

О МЕТАЯЗЫКЕ СИНТАКСИЧЕСКИ УПРАВЛЯЕМОГО ТРАНСЛЯТОРА

И.В. Вельбицкий

Эффективность работы синтаксически управляемого транслятора в значительной степени определяется метаязыком, используемым для описания грамматик [1]. В данной работе обсуждается метаязык, ориентированный на применение в синтаксически управляемых трансляторах.

I<sup>0</sup>. Пусть  $S$  - словарь синтермов [2], элементы которого в дальнейшем, где это не оговорено, будут обозначаться малыми начальными буквами латинского алфавита;  $\alpha$  - некоторый вспомогательный словарь, содержащий пустой символ,  $\alpha = \{\emptyset, \nu_i\}_{i \in I}$ ,  $I = \{1, \dots, k\}$ ;  $R$  - множество имен некоторых конечных множеств (возможно, пустых) подмножеств правил вида:

$$\alpha \xrightarrow{W_m} \gamma \quad (I)$$

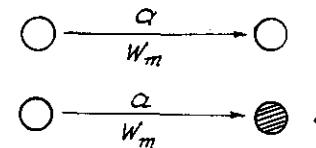
или

$$\alpha \xrightarrow{W_m} \gamma , \quad (2)$$

где  $\alpha \in S$ ;  $\gamma \in R$ ;  $R = \{\emptyset, \gamma_i\}_{i \in I}$ ,  $I = \{1, \dots, l\}$ ,  $\emptyset$  - имя пустого множества правил;  $W_m$  или  $W_m$  - многоместное отношение, означающее, что синтерм слева от него некоторым образом конкатенирует с синтермом левой части любого правила из множества, имя

которого указано справа от  $W_m$  - отношения;  $W_m = W_m(\alpha_1, \alpha_2, \dots, \alpha_k)$ ,  $\alpha_i$  - слово в словаре  $OLUR$ ,  $i = 1, \dots, n$ ; конкатенация синтермов определяется типом  $m$  ( $m = 0, 1, 2$ )  $W_m$ -отношения. Число  $n$  аргументов  $W_m$  - отношения соответствует числу стеков, используемых при интерпретации правил грамматики. При этом предполагается, что на вершине пустых стеков помещены символы  $\emptyset$ ; наличие символа  $\emptyset$  в качестве аргумента  $W_m$  - отношения означает, что соответствующий этому аргументу стек не используется в интерпретации любого правила.

Для приложений удобно правила (I), (2) изображать соответственно в виде графов:

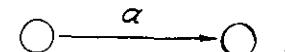


В этих графах дуга соответствует синтерму правила и несет информацию о  $W_m$  - отношении; граничные вершины дуги - именам множества правил; начало дуги - имени множества правил, которому принадлежит данное правило; конец дуги - имени множества правил преемников [3].

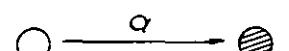
Правило вида (I) интерпретируется следующим образом:

A. При  $m = 0$  (отношение нулевого типа) синтерм  $\alpha$  может конкатенировать с синтермами, встречающимися в левых частях правил из множества с именем  $\gamma$ ; тем самым данное правило определяет свои возможные правила - преемники из множества  $R$ .

$W_0$  - отношение изображается графиком:



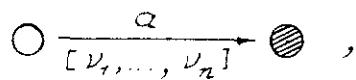
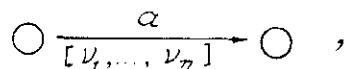
Если в качестве  $\gamma$  стоит символ  $\emptyset$ , то  $\alpha$  соответствует последнему символу предложения языка. В этом случае конец дуги графа выделенется (например, штриховкой):



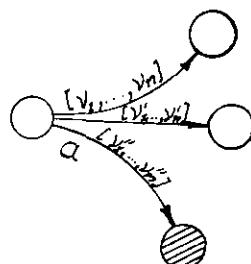
В. При  $m = 1$  (отношение первого типа) указанная конкатенация и преемственность правил имеет место при одновременной записи непустых  $\alpha_i$  в соответствующие  $i$ -е стеки. Очевидно,  $W_1 = W_0$ , если все  $\alpha_i$  из  $W_1$  - пустые символы.  $W_1$ -отношение изображается графом:



С. При  $m = 2$  (отношение второго типа) проверяется условие совпадения всех непустых символов  $v_i \in \Omega$  с символами в вершинах соответствующих стеков. При выполнении условия происходит конкатенация, преемственность правил и выборка этих стеков. При отсутствии такого совпадения ни конкатенация, ни преемственность, ни выборка из соответствующих стеков не имеют места.  $W_2$  - отношения изображаются графами:



если  $\alpha$  соответствует конечному символу некоторого предложения языка. Если дуги, выходящие из одной вершины, отличаются лишь аргументами при  $W_2$  - отношении, то они изображаются в виде метелки:



Все сказанное о правилах вида (1) справедливо и для правил вида (2) с тем отличием, что синтерм  $\alpha$  может быть конечным (может соответствовать конечному символу) некоторого предложения языка даже при  $\zeta$ , отличном от  $\emptyset$ . При этом на графике конец соответствующей дуги выделяется (штриховкой).

Таким образом, интерпретация правил (1),(2) включает определения как допустимого множества конкатенаций с данным синтермом правила, так и множества правил преемников, а также указание на то, что данный синтерм может быть конечным.

1.1. Несколько слов о соответствии термов синтермам. Синтерм является элементом метаязыка, терм - языка. Соответствие между ними зададим функцией  $F = F(t, \tau)$ , где  $t \in T$  - терм,  $T$  - терминалный словарь языка,  $\tau$  - некоторый параметр соответствующего множества из  $\mathcal{R}$ . Каждому имени этого множества поставлен в соответствие некоторый индекс  $\tau$  из множества индексов  $M$ , который в дальнейшем будем записывать в виде верхнего индекса соответствующего имени из  $\mathcal{R}$ . На графике индекс  $\tau$  записывается в вершине, инцидентной началу дуги.



Содержательно: каждому терму в зависимости от его места в предложении языка могут соответствовать различные синтермы. И, наоборот, в каждой локальной конструкции языка одному синтерму могут соответствовать несколько термов. Например, в АЛГОЛЕ - 60 символу "+" можно поставить в соответствие три различных синтерма, которые условно обозначим: "+<sub>1</sub>", "+<sub>2</sub>", "+<sub>3</sub>".

Синтерму "+<sub>1</sub>" соответствует конструкция языка, которую неформально можно назвать как: "перед арифметическим выражением". Этому синтерму соответствует два терма: "+", "-".

Синтерму "+<sub>2</sub>" соответствует конструкция языка, которую можно определить как: "внутри арифметического выражения". Этому синтерму соответствует шесть термов: "+", "-", "x", "/", ":", "/".

Синтерму "+<sub>3</sub>" соответствует конструкции: "внутри комментария". Этому синтерму соответствует 115 термов: 52 буквы латинского алфавита (заглавные и прописные), 10 цифр, 29 ограничителей, 24 буквенных выражения.

Очевидно, что конструкции языка могут быть локализованы

лишь с помощью правил его грамматики. Естественно каждому множеству правил, соответствующему определенной конструкции языка, поставить в соответствие некоторый параметр  $\tau$ , благодаря которому термы, имеющие одинаковое продолжение, можно обозначать одним синтермом и представлять одним правилом в соответствующем множестве правил  $\Sigma^{\tau} \in \mathcal{R}$ .

Таким образом, использование словаря синтермов, определяемого с помощью некоторой функции  $F=F(t, \tau), t \in T, \tau \in \mathcal{M}$ , позволяет существенно уменьшить число правил в каждом множестве  $\Sigma^{\tau}$ , следовательно, во всем множестве  $\mathcal{R}$ . Это на практике позволяет компактно записывать грамматики исходных языков.

