

Лекция 4

# Синтаксический анализ

## §17. Постановка задачи синтаксического анализа

**Определение.** Порождающую грамматику  $G = \langle N, T, P, S \rangle$  называют *контекстно-свободной* (КС-грамматикой), если каждое правило вывода из множества  $P$  имеет вид

$$A \rightarrow \gamma,$$

где  $A \in N$  – нетерминал, а  $\gamma \in (N \cup T)^*$  – цепочка в объединённом алфавите грамматики.

Будем говорить, что  $uxv$  выводится за один шаг из  $uAv$  (и записывать это как  $uAv \Rightarrow uxv$ ), если  $A \rightarrow x$  – правило вывода, и  $u$  и  $v$  – произвольные цепочки из  $(N \cup T)^*$ .

Если  $u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_n$ , будем говорить, что из  $u_1$  выводится  $u_n$ , и записывать это как  $u_1 \Rightarrow^* u_n$ .

Для отношения  $\Rightarrow^*$  справедливы утверждения:

1.  $u \Rightarrow^* u$  для любой цепочки  $u$ ;
2. если  $u \Rightarrow^* v$  и  $v \Rightarrow^* w$ , то  $u \Rightarrow^* w$ .

Аналогично,  $\Rightarrow^+$  означает «выводится за один или более шагов».

**Определение.** Цепочка  $u$  называется *сентенциальной формой* в грамматике  $G = \langle N, T, P, S \rangle$ , если  $S \Rightarrow^* u$ .

**Определение.** Сентенциальная форма  $u$  называется *левой сентенциальной формой*, если на каждом шаге вывода  $u$  из аксиомы грамматики  $S$  осуществляется подстановка самого левого нетерминала. При этом такой вывод называется *левосторонним*.

Аналогично, *правая сентенциальная форма* и *правосторонний вывод*.

**Определение.** *Предложение* – это сентенциальная форма, не содержащая нетерминалов.

**Определение.** Упорядоченным графом называется пара  $\langle V, E \rangle$ , где  $V$  обозначает множество вершин, а  $E$  – множество линейно упорядоченных списков дуг, каждый элемент которого имеет вид  $\langle \langle v, e_1 \rangle, \langle v, e_2 \rangle, \dots, \langle v, e_n \rangle \rangle$ .

Этот элемент указывает, что из вершины  $v$  выходят  $n$  дуг, причем первой из них считается дуга, входящая в вершину  $e_1$ , второй – дуга, входящая в вершину  $e_2$ , и т.д.

**Определение.** Дерево вывода в КС-грамматике  $G = \langle N, T, P, S \rangle$  – это упорядоченное дерево, каждая вершина которого помечена символом из множества  $N \cup T \cup \{\varepsilon\}$ . При этом:

1. корень дерева помечен аксиомой  $S$ ;
2. каждый лист помечен либо терминалом, либо  $\varepsilon$ ;
3. каждая внутренняя вершина помечена нетерминалом;
4. если  $N$  – нетерминал, которым помечена нелистовая вершина, и  $X_1, \dots, X_n$  – метки её прямых потомков в указанном порядке, то  $N \rightarrow X_1 \dots X_n$  – правило из множества  $P$ .

**Задача синтаксического анализа.** Дано:

$G = \langle N, T, P, S \rangle$  – КС-грамматика языка;

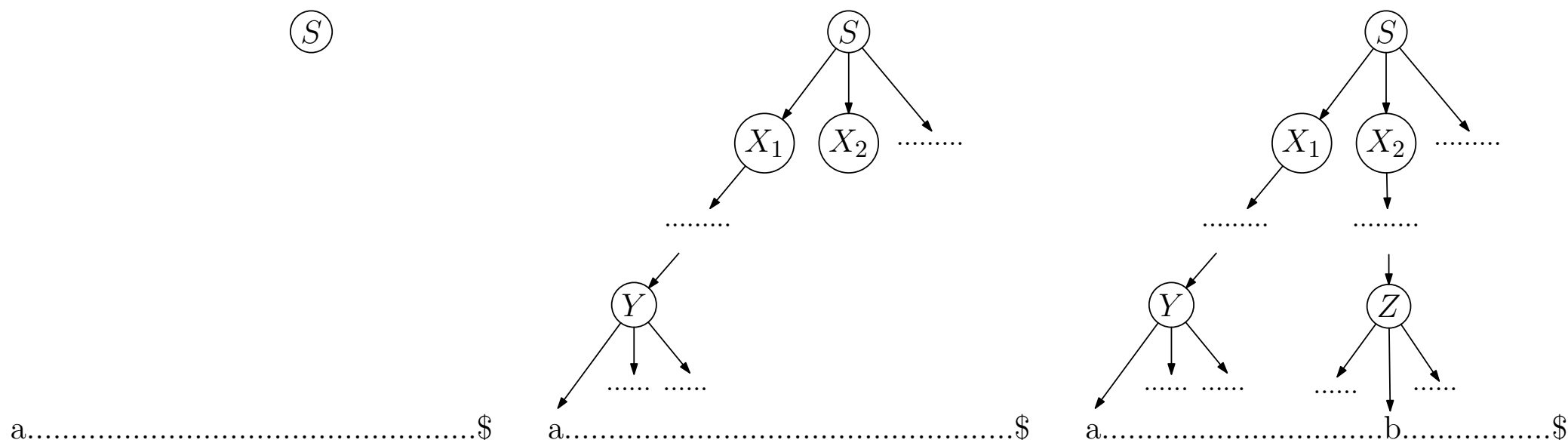
$x \in T^*$  – цепочка.

