

Эквивалентные преобразования контекстно-свободных грамматик

С.Ю.Соловьев

МГУ имени М.В.Ломоносова, факультет ВМК, Москва, Россия

Поступила в редколлегию 25.09.2010

Аннотация—В работе описывается и обосновывается эквивалентное преобразование контекстно-свободных грамматик, позволяющее удалять общие префиксы и суффиксы в терминальных реализациях нетерминальных символов. Кроме того, в работе предлагается универсальный метод совместного применения нескольких эквивалентных преобразований грамматик.

1. ВВЕДЕНИЕ

В теории формальных языков известно большое количество эквивалентных преобразований контекстно-свободных грамматик (КС-грамматик). Помимо прочего это означает, что один и тот же КС-язык может порождаться весьма непохожими грамматиками. В настоящей работе эквивалентные преобразования грамматик рассматриваются с точки зрения удаления конструкций “несущественных” для порождаемого языка.

Контекстно-свободной грамматикой [1] называется четверка $G = \langle N, \Sigma, P, S \rangle$, где N – алфавит нетерминальных символов (нетерминалов); Σ – непересекающийся с N алфавит терминальных символов (терминалов); P – конечное множество правил вывода вида $A \rightarrow \alpha$, где $A \in N$, α – цепочка символов из $N \cup \Sigma$; S – выделенный символ из N , именуемый начальным символом.

В последующих выкладках будем полагать, что действуют следующие соглашения:

- A, B, C, D – нетерминальные символы; S – начальный символ;
- a, b, c, d – терминальные символы;
- $\alpha, \beta, \gamma, \sigma$ – цепочки символов из $N \cup \Sigma$;
- x, y, z – цепочки символов (предложения) из Σ ;
- ϵ – пустая цепочка нулевой длины;
- запись $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$, означает множество правил $\{ A \rightarrow \alpha_1, \dots, A \rightarrow \alpha_n \}$;
- запись $\alpha \Rightarrow_G \beta$ означает, что $\alpha = \alpha_1 A \alpha_2$, $\beta = \alpha_1 \gamma \alpha_2$ и $A \rightarrow \gamma \in P$; в этом случае будем говорить, что цепочка β непосредственно выводима из цепочки α в грамматике G ;
- $L(G) = \{ x \mid S \Rightarrow_G^* x \}$ – язык, порождаемый грамматикой G .

Принятые соглашения позволяют, в частности, задавать КС-грамматики простым перечислением правил вывода.

2. КС#ГРАММАТИКИ

Каждая контекстно-свободная грамматика $G = \langle N, \Sigma, P, S \rangle$ порождает семейство грамматик $G(A) = \langle N, \Sigma, P, A \rangle$, где $A \in N$. Для начального символа S имеем: $G(S) = G$ и $L(G(S)) = L(G)$. В общем случае КС-язык $L(G(A))$ есть реализация нетерминала A в классе терминальных цепочек.

В дальнейшем изложении будем рассматривать только такие КС-грамматики $\langle N, \Sigma, P, S \rangle$, в которых:

- отсутствуют правила с пустой правой частью; и
- имеется правило $S \rightarrow \#$, причем терминальный символ $\#$ в других правилах не встречается.

Для КС-грамматик, удовлетворяющих перечисленным свойствам, будем использовать обозначение КС $\#$ грамматики.

С точки зрения порождаемых языков приведенные ограничения не являются существенными. Любой КС-язык L может быть получен из КС $\#$ языка $L_{\#}$ одним из двух способов: либо $L = L_{\#} \setminus \{\#\}$, либо $L = (L_{\#} \setminus \{\#\}) \cup \{e\}$; соответствующая КС-грамматика получается либо посредством удаления правила $S \rightarrow \#$, либо посредством его замены на правило $S \rightarrow e$. Вместе с тем, использование дополнительного символа $\#$ позволяет естественным образом вывести начальный символ S из-под действия многих эквивалентных преобразований.

