HOCTPOEME TR -- TPAMMATHK HO FPAMMATHKAM B EEKYCOBO-HAYPOBCKON GOPME

В.И. Константинов

Вопросы синтаксического анализа входных языков занимают одно из центральных мест в современной теории и практике автомативации программирования. Процесс синтаксического анализа сущест венно зависит от метаязыка, используемого для задания грамматик входных языков.

В [1,2] разработан метаязык \mathcal{R} -грамматик, допускакций простне и быстро работающие алгоритмы синтаксического анализа. Однако процесс перехода от общепринятой бэкусово-науровской формы (НФ) записи грамматик языков к задающим их \mathcal{R} -грамматикам оказался трудоемок.

Действительно, алгориты синтеза R -грамматик [3] по НФ виличает следующие этапы:

- переход от НФ к модифицированной НФ (МНФ);
- выделение канонической формы МНФ (КФ);
- введение А-термов;

1

- определение базовых А-термов КФ (МНФ) грамматики;
- построение для каждого начального терма множества θ -це почек и т.д.

Ввиду этого была предложена [4] новая форма задания грамматик языков — запись их TR-грамматиками, которые являются рас — ширением R-грамматик на нетерминальный алфавит. Переход от НФ к TR-грамматикам проце, чем к R-грамматикам.

І. Под грамматикой в НФ будем понимать четверку

$$G = (A, V_{\kappa}, S_{o}, P),$$

где $A = \{\alpha_1, \alpha_2, ..., \alpha_n\}$ терминальный алфавит,

 $\bigvee_{N} = \{S_{o}, S_{1}, \dots, S_{m}\}$ - нетерминальный алфавит, $A \cap \bigvee_{N} = \emptyset$; \mathcal{S}_{o} — аксиома грамматики языка (начальный символ);

 \mathcal{P} — множество правил вида $\mathcal{S}_i ::= \beta$, где β - цепочка над алфавитом AUV_N , то есть $\beta \in (AUV_N)^*$, или множество цепочек над А U V_N , разделенных знаком I (исключительно "или"), то есть $\beta = \beta_1 | \beta_2 | \dots | \beta_p; \quad \beta_i \in (A \cup V_N)^*, 1 \le i \le p.$ Пусть $\alpha, y \in (A \cup V_N)^*$. Скажем, что из α следует $y \in (A \cup V_N)^*$.

если существуют $z_1, z_2, \beta \in (A \cup V_N)^*$ и $S \in V_N$ $x = z_1 S z_2, y = z_1 \beta z_2, S := \beta$ — правило из P .

Скажем, что из x выводимо y $(x \Rightarrow y)$,если существует цепочка w_0, \dots, w_n над алфаеитом AUV_N такая, что $w_0 \rightarrow w_1, \dots, w_{n-1} \rightarrow w_n$, THE $w_0 = \infty$, $w_0 = y$.

Множество $L(G) = \{ x \in A^*/S_o \Rightarrow x \}$ назовем языком, порожденным грамматикой G .

Ограничим класс грамматик в НФ, для которых строятся ТК грамматики, грамматиками, не содержащими правил вида $S := \mathcal{E}, \mathcal{E}$ -пустая цепочка. Ограничение это несущественно, т.к. в [5] показана вквивалентность контекстно-свободных (КС-языков) и языков, порождаемых грамматиками в НФ (в [5] они названы языками типа АЛГОЛ); там же показано, что если G - KC-грамматика и $\varepsilon \in L(G)$, то по грамматике С может быть эффективно построена КС-грамматика, порождающая язык L(G) , но не содержащая правил вида $S::=\varepsilon$. Грамматини в $H\Phi$, не содержащие правил вида $S::=\mathcal{E}$, назовем приведенными. Исключение правил вида $\mathcal{S} ::= \mathbf{s}$ связано с необходимостью устранения недетерминированности, возникающей в таких граммати ках при анализе входной строки.

