающих символов. Получим множество R_0 .

4. Конструируем новое множество продукций грамматики по правилу

$$R' = \{R \setminus R_\varepsilon\} \cup \{R_0\}.$$

Пример.

$$\begin{aligned} S &\rightarrow [E] | E, \\ E &\rightarrow D | D+E | D-E, \\ D &\rightarrow F | D * F | D / F, \\ F &\rightarrow a | b | c | \varepsilon. \end{aligned}$$

1. Множество $R_\varepsilon = \{F \rightarrow \varepsilon\}$.

2. Поиск ε -порождающих символов $F \Rightarrow \varepsilon$, $S \Rightarrow^* \varepsilon$, $E \Rightarrow^+ \varepsilon$, $D \Rightarrow^+ \varepsilon$.

3 Построение множества R_0 .

$$\begin{aligned} S &\rightarrow [], \\ E &\rightarrow D+|+E|D-|-E|+|- , \\ D &\rightarrow *F|D*|D/|/F|/|* . \end{aligned}$$

4. Окончательный набор продукций для $G'[S]$.

$$\begin{aligned} S &\rightarrow [E] | E|[], \\ E &\rightarrow D|D+E|D-E| D+|+E|D-|-E|+|- , \\ D &\rightarrow F|D * F|D / F| * F|D *|D /|/F|/|* , \\ F &\rightarrow a|b|c. \end{aligned}$$

ПРЕОБРАЗОВАНИЕ ГРАММАТИКИ В НОРМАЛЬНУЮ ФОРМУ ХОМСКОГО (CHOMSKY NORMAL FORM)

Продукция КС-грамматики $G[S] = (V_T, V_N, S, R)$ называется **релевантной** тогда и только тогда, когда существует вывод некоторой цепочки $x \in L\{G[S]\}$, в которой эта продукция используется.

То есть, если $A ::= \alpha$, $\alpha \in V^*$, $A \in V_N$ – релевантная продукция, то

$$\forall x \in \{L(G[S]) \mid S \Rightarrow^* \beta A \gamma \Rightarrow \beta \alpha \gamma \Rightarrow^+ x\}.$$

Теорема. Любой ε -свободный КС-язык может быть порождён некоторой КС-грамматикой в нормальной форме Хомского, продукции которой имеют вид: $A \rightarrow BC$ или $A \rightarrow \alpha$, где $A, B, C \in V_N$, $\alpha \in V_T$.

АЛГОРИТМ ПРЕОБРАЗОВАНИЯ

1. Среди продукций грамматики отыщем все **первичные продукции** вида

$$A_i \rightarrow B_j, A_i, B_j \in V_N.$$

правая часть которых содержит единственный нетерминальный символ.

При этом всё множество продукций с нетерминалом A_i в левой части разбивается на подмножества: $U\{A_i\}$ – первичных продукций и $N\{A_i\}$ – непервичных продукций.

Для $\forall A_i$, где $U\{A_i\} \neq \{\}$, поставим в соответствие множество правил

$$\{A_i \Rightarrow^+ \alpha \mid A_i \rightarrow B_j, B_j \Rightarrow^+ \alpha, \alpha \in V^+, B_j \in N(B_j)\}.$$

В результате не останется первичных продукций.

2. Отыщем все **вторичные** продукции. Правые части таких продукций представляют собой **смешанные цепочки** терминальных и нетерминальных символов.

Каждый терминальный символ такой цепочки η заменяется нетерминалом A_η . Одновременно множество продукций пополняется правилом $A_\eta \rightarrow \eta$, определяющим данный нетерминал.

Для примера, пусть правая часть вторичной продукции $x_1x_2...x_m, x_i \in (V_T \cup V_N), i=1, \overline{m}$. Тогда после преобразования она примет вид $y_1y_2...y_m, y_i \in V_N, i=1, \overline{m}$ при соответствии между цепочками

$$y_i = \begin{cases} x_i, & x_i \in V_N, \\ X_{x_i}, & X_{x_i} \in V_N, x_i \in V_T, \\ R^* = R \cup \{X_{x_i} ::= x_i\} \end{cases}$$

Будут устранены все вторичные продукции.