Остановимся на одном из преобразований КС-грамматик, связанным с изменением языков $L(G(A))$, $A \neq S$. В качестве неформально введения в задачу рассмотрим цепочки $abcs$, $abcsfg$ и $abcdcdg$. Эти цепочки имеют общий префикс ab и общий суффикс g . В более изоциренных случаях необходимо точно оговорить, что считать общим префиксом и/или суффиксом:

- для $\{a, ab, abb\}$ префикс и суффикс отсутствуют;
- для $\{abc, abcs, abcss\}$ префикс = ab , суффикс отсутствует;
- для $\{abc, abbc, abbbc\}$ префикс = ab , суффикс отсутствует и пр.

В общем случае, после удаления общих префикса и суффикса не должна возникать пустая цепочка. Кроме того, будем исходить из того, что сначала определяется префикс, а затем – суффикс.

Если язык $L(G(A))$ имеет непустой общий префикс x ,

то будем говорить, что нетерминал A имеет префикс x .

Если язык $L(G(A))$ имеет непустой общий суффикс z ,

то будем говорить, что нетерминал A имеет суффикс z .

Продолжая неформальное введение, рассмотрим КС-грамматику

$$G : \begin{array}{l} S \rightarrow A + A \mid gb \\ A \rightarrow abcg \mid abcfg \mid abddg \end{array}$$

Здесь язык $L(G(A)) = \{abcs, abcsfg, abcdcdg\}$ задан в явном виде. У нетерминала A префикс ab и суффикс g можно изъять и передать “вышестоящему” нетерминалу S , при этом грамматика примет следующий вид

$$G' : \begin{array}{l} S \rightarrow abA'g + abA'g \mid gb \\ A' \rightarrow c \mid cf \mid dd \end{array}$$

В грамматике G' нетерминал A' не имеет ни префикса ни суффикса. Понятно, что грамматики G и G' эквивалентны. Возникает вопрос о возможности изъятия префиксов и суффиксов нетерминалов для КС-грамматик общего вида. Отметим, что в данном случае изъятие рассматривается как эквивалентное преобразование грамматик.

Нетерминалы, у которых можно забирать префиксы и суффиксы, удовлетворяют естественным условиям:

$$\begin{array}{ll} L(G(A)) = \text{конкатенация}(a, L) & (L) \\ \text{и/или } L(G(A)) = \text{конкатенация}(L, a), & (R) \end{array}$$

где L – некоторый язык, $e \notin L$.

Однако приведенные условия являются слишком общими; можно показать алгоритмическую неразрешимость задачи изъятия префиксов и суффиксов у всех нетерминалов, удовлетворяющих условиям (L) и (R) . Пойдем по пути уточнения класса нетерминалов, у которых можно “безнаказанно” забирать префиксы и суффиксы.

3. LD-ПРЕОБРАЗОВАНИЕ КС#ГРАММАТИК

Подмножество нетерминалов КС#грамматики $G = \langle N, \Sigma, P, S \rangle$ будем называть делимым-слева относительно терминала a , и будем его обозначать $LD(a)$,

если множество правил вывода P образуют правила трех типов:

- тип 1: $A \rightarrow a\alpha$, где $A \in LD(a)$ и $\alpha \neq \epsilon$;
 тип 2: $A \rightarrow A_0\alpha$, где $A \in LD(a)$ и $A_0 \in LD(a)$
 тип 3: $B \rightarrow \alpha$, где $B \notin LD(a)$

Понятно, что если $A \in LD(a)$, то $L(G(A))$ обязательно удовлетворяет условию (L) ; если $B \notin LD(a)$, то $L(G(B))$ может удовлетворять или не удовлетворять условию (L) . Например:

$$G_T : \begin{array}{lcl} S & \rightarrow & A + A \quad | \quad Bg \\ A & \rightarrow & ab \quad | \quad Aa \\ B & \rightarrow & Dd \quad | \quad abd \\ D & \rightarrow & a \quad | \quad aD \end{array}$$

$\forall x \in \{ a, b, d, g, + \}$ $D \notin LD(x)$, и, следовательно, $B \notin LD(x)$.

$LD(a) = \{ A \}$,

правила типа 1: $A \rightarrow ab$;

правила типа 2: $A \rightarrow Aa$;

правила типа 3: $S \rightarrow A+A$, $S \rightarrow Bg$, $B \rightarrow Dd$, $B \rightarrow abd$, $D \rightarrow a$, $D \rightarrow aD$.