1.2. Выше было определено, что  $\Sigma^{\tau}$  – множество подмножеств правил вида (1), (2). Будем в дальнейшем подмножество правил, совпадающее с некоторым множеством правил, имеющим имя  $\tau_i^{\tau}$ , обозначать именем  $\tau_i^{\tau}$ . Следовательно, записи множеств:

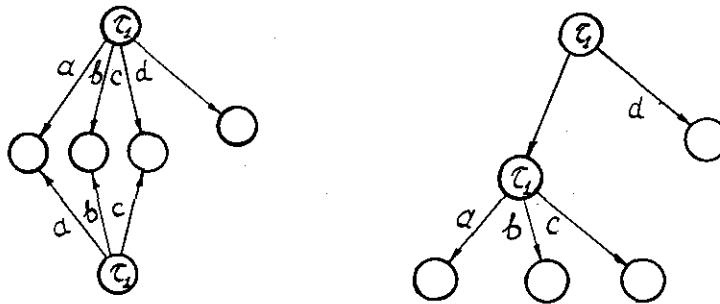
$$\tau_i^{\tau} : \{A_1, \dots, A_p, A_{p+1}, \dots, A_{p+q}\},$$

$$\tau_j^{\tau} : \{A_1, \dots, A_p\},$$

$$\tau_i^{\tau} : \{\tau_j^{\tau}, A_{p+1}, \dots, A_{p+q}\},$$

$$\tau_j^{\tau} : \{A_1, \dots, A_p\}$$

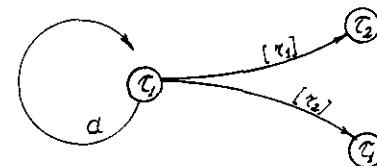
эквивалентны между собой, где  $A_1, \dots, A_{p+q}$  – правила вида (1), (2). На графике указанная ситуация изображается дугой без обозначений, направленной к вершине  $\tau_j^{\tau}$ . Например, изображенные ниже графы эквивалентны:



### I.3. Множество правил

$$\Sigma^{\tau} : \{A_1, \dots, A_q\},$$

элементами которого являются правила множества  $\{A_1, \dots, A_q, \tau\}$ , где  $\tau$  – расположение на вершине стековой памяти имя некоторого множества правил, назовем разомкнутым. На графике разомкнутым множествам соответствуют дуги, в обозначении которых отсутствуют синтермы, например:



Подводя итог, можно сказать, что грамматикой языка, неориентированной для синтаксического анализа (не содержащей определений семантики), понимается семерка:

$$G = (T, S, M, \alpha, F, R, \tau_0), \quad (3)$$

где

$T$  – терминальный словарь языка;

$S$  – словарь синтермов;

$M, \alpha$  – вспомогательные словари,  $M \cap \alpha = \emptyset$ ;

$F = F(t, \tau)$  – функция соответствия термов ( $t \in T$ ) синтермам в зависимости от параметра  $\tau \in M$  грамматики языка;

$R = \{\tau_i^{\tau}\}_{i \in I, \tau \in M}$  – множество имен некоторых конечных множеств (возможно, пустых) подмножеств правил вида (1), (2);

$\tau_0$  – аксиома грамматики, некоторое выделенное начальное множество правил.

ПРИМЕР. Простое арифметическое выражение языка АЛГО-60, в определении [4] которого: <первичное выражение> ::= <целое> <простая переменная> <(простое арифметическое выражение)> – задаётся следующей  $R$  – грамматикой

$$G = (T, S, M, \alpha, F, R, \tau_0),$$

$$T = \{ \alpha, \dots, z, A, \dots, Z, 0, 1, \dots, 9, +, -, \times, /, \div, (, ), \},$$

$$S = \{ \delta, \delta u, u, +, \times, ), ( \},$$

$$M = \{ 1, 2, 3 \}, \quad \Omega = \emptyset,$$

$f$  - функция, которую можно задать таблицей вида:

$\tau^T$	$a$	$b$	$\dots$	$y$	$z$	$A$	$B$	$\dots$	$Y$	$Z$	$0$	$1$	$\dots$	$9$	$+$	$-$	$\times$	$/$	$\div$	$($	$)$
I	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$u$	$u$	$\dots$	$u$	$+$	$+$	$x$	$x$	$x$	$)$	$($
2	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$u$	$u$	$\dots$	$u$	$x$	$x$	$x$	$x$	$x$	$)$	$($
3	$\delta u$	$\delta u$	$\dots$	$\delta u$	$\delta u$	$\delta u$	$\delta u$	$\dots$	$\delta u$	$\delta u$	$\delta u$	$\delta u$	$\dots$	$\delta u$	$x$	$x$	$x$	$x$	$x$	$)$	$($

$$R = \{ z'_o, z'_1, z'_2, z'^{3*}_3, z^{2*}_4, z^2_5, z'_6, z'_7, z^{3*}_8, z^{2*}_9, z^2_{10}, z'_{11} \},$$

$$z'_o: \{ + \xrightarrow{W_o} z, z, \},$$

$$z'_1: \{ \delta \xrightarrow{W_i(z_4)} z_2, u \xrightarrow{W_i(z_4)} z_3, (\xrightarrow{W_i(z_5)} z_6) \},$$

$$z'^{3*}_2: \{ \delta u \xrightarrow{W_o} z_2 \},$$

$$z^{2*}_3: \{ u \xrightarrow{W_o} z_3 \},$$

$$z^2_4: \{ x \xrightarrow{W_o} z_4 \},$$

$$z'_5: \{ ) \xrightarrow{W_o} z_4 \},$$

$$z'_6: \{ + \xrightarrow{W_o} z_7, z_7 \},$$

$$z'_7: \{ \delta \xrightarrow{W_i(z_{10})} z_8, u \xrightarrow{W_i(z_{10})} z'_9, (\xrightarrow{W_i(z_{11})} z_6) \},$$

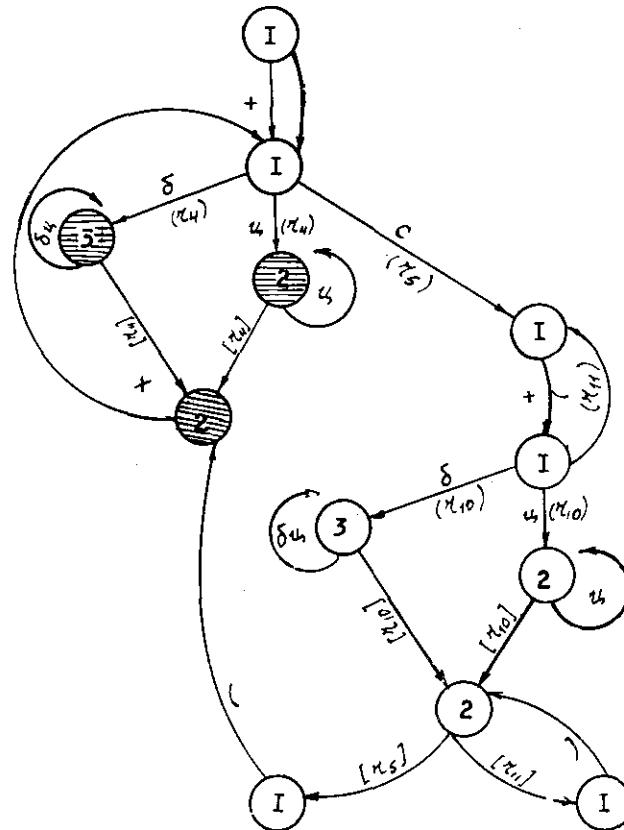
$$z'^{3*}_8: \{ \delta u \xrightarrow{W_o} z_8 \},$$

$$z^{2*}_9: \{ u \xrightarrow{W_o} z_9 \},$$

$$z'^{2*}_{10}: \{ x \xrightarrow{W_o} z_7 \},$$

$$z'_{11}: \{ ) \xrightarrow{W_o} z_{10} \}.$$

Графически эта же грамматика изображена на рис. I.



$\tau^T$	$a$	$b$	$\dots$	$y$	$z$	$A$	$B$	$\dots$	$Y$	$Z$	$0$	$1$	$\dots$	$9$	$+$	$-$	$\times$	$/$	$\div$	$($	$)$
I	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$u$	$u$	$\dots$	$u$	$+$	$+$	$x$	$x$	$x$	$)$	$($
2	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$\delta$	$\delta$	$\dots$	$\delta$	$\delta$	$u$	$u$	$\dots$	$u$	$x$	$x$	$x$	$x$	$x$	$)$	$($
3	$\delta u$	$\delta u$	$\dots$	$\delta u$	$\delta u$	$\delta u$	$\delta u$	$\dots$	$\delta u$	$\delta u$	$\delta u$	$\delta u$	$\dots$	$\delta u$	$x$	$x$	$x$	$x$	$x$	$)$	$($

Рис. I.  
Граф фрагмента простого арифметического выражения языка АЛГОЛ-60.