3. В модифицированном, в ходе п.п. 2 и 3 настоящего алгоритма, множестве продукций, отыщем **третичные** продукции, в правых частях которых содержится **не менее трёх** нетерминалов подряд.

Для каждой найденной третичной продукции строится разложение. Если исходная третичная продукция имеет вид:

$$A_r \rightarrow S_1S_2...S_m, S_i \in V_N, i=1, \overline{m}, m \geq 2.$$

Её разложение будет выглядеть так:

$$\begin{aligned} A_r &\rightarrow S_1 N_{B1}, \\ N_{B1} &\rightarrow S_2 N_{B2}, \\ &\dots \\ N_{B_{m-1}} &\rightarrow S_{m-1} S_m. \end{aligned}$$

Нетерминалы N_{Bi} , $i=1, m-1$ суть новые, нетерминалы, *не содержащиеся более ни в одной продукции*, добавленные к словарию V_N в ходе разложения.

Пример.

$$\begin{aligned} S &\rightarrow ABA|A, \\ A &\rightarrow a|aA|B, \\ B &\rightarrow b|bB. \end{aligned}$$

1. Первичные продукции суть

$$\begin{aligned} S &\rightarrow A, \\ A &\rightarrow B. \end{aligned}$$

Имеем следующие множества первичных и непервичных продукций.

$$\begin{aligned} U(A) &= \{A \rightarrow B\}, & N(A) &= \{A \rightarrow aA|a\}, \\ U(B) &= \{\}, & N(B) &= \{B \rightarrow bB|b\}, \\ U(S) &= \{S \rightarrow A\}. & N(S) &= \{S \rightarrow ABA\}. \end{aligned}$$

Правилу $A \rightarrow B$ поставим в соответствие заменяющие правила $A \rightarrow bB|b$, а продукции $S \rightarrow A$ – правила $S \rightarrow aA|bB|a|b$.

Получим новый набор правил изменённой грамматики.

$$S \rightarrow ABA | aA | bB | a | b,$$

$$A \rightarrow aA | bB | a | b,$$

$$B \rightarrow b | bB.$$

2 Ко вторичным продукциям относятся все продукции с правыми частями aA и bB . Введением правил: $C \rightarrow a$, $D \rightarrow b$ эти правые части преобразуются к виду **CA** и **DB**. В результате:

$$S \rightarrow ABA | CA | DB | a | b,$$

$$A \rightarrow CA | DB | a | b,$$

$$B \rightarrow DB | b,$$

$$C \rightarrow a,$$

$$D \rightarrow b.$$

3. Единственная третичная продукция: $S \rightarrow ABA$. Заменим её продукциями $S \rightarrow AE$, $E \rightarrow BA$. Получим грамматику в нормальной форме Хомского:

$$S \rightarrow AE | CA | DB | a | b,$$

$$A \rightarrow CA | DB | a | b,$$

$$B \rightarrow DB | b,$$

$$E \rightarrow BA,$$

$$C \rightarrow a,$$

$$D \rightarrow b.$$

ИССЛЕДОВАНИЕ НА РЕЛЕВАНТНОСТЬ

Свойство А. Если все нетерминальные символы правой части правила продуктивны, то продуктивен и символ, стоящий в левой части правила.

Алгоритм поиска непродуктивных нетерминалов

1. Составить список нетерминалов, для которых найдётся хотя бы одна продукция, правая часть которого состоит только из терминалов.
2. Если найдено такое правило, что все нетерминальные символы, стоящие в его правой части, уже помещены в список, то добавить в список нетерминал, стоящий в левой части данного правила.
3. Если список не пополняется, то эти нетерминальные символы продуктивны, остальные элементы V_N – не продуктивны.