$A \in LD(a)$; $L(G_T(A)) = \{ aba^n \mid n \geq 0 \} = \{ ax \mid x = ba^n, n \geq 0 \}$.

$B \notin LD(a)$; $L(G_T(B)) = \{ a^n d \mid n \geq 1 \} \cup \{ abd \} = \{ ax \mid x = a^n d / bd, n \geq 2 \}$.

Предложения из $L(G_T(D))$ не сводимы к виду ax , где $x \neq \epsilon$.

Конец примера.

Каждому делимому-слева подмножеству нетерминалов $LD(a) = \{ A_1, A_2, \dots \}$ поставим во взаимнооднозначное соответствие подмножество ранее не использовавшихся нетерминальных символов $LD'(a) = \{ A'_1, A'_2, \dots \}$.

Если в контексте некоторой КС-грамматики известно множество $LD(a)$, то можно рассматривать преобразование цепочек W , заключающееся в выполнении всевозможных подстановок aA'_i вместо A_i . Цепочка $W(\alpha)$ получается из цепочки α посредством замены всех символов A из $LD(a)$ на цепочки из двух символов aA' , где $A' \in LD'(a)$. Обратное преобразование W^{-1} заменяет в заданной цепочке все вхождения aA'_i на соответствующие нетерминалы A_i . Обратное преобразование полностью удаляет из заданной цепочки α все символы $LD'(a)$ только в том случае, когда непосредственно перед символом из $LD'(a)$ располагается символ a .

Например, в грамматике G_T для $LD(a) = \{ A \}$ имеем:

$$\begin{aligned} W(Aa + A + aBg) &= aA'a + aA' + aBg, \\ W^{-1}(aA'a + aA' + aBg) &= Aa + bA' + aBg, \\ W^{-1}(aA'a + aA') &= Aa + A \end{aligned}$$

Пусть $LD(a)$ – подмножество делимых-слева нетерминалов некоторой КС#грамматики $G = \langle N, \Sigma, P, S \rangle$, определим грамматику $G_W = \langle N_W, \Sigma, P_W, S \rangle$ следующим образом:
 $N_W = (LD'(a) \cup N) \setminus LD(a)$, а P_W построено так:

- если $A \rightarrow a\alpha$ есть правило типа 1 грамматики G , то $A' \rightarrow W(\alpha) \in P_W$;
- если $A \rightarrow A_0\alpha$ есть правило типа 2 грамматики G , то $A' \rightarrow A'_0W(\alpha) \in P_W$;
- если $B \rightarrow \alpha$ есть правило типа 3 грамматики G , то $B \rightarrow W(\alpha) \in P_W$.

Например, для множества $LD(a) = \{ A \}$ грамматики

$$G : \begin{array}{l} S \rightarrow A \mid A + A \mid \# \\ A \rightarrow ab \mid A + A \mid A - A \end{array}$$

имеем:

$$G_W : \begin{array}{l} S \rightarrow aA' \mid aA' + aA' \mid \# \\ A' \rightarrow b \mid A' + aA' \mid A' - aA' \end{array}$$

Заметим, что преобразование грамматики G в грамматику G_W (LD-преобразование) возможно только в том случае, когда в G найдется хотя бы одно непустое множество нетерминалов $LD(a)$, более того, $LD(a)$ является существенным аргументом LD-преобразования.

Утверждение 1. $L(G) \subseteq L(G_W)$ для КС#грамматики $G = \langle N, \Sigma, P, S \rangle$.

Доказательство. Рассмотрим произвольное предложение x из $L(G)$ и некоторый левый вывод x в грамматике G .

$$S = \sigma_0 \Rightarrow_G \sigma_1 \Rightarrow_G \dots \Rightarrow_G \sigma_{k-1} \Rightarrow_G \sigma_k = x.$$

Докажем по индукции, что последовательность цепочек $W(\sigma_0), \dots, W(\sigma_k)$ есть левый вывод предложения x в грамматике G_W .

Из-за наличия правила $S \rightarrow \#$ основной символ S не может входить ни в одно множество делимых-слева нетерминалов. Поэтому $W(\sigma_0) = W(S) = S$, то есть $W(\sigma_0)$ – цепочка, выводимая в G_W .