2<sup>o</sup>. Рассмотрим некоторые свойства  $R$  - грамматики. Пусть для простоты  $S = T$ . Цепочка символов  $\mathcal{Y} = a_1, a_2, \dots, a_n$  является  $W$  - выводом  $\mathcal{Y}$  (обозначение  $a_i \Rightarrow \mathcal{Y}$ ), если существует последовательность правил преемников, порождающих эту цепочку, начиная от  $a_i$ . Если при этом  $a_k$  порождается правилом вида (1) с символом  $\phi$  в правой части или правилом вида (2), то  $W$  - вывод  $\mathcal{Y}$  называется **законченным** и обозначается  $a_i \Rightarrow^* \mathcal{Y}$ .

Язык, порождаемый  $R$  - грамматикой (3):

$$\mathcal{L}(G) = \{\mathcal{Y} | \mathcal{Y} = a_1, \dots, a_n; a_i \in T, 1 \leq i \leq n, a_i \in \Sigma, a_i \xrightarrow{W} \mathcal{Y}\}.$$

Терминальные символы  $a_j$  и  $a_{j+k}$  ( $k > 0$ ) в цепочке  $\mathcal{Y} = (a_1, a_2, \dots, a_n)$ , порождаемые правилами с  $W_1 = (\mathcal{Y}_1, \dots, \mathcal{Y}_k)$  и  $W_2 = (\mathcal{Y}_{k+1}, \dots, \mathcal{Y}_n)$ , соответственно, называются базовым начальным и базовым конечным для  $i$ -го стека, если  $y_i = y'_i, 1 \leq i \leq n$ .

Правило  $R$  - грамматики называется **фактивным**, если не существует законченного  $W$  - вывода  $\mathcal{Y}$ , в котором данное правило применялось по крайней мере один раз. Грамматика, не содержащая фактивных правил, называется **приведенной**.  $R$  - грамматика является **детерминированной**, если терминальные символы каждого множества правил совпадают лишь в правилах, содержащих  $W_2$  с различными наборами аргументов.

Выделим интересные для практических приложений классы детерминированных  $R$  - грамматик:

1. **Монарные**, мощность всех множеств правил которых равна 1;

2. **односторонние линейные**, для всех правил которых  $m = 0$ ;

3. **левоконтекстные**, для всех правил которых  $m \neq 2$ ;

4. **с ранним определением синтерма**, у которых все подмножества любого из множеств  $\tau_i$  правил имеют одинаковый признак  $\mathcal{C}$ ;

5. **с поздним определением синтерма**, у которых хотя бы два подмножества любого из множеств  $\tau_i$  правил имеют различные признаки  $\mathcal{C}$ ;

6. **последовательные**, у которых элементы множества  $R$  могут быть упорядочены так, что при любом порож-

дении преемников может быть либо правило из этого же множества, либо из любого следующего в упорядоченной последовательности множеств правил.

Монарные грамматики допускают построение универсального алгоритма исправления одиночных ошибок, однако они порождают весьма узкий класс языков.

Односторонние линейные  $R$  - грамматики позволяют проводить синтаксический анализ без стековой памяти.

Левоконтекстные  $R$  - грамматики допускают построение эффективного алгоритма локализации синтаксических ошибок и позволяют компактно описать довольно широкий класс языков.

$R$  - грамматики с ранним определением синтерма позволяют строить более быстрые схемы синтаксического анализа по сравнению с  $R$  - грамматиками с поздним определением синтерма.

Последовательные  $R$  - грамматики допускают весьма простые технические средства для реализации синтаксического анализа.