Свойство Б. Если все нетерминальный символ в левой части правила является достижимым, то достижимы и все символы, стоящие в правой части этого правила.

Алгоритм поиска недостижимых нетерминалов

1. Составить множество, первоначально, содержащее только аксиому грамматики.
2. Если найдено правило, нетерминал из левой части которого уже принадлежит множеству, то поместить во множество нетерминалы из правой части правила, которые в нём отсутствуют.
3. Если множество не пополняется, то нетерминалы, его составляющие, являются достижимыми, а прочие нетерминалы из V_N – нет.

ПРЕОБРАЗОВАНИЕ ГРАММАТИКИ В НОРМАЛЬНУЮ ФОРМУ ГРЕЙБАХ

КС-грамматика, все productions которой имеют вид

$$A \rightarrow a\beta, \quad a \in V_T, \beta \in (V_N \cup V_T)^*$$

называется представленной в форме Грейбах (Greibach S.A. normal form).

Теорема. Для каждого ε -свободного КС-языка существует КС-грамматика в нормальной форме Грейбах, такая, что $L\{G_0[S]\} = L\{G_G[S']\}$.

Теорема 1. Если $A \rightarrow \alpha B \gamma$, $A, B \in V_N$, $\alpha, \gamma \in V^*$ правило КС – грамматики, а $B \rightarrow \beta_1 | \beta_2 | \dots | \beta_m$ – все productions данной грамматики с нетерминалом B в своих левых частях, то указанные productions можно заменить правилами вида $A \rightarrow \alpha \beta_1 \gamma | \alpha \beta_2 \gamma | \dots | \alpha \beta_m \gamma$, причём произведённая замена никак не отражается на языке, порождаемом этой грамматикой.

Теорема 2. Если productions некоторой контекстно-свободной грамматики являются рекурсивными слева productions $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m$, $A \in V_N$, $\alpha_i \in (V \setminus A)^*$, а остальные productions с нетерминалом A в левой части суть $A \rightarrow \beta_1 | \beta_2 | \dots | \beta_k$, то новая грамматика, эквивалентная исходной, может быть построена путём добавления нетеминала A^∇ к множеству V_N и заменой этих productions исходной грамматики productions вида

$$\begin{aligned} &A \rightarrow \beta_1 | \beta_2 | \dots | \beta_k | \beta_1 A^\nabla | \beta_2 A^\nabla | \dots | \beta_k A^\nabla \text{ и} \\ &A^\nabla \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_m | \beta_1 A^\nabla | \beta_2 A^\nabla | \dots | \beta_k A^\nabla. \end{aligned}$$