Предположим, что для некоторого $i, i < k$, установлено, что

$$S = W(\sigma_0) \Rightarrow_{G_W} W(\sigma_1) \Rightarrow_{G_W} \dots \Rightarrow_{G_W} W(\sigma_i)$$

Покажем, что в этом случае $W(\sigma_i) \Rightarrow_{G_W} W(\sigma_{i+1})$.

Рассмотрим $i+1$ -й этап левого вывода предложения x в грамматике G : $\sigma_i \Rightarrow_G \sigma_{i+1}$. По определению левого вывода:

$$\sigma_i = zC\gamma \Rightarrow_G z\beta\gamma = \sigma_{i+1} \text{ причем } C \rightarrow \beta \text{ есть правило грамматики } G.$$

Из-за наличия в грамматике G множества делимых-слева нетерминалов $LD(a)$ правило $C \rightarrow \beta$ может быть:

▷ либо типа 1: $A \rightarrow a\alpha$ и тогда

$$\begin{aligned} W(\sigma_i) &= W(zA\gamma) = zaA'W(\gamma), \\ W(\sigma_{i+1}) &= W(za\alpha\gamma) = zaW(\alpha)W(\gamma), \end{aligned}$$

то есть $zaA'W(\gamma) \Rightarrow_{G_W} zaW(\alpha)W(\gamma)$ посредством правила $A' \rightarrow W(\alpha)$ из P_W ;

▷ либо типа 2: $A \rightarrow A_0\alpha$ и тогда

$$\begin{aligned} W(\sigma_i) &= W(zA\gamma) = zaA'W(\gamma), \\ W(\sigma_{i+1}) &= W(zA_0\alpha\gamma) = zaA'_0W(\alpha)W(\gamma), \end{aligned}$$

то есть $zaA'W(\gamma) \Rightarrow_{G_W} zaA'_0W(\alpha)W(\gamma)$ посредством правила $A' \rightarrow A'_0W(\alpha)$ из P_W ;

▷ либо типа 3: $B \rightarrow \alpha$ и тогда

$$\begin{aligned} W(\sigma_i) &= W(zB\gamma) = zBW(\gamma), \\ W(\sigma_{i+1}) &= W(z\alpha\gamma) = zW(\alpha)W(\gamma), \end{aligned}$$

то есть $zBW(\gamma) \Rightarrow_{G_W} W(\alpha)W(\gamma)$ посредством правила $B \rightarrow W(\alpha)$ из P_W .

Во всех трех случаях $W(\sigma_i) \Rightarrow_{G_W} W(\sigma_{i+1})$ и, следовательно, предложение $x = W(\sigma_k)$ выводимо в грамматике G_W , а значит $L(G) \subseteq L(G_W)$. Утверждение 1 доказано.

Утверждение 2. $L(G) \supseteq L(G_W)$ для $KC\#$ грамматики G .

Доказательство проведем индукцией по длине левого вывода некоторого произвольного предложения x из $L(G_W)$.

$$S = \sigma_0 \Rightarrow_{G_W} \sigma_1 \Rightarrow_{G_W} \dots \Rightarrow_{G_W} \sigma_{k-1} \Rightarrow_{G_W} \sigma_k = x$$

Покажем, что:

- 1и) в цепочках $\sigma_0, \sigma_1, \dots, \sigma_{k-1}, \sigma_k$ перед каждым символом из N'_A размещается символ a ; и
2и) последовательность цепочек $W^{-1}(\sigma_0), W^{-1}(\sigma_1), \dots, W^{-1}(\sigma_{k-1}), W^{-1}(\sigma_k)$ есть левый вывод предложения x .

Цепочка $W^{-1}(\sigma_0)$ состоит из единственного символа S и поэтому она удовлетворяет условиям 1и) и 2и).

Предположим, что для некоторого $i, i < k$, установлено, что:

- 1п) в цепочках $\sigma_0, \sigma_1, \dots, \sigma_i$ перед каждым символом из N'_A размещается символ a ; и
2п) последовательность цепочек $S = W^{-1}(\sigma_0) \Rightarrow_G W^{-1}(\sigma_1) \Rightarrow_G \dots \Rightarrow_G W^{-1}(\sigma_i)$ есть левый вывод предложения $W^{-1}(\sigma_i)$.