3. Пусть  $\mathcal{U}$  - множество информатив, некоторое множество стандартных подпрограмм и таблиц, однозначно описывающих некоторый язык  $L_i$ ,  $A$  - множество правил  $R$  - грамматики, задающей язык  $L_{i-1}$ ;  $D$  - множество директив, задающих отображение множества  $A$  в  $\mathcal{U}$ . Тогда под семантическим описанием правил  $R$  - грамматики будем понимать некоторое множество  $C \subseteq (D \times \mathcal{U})$ , элементы которого поставлены во взаимно однозначное соответствие элементам множества  $A$ . Под полным описанием  $R$  - грамматики языка понимается восьмерка

$$G = (\mathcal{C}, S, M, D, R, \tau_0, C),$$

каждое правило которой имеет вид:

$$a \xrightarrow{W_m} \tau_1(d, u) \quad (4)$$

или

$$a \xrightarrow{W_m} \tau_1(d, u), \quad (5)$$

где  $d \in D$ ,  $u \in \mathcal{U}$  из  $C$ .

Приведенная  $R$  - грамматика весьма эффективна при распознавании в синтаксически управляемых устройствах. Действительно, последовательность символов является предложением языка, если синтерм, соответствующий ее первому символу, является символом, встречающимся в левой части одного из правил - преемников примененного правила, а последний символ встречается в одном из

правил вида (5) или в правиле вида (4) с символом  $\phi$  в правой части. На каждом этапе распознавания текущему символу исходной последовательности, или текущему символу, и некоторой последовательности предшествующих символов по правилам, предписанным в  $d$ , ставится в соответствие последовательность символов из  $u$ , где  $d$  и  $u$  - директивное и информативное описание семантики правила  $R$  - грамматики, применяемого к текущему символу исходной последовательности.

Транслятор строится таким образом, что каждому его блоку ставится в соответствие как некоторый параметр  $R$  - грамматика входного для этого блока языка. Блоки работают последовательно. Каждый символ исходного предложения последовательно "проходит" через блоки транслятора до блока, в котором либо заканчивается его обработка, либо происходит накопление для последующих этапов трансляции. Затем выбирается следующий символ исходного предложения (управление передается первому блоку) и т. д. Все блоки на всех этапах трансляции работают по одной и той же программе синтаксического анализа, которая является основной программой всего транслятора. По окончании работы  $i$ -го блока, и включении ( $i + 1$ )-го или первого в  $i$ -м блоке формируется управляющее слово о состоянии программы синтаксического анализа в момент включения. Включение ( $i + 1$ ) -го блока - это настройка общей программы синтаксического анализа с помощью соответствующего этому блоку управляющего слова. Каждый блок, кроме синтаксического анализа, может доопределить текущий символ семантически. Благодаря этому сокращаются  $R$  - грамматики последующих блоков, упрощаются соответствующие им информативные и директивные стандартные подпрограммы.

При многопроходной схеме трансляция ведется последовательно, несколькими итеративно работающими блоками. Важно отметить, что и в этом случае все блоки при всех проходах работают по одной и той же программе синтаксического анализа.

Описанный метаязык позволяет весьма компактно представлять грамматики существующих алгоритмических языков. Так, например, грамматика языка АЛГОЛ - 60 задается приблизительно 250 правилами (машинными словами). Программа синтаксического анализа, реализующая распознавание правил  $R$  - грамматики в процессе работы синтаксически управляемого транслятора, содержит около 100

команд и работает со скоростью нескольких десятков команд на каждый символ исходного предложения.

## ЛИТЕРАТУРА

1. И.В.ВЕЛЬБИЦКИЙ, Ю.Л.ЮЩЕНКО. Метаязык, ориентированный для синтаксического анализа и контроля. Кибернетика, № 2, 1970 .
2. И.В.ВЕЛЬБИЦКИЙ. Первичные средства подготовки синтаксически правильных программ-Труды II Всесоюзной конференции по ВС., М., 1969.
3. И.БЕРЖ. Теория графов и её применение. Изд-во Иностранной литературы, М., 1962 .
4. Алгоритмический язык АЛГОЛ-60 (пересмотренное сообщение). М., Изд-во "Мир", 1965 .

Поступила в редакцию  
20 июня 1970 г.