Синтаксические деревья и графы

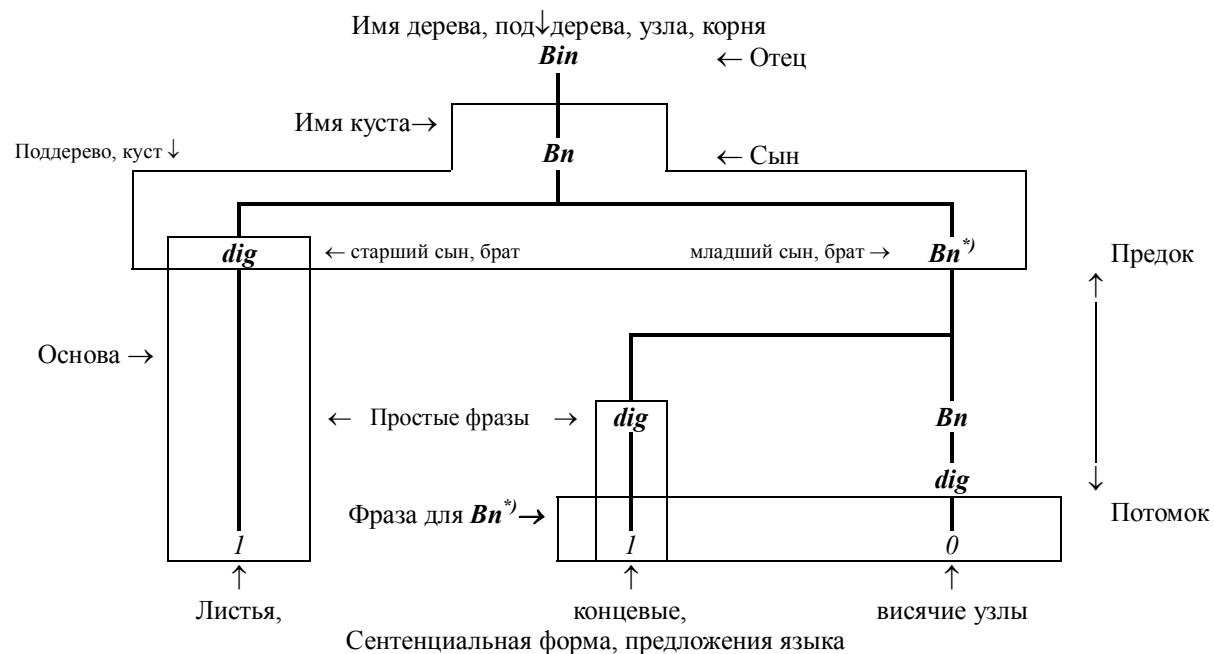


Рисунок 1 – Синтаксическое дерево разбора предложения 110 грамматики $G[Bin]$ (иллюстрация терминологии)

СВОЙСТВА СИНТАКСИЧЕСКИХ ДЕРЕВЬЕВ

1. Для каждого синтаксического дерева существует, по крайней мере, один вывод.

2. Для каждого вывода существует *соответствующее* синтаксическое дерево.

3. Куст дерева указывает непосредственный вывод, в котором имя куста заменяется узлами куста. То есть, среди продукций грамматики существует *правило*, *левая* часть которого – *имя куста*, а *правая* – цепочка из узлов куста, например $\langle Bn \rangle \rightarrow \langle dig \rangle \langle Bn \rangle$.

4. Концевые узлы дерева образуют sentential form на каждом шаге вывода или приведения $\langle Bin \rangle \Rightarrow^* 110$, $\langle Bin \rangle \Rightarrow^* \langle dig \rangle 10$ и т.д

5. Если U – корень поддеревы для sentential form $w=xu$, где u образует цепочку концевых узлов этого поддерева, то u – фраза sentential form w для U . Если поддерево представлено единственным кустом, то это простая фраза, например

$w = 110$, $u \sim 10$, $U \sim \langle Bn \rangle$ – фраза, $u \sim 1$, $U \sim \langle dig \rangle$ – простая фраза.

6. Синтаксическое дерево не отражает порядок использования правил. Для синтаксического анализа это несущественно, главное, что дерево построено с их использованием.

ИССЛЕДОВАНИЕ ГРАММАТИК НА ОДНОЗНАЧНОСТЬ

Предложение является *синтаксически неоднозначным*, если для его вывода существует *несколько* деревьев. Грамматика *неоднозначна*, если допускает *неоднозначные* предложения.

Примеры. 1. “They are flying planes”. 2. “Комиссия по изучению четвертичного периода АН СССР”.

3. Грамматика $G[V]$.

$$V ::= V + V \mid V - V,$$
$$V ::= V * V \mid V / V,$$
$$V ::= \langle iden \rangle \mid \langle data \rangle \mid (V)$$

Порождает сентенциальные формы:

$V \Rightarrow^* \langle iden \rangle + \langle data \rangle * \langle iden \rangle$ и $V \Rightarrow^* \langle iden \rangle * (\langle iden \rangle - \langle iden \rangle)$, которые соответствуют в программах на алгоритмических языках конструкции выражений, например

$x + 2.0 * y$ или $alpha * (beta - gamma)$.