Рассмотрим переход $\sigma_i \Rightarrow_{G_W} \sigma_{i+1}$ в левом выводе предложения x . По определению левого вывода $\sigma_i = zC\gamma$, $\sigma_{i+1} = z\beta\gamma$ и $C \rightarrow \beta$ есть правило грамматики G_W . Отметим два обстоятельства.

Во-первых, в цепочке γ перед каждым символом из $N'(a)$ размещается символ a . Это следует из того, что

- в $zC\gamma$ по индуктивному предположению 1п) перед каждым символом из N'_A размещается символ a ; и
- в $zC\gamma$ цепочка γ располагается непосредственно за нетерминальным символом C и поэтому не может начинаться символом из $N'(a)$.

Во-вторых, по построению грамматики G_W правило $C \rightarrow \beta$ может быть

▷ либо $A' \rightarrow \alpha$, если оно получено из правила первого типа $A \rightarrow aW^{-1}(\alpha)$ из G , и тогда $\sigma_i = zC\gamma = zA'\gamma = yaA'\gamma$ на основании индуктивного предположения 1п),

$$\begin{aligned} W^{-1}(\sigma_i) &= W^{-1}(yaA'\gamma) = W^{-1}(y)W^{-1}(aA')W^{-1}(\gamma) = yAW^{-1}(\gamma), \\ W^{-1}(\sigma_{i+1}) &= W^{-1}(ya\alpha\gamma) = W^{-1}(y)W^{-1}(a\alpha)W^{-1}(\gamma) = yaW^{-1}(\alpha)W^{-1}(\gamma), \end{aligned}$$

то есть $yAW^{-1}(\gamma) \Rightarrow_G yaW^{-1}(\alpha)W^{-1}(\gamma)$ посредством правила $A \rightarrow aW^{-1}(\alpha)$;

▷ либо $A' \rightarrow A'_0\alpha$, если оно получено из правила второго типа $A \rightarrow A_0W^{-1}(\alpha)$ из G , и тогда $\sigma_i = zC'\gamma = zA'\gamma = yaA'\gamma$ на основании индуктивного предположения 1п),

$$\begin{aligned} W^{-1}(\sigma_i) &= W^{-1}(yaA'\gamma) = W^{-1}(y)W^{-1}(aA')W^{-1}(\gamma) = yAW^{-1}(\gamma), \\ W^{-1}(\sigma_{i+1}) &= W^{-1}(yaA'_0\alpha\gamma) = W^{-1}(y)W^{-1}(aA'_0\alpha)W^{-1}(\gamma) = yA_0W^{-1}(\alpha)W^{-1}(\gamma), \end{aligned}$$

то есть $yAW^{-1}(\gamma) \Rightarrow_G yA_0W^{-1}(\alpha)W^{-1}(\gamma)$ посредством правила $A \rightarrow A_0W^{-1}(\alpha)$;

▷ либо $B \rightarrow \alpha$, если оно получено из правила третьего типа $B \rightarrow W^{-1}(\alpha)$ из G , и тогда $\sigma_i = zC'\gamma = zB\gamma$,

$$\begin{aligned} W^{-1}(\sigma_i) &= W^{-1}(yB\gamma) = W^{-1}(y)W^{-1}(B)W^{-1}(\gamma) = yBW^{-1}(\gamma), \\ W^{-1}(\sigma_{i+1}) &= W^{-1}(y\alpha\gamma) = W^{-1}(y)W^{-1}(\alpha)W^{-1}(\gamma) = yW^{-1}(\alpha)W^{-1}(\gamma), \end{aligned}$$

то есть $yBW^{-1}(\gamma) \Rightarrow_G yW^{-1}(\alpha)W^{-1}(\gamma)$ посредством правила $B \rightarrow W^{-1}(\alpha)$.

Во всех трех случаях $W(\sigma_i) \Rightarrow_G W(\sigma_{i+1})$ и, следовательно, предложение $x = W(\sigma_k)$ выводимо в грамматике G , а значит $L(G) \supseteq L(G_W)$. Утверждение 2 доказано.

Окончательно имеем следующее

Утверждение 3. $L(G) = L(G_W)$ для КС#грамматики G .