Действительно, пусть G содержит правила:

$$S_4 := \varphi S_2 \varphi,$$

$$S_p := \varepsilon |\alpha_1| \dots |\alpha_n$$

причем одно из α_i таково, что $\alpha_i \Rightarrow \psi \beta; \delta_{ij} \delta_{j} \delta_{k} \delta_{j} \in V_N, \psi, \alpha_{ij}, \dots, \alpha_{in}, \beta \in (AUV_N)^*$ $\psi \in A^*$, ψ , \mathcal{L}_i , $\psi \neq \mathcal{E}$. Тогда при анализе цепочки $\psi \psi$ в момент прихода $A \ni \alpha$ -первого алемента цепочки ψ , нельзя сказать, следует ли распознавать его в S_2 , или остаться в S_4 .

П. В качестве исходной формы задания грамматик входных языков предлагается, как и в [6], модифицированная НФ (МНФ). Задача построения по произвольной НФ соответствующей МНФ алгоритмически разрешима, алгоритм перехода подробно описан в [5]. От метим, что по сравнению с ЕНФ:

- а) в МЕНФ отсутствует леволинейная рекурсия (то есть не существует $S \in V_N$ и непустой цепочки $\varphi \in (A \cup V_N)^*$ такой, что $S \Rightarrow S \varphi$; она исключена с помощью итерационных скобок):
 - б) никакие два правила не имеют совпадающих девых частей.

. ПРИМЕР применения итерационных скобок. $S ::= S \varphi | \psi$ с помощью итерационных скобок можно записать в виде $\mathcal{S} := \psi \ \{ \ \psi \}$ означает, что φ может быть повторено нуль или более раз.

Ш. Рассмотрим грамматику в МНФ, не содержащую правил с вложенными итерационными скобками (то есть правил вида $\{\{\varphi\},\{\psi\}\}$). Пусть она задается множеством правил вида:

$$S_o := \alpha_o, S_1 := \alpha_1, \ldots, S_m := \alpha_m,$$

где $S_i \in V_N$, $0 \le i \le m$, α_i есть

21

а) либо цепочка, принадлежащая $(AUV_{N})^{*} \setminus \epsilon$;

б) либо цепочка, принадлежащая $(AUV_{\nu})^* \setminus g$, некоторые полцепочки которой заключены в { } ;

в) либо цепочки вида а)-б), разделенные знаком

С каждым негерминальным символом S связана атомная TR-грамматика [4]: $G(S)=(A\,,\,\bigvee_N\,\mathcal{P},\,\mathcal{R}^{G(S)}_{\circ},\,\mathcal{R}^{G(S)}_{\circ},\,\mathcal{R}^*_{G(S)})$.

Интерпретация правил ТР -грамматик подробно описана в [4]. Поскольку каждое правило с нетерминальной левой частью интерпретируется с помощью стека, переход от грамматик в НФ к грамматикам в МЕНФ сокращает стековые операции анализатора и гарантирует от зацикливания (между любыми двумя обращениями к одной и той же грамматике будет считан по крайней мере один символ входной стро-

Исходной грамматике G , заданной системой (I), соответст вует система атомных TR-грамматик $G(S_o), \ldots, G(S_m)$. Построить TRграмматику G - значит построить образующие её грамматики $G(S_s)$ O≤i≤m.

Возьмем любое правило из (I). Пусть это правило имеет вид $S := \beta_1 | \beta_2 | \dots | \beta_n$, β_i — цепочки из $(A \cup V_v)^* \setminus \varepsilon$, некоторые подцепочки которых могут быть заключены в итерационные скобки.

Алгоритм построения атомной TR-грамматики G(S) есть следующая последовательность действий.

І. Определение имени формируемого множества правил. Для грамматики G(S) на первом шаге это имя есть $z_o^{G(S)}$ (аксиома грамматики G(S)).

2. Выбор первых элементов цепочек β_1,\ldots,β_m ; пусть это будут x_1,\ldots,x_k . Так как для цепочек β_i вида $\{\theta_i\}\ldots\{\theta_p\}$ ψ первыми элементами являются начальные элементы цепочек $\theta_n,\ldots,\theta_p,\psi$ число k может превосходить n.