Сентенциальной форме $\langle iden \rangle + \langle data \rangle * \langle iden \rangle$.

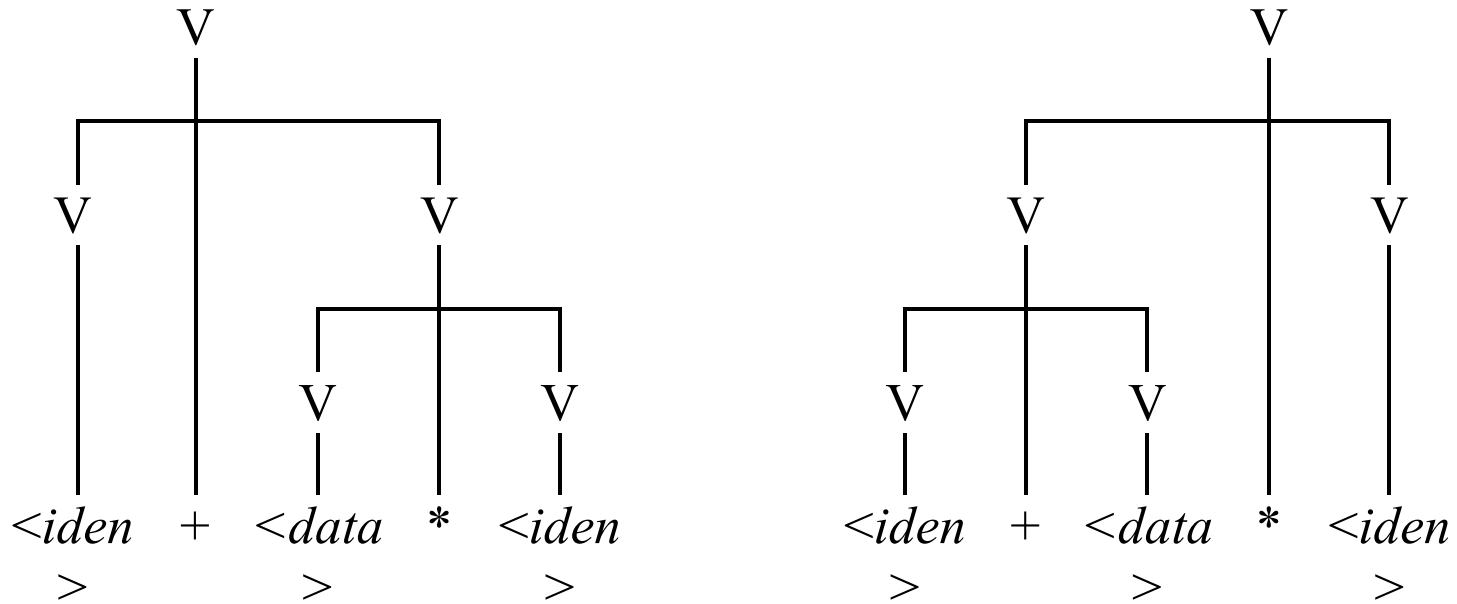


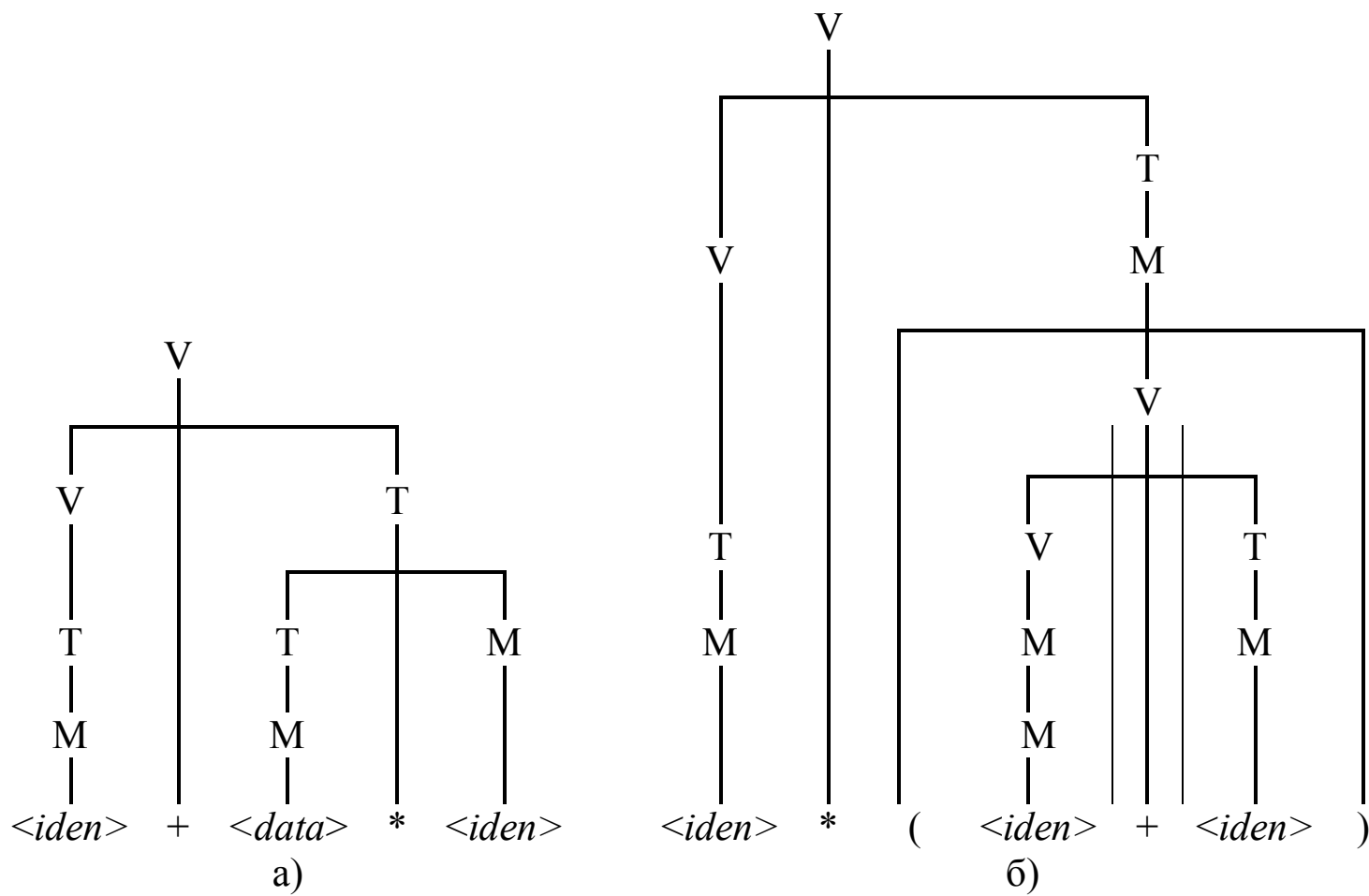
Рисунок 2 – Неоднозначные синтаксические деревья грамматики $G[V]$

ПРЕОБРАЗОВАННАЯ ГРАММАТИКА

$V ::= T \mid V + T \mid V - T,$

$T ::= M \mid V * M \mid V / M,$

$M ::= \langle iden \rangle \mid \langle data \rangle \mid (V).$



ПОСТРОЕНИЕ ДВОИЧНЫХ ДЕРЕВЬЕВ ПО ПРОИЗВОЛЬНОМУ ДЕРЕВУ

Правила преобразования

Пусть T – произвольное дерево, а $B(T)$ – соответствующее ему двоичное дерево.

1. Вершина v является левым потомком (страшим сыном) вершины a дерева $B(T)$ тогда и только тогда, когда v самый левый потомок a в T .