Другими словами, LD-преобразование грамматик является эквивалентным преобразованием.

Если КС#грамматика G одновременно является LL(1)-грамматикой [1], то в LD-множества могут попасть только простые¹ нетерминалы. Отсюда следует, что LD-преобразование сохраняет класс LL(1)-грамматик.

4. RD-ПРЕОБРАЗОВАНИЕ КС#ГРАММАТИК

Аналогично LD-преобразованию вводится RD-преобразование КС#грамматик, основанное на множестве RD-нетерминалов делимых-справа.

Подмножество нетерминалов КС#грамматики $G = \langle N, \Sigma, P, S \rangle$ будем называть делимым-справа относительно терминала a , и будем его обозначать $RD(a)$, если множество правил вывода P образуют правила трех типов:

$$\begin{aligned} \text{тип } 1': & A \rightarrow \alpha a, \quad \text{где } A \in RD(a) \text{ и } \alpha \neq \epsilon; \\ \text{тип } 2': & A \rightarrow \alpha A_0, \quad \text{где } A \in RD(a) \text{ и } A_0 \in RD(a) \\ \text{тип } 3': & B \rightarrow \alpha, \quad \text{где } B \notin RD(a) \end{aligned}$$

Если в КС-грамматике зафиксировано некоторое подмножество делимых-справа нетерминалов $RD(a)$, то в такой грамматике можно рассматривать преобразование цепочек V . Цепочка $V(\alpha)$ получается из цепочки α посредством замены всех символов A из $RD(a)$ на цепочки из двух символов $A'a$, где $A' \in RD'(a)$.

Пусть $RD(a)$ – подмножество делимых-справа нетерминальных символов КС#грамматики $G = \langle N, \Sigma, P, S \rangle$, определим грамматику $G_V = \langle N_V, \Sigma, P_V, S \rangle$ следующим образом:

¹ Простым называется нетерминал A , для которого в грамматике имеется ровно одно правило $A \rightarrow \alpha$ (A -правило).

$N_V = (RD'(a) \cup N) \setminus RD(a)$, а

P_V построено так:

- если $A \rightarrow \alpha a$ — это правило типа 1' грамматики G , то $A' \rightarrow V(\alpha) \in P_V$;
- если $A \rightarrow \alpha A_0$ — это правило типа 2' грамматики G , то $A' \rightarrow V(\alpha)A'_0 \in P_V$;
- если $B \rightarrow \alpha$ — это правило типа 3' грамматики G , то $B \rightarrow V(\alpha) \in P_V$.

Можно показать, что $L(G) = L(G_V)$, то есть RD-преобразование КС#грамматики G в КС#грамматику G_V относительно некоторого непустого множества RD-нетерминалов является эквивалентным преобразованием.

5. РЕАЛИЗАЦИЯ ЭКВИВАЛЕНТНЫХ ПРЕОБРАЗОВАНИЙ КС-ГРАММАТИК

Разнообразие преобразований КС-грамматик, а также необходимость их многократного повторения порождает задачу совместного использования эквивалентных преобразований. Как организовать процесс трансформации заданной грамматики с тем, чтобы результирующая грамматика уже не допускала ни одного заданного преобразования? Сложность состоит в том, что одно преобразование может удалять из грамматики А-конструкции и пополнять грамматику Б-конструкциями, а другое преобразование может поступать ровно наоборот.

Прежде всего, отметим, что с алгоритмической точки зрения каждое преобразование можно представить в виде, приведенном на рис. 1.



Рис. 1. Схема преобразования КС-грамматики

Например, в LD-преобразовании

- этап “Выявить” состоит в нахождении некоторого непустого множества $LD(a)$;
- этап “Преобразовать” состоит в построении грамматики G_W относительно найденного множества $LD(a)$.

Фактически по успешной ветке может передаваться некоторая информация, существенная для преобразования.

Не вдаваясь в подробности этапов, будем изображать преобразование в виде прямоугольника со сглаженными углами (рис. 2), в который:

- сверху входит стрелка, соответствующая исходной грамматике;
- направо выходит стрелка, соответствующая измененной грамматике; и
- вниз выходит стрелка, соответствующая исходной грамматике, оставшейся без изменений.