Сформируем множество \mathcal{M} , в которое попадут все различные элементы из x_1,\dots,x_k . Множество \mathcal{M} назовем множеством левых частей формируемого множества правил (на первом маге это множе — ство левых частей правил из $z_o^{\mathcal{G}(S)}$, что обозначим $\mathcal{M}(z_o^{\mathcal{G}(S)})$).

Число элементов в M может меняться от I (когда все $x_1, ..., x_k$ совпадают) до k (когда они все различны). Элемент $m \in M$ назовем элементом кратности ℓ , если он ℓ раз, $1 \leqslant \ell \leqslant k$, встречается в последовательности $x_1, ..., x_k$.

- 3. Определение типов правил из $\mathcal{C}_o^{\mathcal{L}(\mathcal{S})}$, для чего анализи руются подцепочки, следующие за m. Возможны следующие случаи:
 - а) m либо $m \mid$, то есть m- единственный алемент цепочки;
 - 6) my, y∈(AUV,)*\E:
 - в) $m\{y_n\}...\{y_n\}$ либо $m\{y_n\}...\{y_n\}|, n \ge 1, y_i \in (AUV_n)^* \setminus E;$
 - r) $m\{y_i\}...\{y_n\}\psi, \ \psi_i, \psi \in (AUV_n)^* \setminus \varepsilon, \ n > 1;$
 - $\Pi \mid m \mid \{ \varphi_i \} \dots \{ \varphi_n \} \; \psi, \quad n \geq 1, \quad \psi, \, \varphi_i \in (A \cup V_N)^* \setminus \varepsilon.$

Нроме того, если m не есть первый алемент цепочки (то есть $m \in \mathcal{M}(z_o^{\mathcal{A}(\mathcal{S})})$), следует добавить еще

е) m $\}$ $\{\mathcal{Y}_n\}$ \dots $\{\mathcal{Y}_n\}$ \dots $\{\mathcal{Y}_n\}$ $\}$, $\mathcal{Y}_i \in (A \cup V_n)^n \setminus \mathcal{E}, n \gg 1$. Если m — кратности I, то в случанх a , b , c , d) правило с левой частью m есть правило первого типа, в остальных случанх — правило второго типа. Если формируемое правило второго типа, мно-жество правил—преемников его разомкнуто.

Если кратность m>1, то имеется несколько подцепочек (число их совпадает с кратностью m), следующих за m. В этом случае тип правила определяется следующим образом:

- если все подцепочки, начинающиеся с m, есть α подцепочки, формируемое правило первого типа, при этом имя множества правил-преемников данного правила 2ϕ ;
- если все подцепочки, начинающиеся с m, есть б-,либо г-, либо д-подцепочки, соответствующее правило есть правило первого типа. Во всех остальных случаях формируемое правило есть правило второго типа, а множество правил-преемников его разомкнуто.
- 4. Определение правил-преемников формируемого множества правил. Если m имеет кратность I, множество правил-преемников его состоит из:

- правила с первым элементом подцепочек ϕ в случае б);
- правил с первыми различными алементами подцепочек $\mathcal{G}_1, \dots, \mathcal{G}_n$ в случае в);
- правил с первыми различными элементами подценочек $g_1, \dots, g_n, \ \Psi$ в случае \mathbf{r}).

В случае а) множество правил-преемников пусто, имя его – z_{\emptyset} . Отдельно рассмотрим ситуацию, когда за m следует $\}$, что возможно в случалх д) и е).

Пусть m - последний элемент цепочки φ_o , тогда множество правил-преемников формируемого правила состоит из:

- + правил с различными элементами подцепочек $\mathcal{G}_0,\mathcal{G}_1,\dots,\mathcal{G}_n$ в случае е).