От отца дуга влево ведёт к старшему сыну. Рисунок 3, а.

2. Вершина v является правым потомком (младшим сыном) вершины a дерева $B(T)$ тогда и только тогда, когда v является ближайшем из правых братьев a дерева T .

Правая дуга ведёт от старших братьев к младшим. Рисунок 3, б.

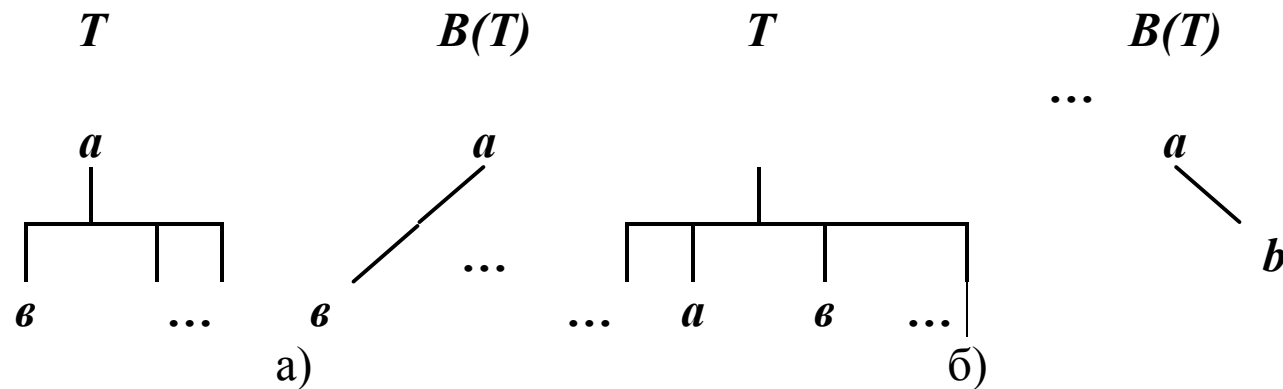


Рисунок 3 – Иллюстрация правил преобразования

СВОЙСТВА ДВОИЧНЫХ ДЕРЕВЬЕВ

Левосторонняя (Правосторонняя) цепочка двоичного дерева есть последовательность *вершин*

$a_1, a_2, a_3, \dots, a_n$, ($n > 1$), где a_i является левым (правым) потомком a_{i-1} .

Цепочка является *максимальной*, если содержится сама в себе. Все остальные её подцепочки имеют *правильные* головы и хвосты.

Для произвольных деревьев левосторонние и правосторонние цепочки определяются аналогично.

Лемма 1. Левосторонние цепочки в $B(T)$ являются левосторонними и в T и наоборот.

Лемма 2. Максимальная правосторонняя цепочка узлов в $B(T)$ является *аргументом* некоторого правила подстановки P . Предок левой вершины в цепочке есть корень P .

Лемма 3. Максимальная левосторонняя цепочка узлов в $B(T)$ содержит только один терминальный символ, являющийся концевым (последней вершиной) в этой цепочке.