Рис. 2. LD-преобразование

С использованием принятых обозначений универсальная схема управления преобразованиями грамматик выглядит так как изображено на рис.3.

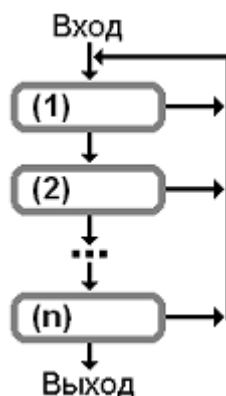


Рис. 3. Универсальная схема преобразования КС-грамматик;
(1), (2), ... (n) – конкретные преобразования КС-грамматик.

Нетрудно видеть, что универсальная схема

- полностью определяется последовательностью эквивалентных преобразований;
- имеет определенные достоинства; и
- порождает некоторые проблемы.

Достоинства. Универсальная схема гарантирует, что в результирующей грамматике более нельзя выполнить ни одного эквивалентного преобразования (1), (2), ... (n).

Проблемы. Для каждого набора преобразований необходимо доказывать корректность универсальной схемы, то есть необходимо доказывать конечность последовательности преобразований. В отдельных, но важных случаях на доказательстве корректности можно “сэкономить”. В этих случаях конечность процесса эквивалентных преобразований основывается на том, что каждое преобразование уменьшает численную характеристику h исходной грамматики [2].

По определению величина h КС-грамматики $G = \langle N, \Sigma, P, S \rangle$ есть

$$h(G) = \sum_{A \in N} \min\{\text{длина}(x) \mid x \in L(G(A))\}$$

Например, для рассмотренной ранее грамматики G_T имеем:

$$\min\{\text{длина}(x) \mid x \in L(G_T(D))\} = 1;$$

$$\min\{\text{длина}(x) \mid x \in L(G_T(B))\} = 2;$$

$$\min\{\text{длина}(x) \mid x \in L(G_T(A))\} = 2;$$

$$\min\{\text{длина}(x) \mid x \in L(G_T(S))\} = 3;$$

и окончательно имеем: $h(G) = 1 + 2 + 2 + 3 = 8$.

Зачастую по результатам эквивалентного преобразования

1 одна часть нетерминалов сохраняет свои реализации неизменными, в том числе, начальный символ S , а

2 другая (непустая) часть нетерминалов:

2.1 либо вообще ликвидируется,

2.2 либо изменяет свою реализацию с $L(G(A))$ на L , где

$$L = \{y \mid xy \in L(G(A))\} \text{ или}$$

$$L = \{y \mid yz \in L(G(A))\} \text{ для некоторых непустых } x \text{ и } z.$$

Преобразование, удовлетворяющее свойствам 1–2, уменьшает величину h . Факт уменьшения величины h будем обозначать h -.

Например, LD-преобразование КС#грамматик относительно некоторого LD(a) подпадает под случай 1–2.2 при $x = a$, $z = \epsilon$. Здесь первую часть нетерминалов составляют $N \setminus LD(a)$, а вторую – $LD(a)$. Нетрудно показать, что $h(G) = h(G_W) - R$, где R – количество нетерминалов в $LD(a)$.

Рассмотрим подвид универсальных схем эквивалентных преобразований, представленный на рис.4. Здесь:

– преобразование (1) – необязательное устранение бесполезных² символов, в том числе:

- нетерминалов, которые не могут породить терминальные цепочки; и
- правил вывода, содержащих недостижимые символы;

– преобразование (2) – необязательное устранение цепных³ правил;

– преобразования (3)..(n) – такие преобразования, которые уменьшают характеристику h .

Преобразования (1) и (2) давно и хорошо изучены, в общем случае они не изменяют характеристику грамматики h , однако их использование совместно или порознь не способны привести к заикливаю. Что касается остальных преобразований, то их выполнение порождает монотонно убывающую последовательность положительных чисел h_1, h_2, h_3, \dots . Такая последовательность не может быть бесконечной, а значит корректность подвида универсальных схем эквивалентных преобразований установлена.