Заметим, что цепочка \mathscr{G} один раз уже пройдена, т.к. m-последний элемент этой цепочки, но в общем случае необходимо пройти по ней еще один раз. Действительно, пусть имеется, например, цепочка $\mathscr{G}\{\mathscr{G}_1\}\{\mathscr{G}_2\}\{\mathscr{G}_3\}$ \mathscr{V} — и пусть m— последний элемент \mathscr{G}_2 .

Положим $\mathcal{G}_{1} = S\mathcal{G}_{1}'$, $\mathcal{G}_{2} = S\mathcal{G}_{2}'$, тогда первое прохождение по цепочке \mathcal{G}_{2} при определении преемников правила с левой частью \mathcal{E}_{3} почтому, хотя цепочки \mathcal{G}_{3}) совмещено с прохождением цепочки \mathcal{G}_{4} , поэтому, хотя цепочка \mathcal{G}_{2} уже пройдена, нельзя вос пользоваться полученной ранее цепочкой правил, так как, например, правило с левой частью S будет иметь преемников и из \mathcal{G}_{4}' , и из \mathcal{G}_{2}' . Эторой проход по цепочке \mathcal{G}_{4} из последовательности \mathcal{G}_{4}' \mathcal{G}_{5} совмещается с прохождением по цепочкам \mathcal{G}_{2+1}' , \mathcal{G}_{2} . При этом множество левых частей правил-преемников последнего элемента цепочки \mathcal{G}_{6} — то же, что и последнего элемента цепочки \mathcal{G}_{6} , \mathcal{G}_{6

Если известно, что в последовательности $\varphi\{\varphi_j\}\dots\{\varphi_n\}$ φ подцепочки φ_j , φ_j , φ_j начинаются с разных алементов, процедура нахождения множеств правил-преемников существенно упрощается. Пусть z_m -имя множества правил-преемников последнего элемента цепочки φ . Тогда множество правил-преемников последнего алемента цепочки φ_j есть $(z_m, 1)$, множество правил-преемников последнего алемента цепочки φ_j есть (z_m, j) и т.д. Индекс φ_j показывает, начиная с какого по счету правила множества правил z_m стоят правила-преемники данного правила.

94

Если $m \in \mathcal{M}$ имеет кратность > I (скажем, ℓ), считаем, что имеется ℓ элементов кратности I,для которых определяем множе — ства правил-преемников так, как это указано в п.4, после чего берем пересечение этих множеств правил и объявляем его множеством правил-преемников с левой частью m.

Выполнив п.4, мы сформировали новые множества M. Присвоив очередные имена из числа $Z_1^{G(S)}$, ..., $Z_n^{G(S)}$ соответствующим множествам правил, переходим на п. 3. Так как правая часть правила $S := \beta_1 \setminus \dots \setminus \beta_m$ состоит из конечного числа элементов, за конечное число шагов T'R—грамматика G(S) будет построена. Затем строим следующую атомную грамматику и т.д., до полного построения $G = \bigcup_{i=1}^n G(S_i)$

IV. В заключение рассмотрим простой пример, иллюстрирующий приведенный алгоритм.

Пусть дана грамматика $G = \{A, V_N, S_o, P\}$, где $A = \{\alpha, \beta, c\}$, $V_N = \{X, Y\}$, $S_o = X$, $P = \{X ::= X\alpha \mid Y, Y ::= \beta \mid c\}$. Соответствующая МНФ имеет вид: $X ::= Y\{\alpha\}, \ Y ::= \beta \mid c$.

Построим TR -грамматику G(X) . I. Определим имя формируемого множества правил - это $z_o^{G(X)}$.

1. Определим ими формиру емого высмество ото состоит из алемента $\mathcal{Y} \in V_N$ кратности I, то есть $Z_o^{G(X)}$ содержит единственное правило, левая часть которого \mathcal{Y} .