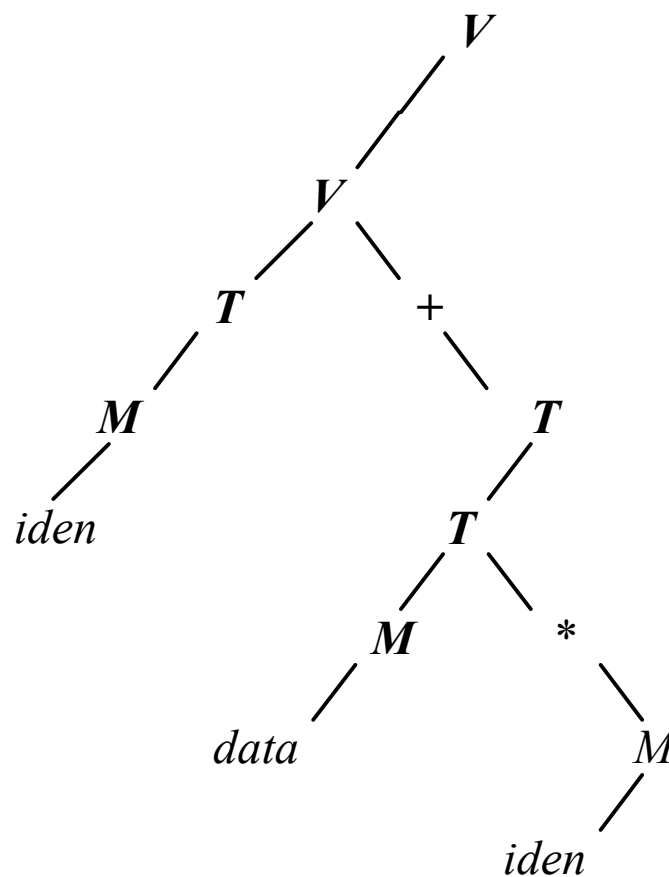


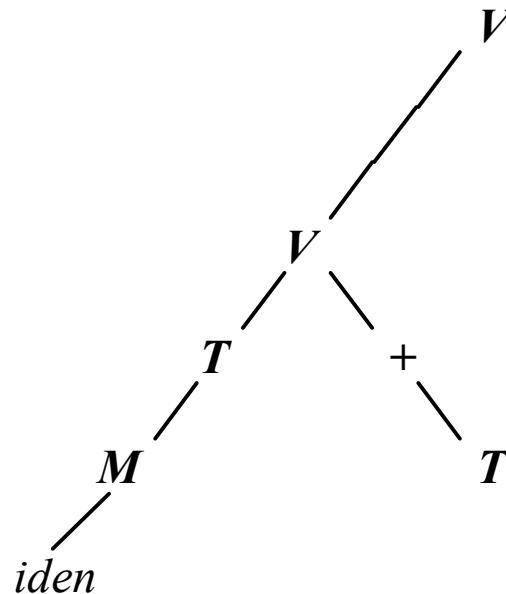
Рисунок 4 – Двоичное дерево синтаксической формы $\langle iden \rangle + \langle data \rangle * \langle iden \rangle$

СКЕЛЕТЫ ДЕРЕВЬЕВ ТРАНСЛЯЦИИ

Пусть имеется двоичное синтаксическое дерево D : начиная с аксиомы S , заканчивая предложением языка $x_1 x_2 \dots x_n$, как это показано на рисунке 4.

Определение 1. Максимальная левосторонняя цепочка от S к x_n называется *позвоночником* этого дерева.

Определение 2. Позвоночник, вместе с максимальными правосторонними цепочками, выходящими из него, называется скелетом дерева подстановок.



РЕЗЮМЕ

Тактика трансляции такова.

Для грамматики $G[S]$, по заданной входной цепочке x , транслятор находит, если это возможно, некоторый вывод цепочки x из аксиомы S

$$S \Rightarrow \alpha_0 \Rightarrow \alpha_1 \Rightarrow \alpha_2 \dots \Rightarrow \alpha_n \equiv x$$

и строит вывод

$$S \Rightarrow (\alpha_0, \beta_0) \Rightarrow (\alpha_1, \beta_1) \Rightarrow (\alpha_2, \beta_2) \dots \Rightarrow (\alpha_n, \beta_n) \equiv (x, y).$$

То есть, *каждая цепочка β_m порождает цепочку β_{m+1}* с помощью соответствующего правила, применяемого при переходе от α_m к α_{m+1} , а цепочка y , представляющая собой результат трансляции, служит выходом для цепочки x .

Пусть V_T — входной алфавит, а Δ — выходной алфавит.

Переводом с языка $L_1 \subseteq V_T^*$ на язык $L_2 \subseteq \Delta^*$ называется отображение T из V_T^* в Δ^* , для которого L_1 — область определения, а L_2 — множество значений.