К преобразованиям (3)..(n), в частности, относятся:

- устранение простых нетерминалов - случай 1+2.1 (h-);
- устранение нерекурсивных⁴ нетерминалов - случай 1+2.1 (h-);
- устранение избыточных нетерминалов [2] - случай 1+2.1 (h-);

² Бесполезным [1] в грамматике $G = \langle N, \Sigma, P, S \rangle$ называется символ $X \in N \cup \Sigma$, для которого в грамматике нет вывода вида $S \Rightarrow *_G yXz \Rightarrow *_G yxz$.

³ Цепным [1] называется правило вывода вида $A \rightarrow B$.

⁴ Нерекурсивным называется нетерминал A , для которого не существует выводов вида $A \Rightarrow^+ \alpha A \beta$.

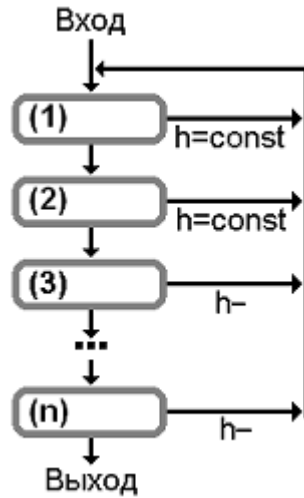


Рис. 4. Частный случай эквивалентных преобразований

- ЛНФ- и ПНФ-преобразования⁵ - случай 1+2.2 (h-);
- LD- и RD- преобразования - случай 1+2.2 (h-).

Особого разговора заслуживает порядок размещения эквивалентных преобразований в универсальной схеме. С точки зрения свойств конечного результата порядок не играет роли, однако иногда частные преобразования имеет смысл выполнять раньше общих преобразований.

Так, в некоторых случаях ЛНФ-преобразования [2] можно рассматривать как частный случай LD-преобразований. Если, например A-правила грамматики имеют вид

$$A \rightarrow aB \mid aC \mid aBA_d,$$

то ЛНФ-преобразование терминала A совпадает с LD-преобразованием относительно $LD(a) = \{ A \}$. Вместе с тем ЛНФ-преобразование способно обрабатывать случаи, когда все правые части A-правил начинаются одним и тем же нетерминалом. При этом ЛНФ-преобразование действует вполне “разумно”, преобразуя

$$A \rightarrow Bb \mid BCc \mid BaAd \quad \text{в} \quad A' \rightarrow b \mid Cc \mid aBA'd.$$

LD-преобразование в этом случае действует более “топорно”, преобразуя

$$A \rightarrow Bb \mid BCc \mid BaAd \quad \text{в} \quad A' \rightarrow B'b \mid B'Cc \mid B'abA'd.$$

6. ЗАКЛЮЧЕНИЕ

Настоящая работа носит исключительно теоретический характер, ее главная цель – расширить спектр возможностей при выдвижении гипотез о строении неизвестной грамматики, породившей некоторые известные предложения. Если, например, известно, что $LL(1)$ грамматика породила два предложения $abcabdd$ и $abcbeddd$, то с определенными оговорками можно считать, что оба эти предложения в искомой грамматике имеют общую сентенциальную форму $abcA''dd$. В самом деле:

⁵ ЛНФ-преобразование (ПНФ-преобразование) [2] заключается в устранении явно указанных общих префиксов (суффиксов) в правых частях нетерминалов.

- быть LL(1) означает наличие общей формы $abcA$, где $A \rightarrow abdd \mid bcddd$, в общем виде известной как “дерево суффиксов” [3];
- доказанная допустимость RD-преобразований фактически означает наличие общей формы $abcA''dd$, где $A'' \rightarrow ab \mid bcd$.

Упомянутые оговорки будут раскрыты в следующих работах.

СПИСОК ЛИТЕРАТУРЫ

1. Ахо А., Ульман Дж. *Теория синтаксического анализа, перевода и компиляции*. М.: Мир, 1978, тт. 1,2.
2. Соловьев С.Ю. Нормализация контекстно-свободных грамматик для целей грамматического вывода. *XII национальная конференция по искусственному интеллекту с международным участием КИИ-2010. Труды конференции*. М: Физматлит, 2010, том 1, стр.218-224.
http://www.park.glossary.ru/serios/read_09.php
3. Смит Б. *Методы и алгоритмы вычислений на строках*. М: Вильямс, 2006.