3. Определим тип формируемого правила. Цепочка $\mathcal{Y}\{\alpha\}$ - типа \mathcal{E} , поэтому формируемое правило – второго типа: $z_o^{\mathcal{E}(X)} = \{\mathcal{Y} - \cdot\}$, множество правил-превиников его разомкнуто. Присвоим ему имя $z_f^{\mathcal{E}(X)}$, $z_o^{\mathcal{E}(X)} = \{\mathcal{Y} - \cdot z_f^{\mathcal{E}(X)}\}$.

4. Сформируем $\mathcal{E}_{q}^{G(X)}$. Как отмечалось, множество это ра — зомунуто, поэтому выделим его в соответствии с [I], в } { . Оно состоит из правила с левой частью α , $\mathcal{E}_{q}^{G(X)} = \alpha$. Поскольку α — последний элемент цепочки, стоящей в $\{ \}$ (случай е), данное правило-второго типа: $\mathcal{E}_{q}^{G(X)} = \}\alpha \rightarrow \{$. Так как α -последний элемент в единственной из цепочек типа $\{\mathcal{G}_{q}\}\dots\{\mathcal{G}_{n}\}$, второй проход по цепочке не нужен, а множество правил-преемников правила с левой частью α — то же, что и \mathcal{G}_{q} . Получаем, таким образом, что \mathcal{G}_{q} —грамматика \mathcal{G}_{q} есть:

$$G(X): \qquad z_o^{G(X)} = \{ y - z_r^{G(X)} \}$$

$$z_r^{G(X)} = \{ a - z_r^{G(X)} \}$$

IІостроим теперь G(Y) .

I. Имя формируемого множества правил — $\mathcal{P}_o^{G(\mathcal{Y})}$. 2. $\mathcal{M}(\mathcal{P}_o^{G(\mathcal{Y})}) = \{\mathcal{B}, \mathcal{C}\}$, это элементы кратности I, $\mathcal{P}_o^{G(\mathcal{Y})}$ состоит из двух правил с девыми частями \mathcal{E} и \mathcal{C} .

3. Определим типы этих правил. Так как подцепочки, начинакщиеся с ℓ , c есть α -подцепочки, данные правила-первого типа, а имена их правил-пресыников – $7_{\mathcal{O}}$.Получаем

$$G(Y): r_o^{GXX} = \{ \beta \rightarrow r_o, c \rightarrow r_o \}.$$

Таким образом, соответствующая G система атомных грамматик имеет вид:

$$G(X): \quad z_o^{G(X)} = \{ y \rightarrow z_1^{G(X)} \}$$

$$z_i^{G(X)} = \{ a \rightarrow z_1^{G(X)} \}$$

$$G(Y): \quad z_o^{G(Y)} = \{ \ell \rightarrow z_{\phi}, c \rightarrow z_{\phi} \}.$$

Литература

- I. ВЕЛЬНИКИИ И.В. О метаязыке синтаксически управляемого транслятора.-"Вычислительные системы," новосибирск, 1970, вып. 42, с. 22-33.
- 2. ВЕЛЬНИКИЙ И.В., КЩЕЖО Е.Л.Метаязык, ориентированный для синтаксического анализа и контроля.— "Кибернетика," Киев, 1970, №2, с. 50-53.
- 3. ВЕЛЬНИКИЙ И.В. К вопросу перехода от бекусово-науровс кой формы записи грамматик к грамматикам R -типа. -Автоматизация программирования, Киев, 1969, № 3, с. 43-59.
- 4. КОНСТАНТИНОВ В.И., НУРИЕВ Р.М. Метаязык трансляции кон текстно-свободных языков. Настоящий сборник, с.74-83.
- 5. ГИНЗБУРГ С. Математический анализ контекстно-свободных языков. М., "Мир", 1970.
- 6. ВЕЛЬНИНИЙ И.В., МЕЛВЕЛЕВА В.Н., ХИЛЬЧЕНКО В.И. О фор мальном описании синтаксиса языка АЛГОЛ—60 и транслятора ТА 2, —Автоматизация программирования, Киев, I969, № 3, с. 85-92.

Поступила в ред.-изд.отд. 6 февраля 1973 года

1