Требуется определить, является ли цепочка  $x$  предложением в грамматике  $G$ . В случае положительного ответа на этот вопрос следует построить дерево вывода цепочки  $x$  в грамматике  $G$ .

**Замечание.** Следует понимать, что реальный синтаксический анализатор вовсе не обязан строить настоящее дерево вывода. Достаточно, чтобы он в процессе работы формировал некие данные, по которым можно построить дерево вывода.

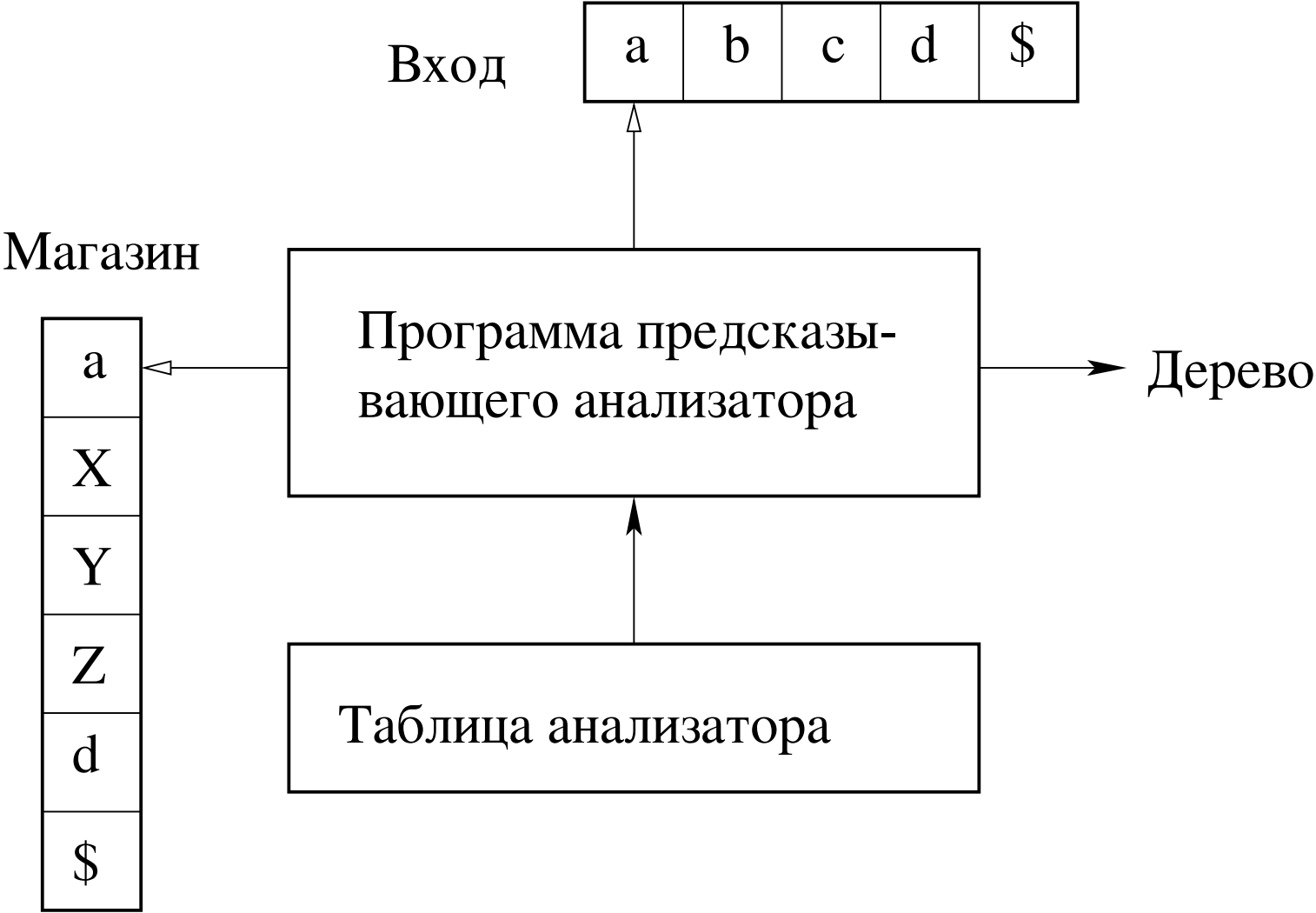
## §18. Алгоритм предсказывающего синтаксического разбора

Последовательность формирования дерева вывода в процессе предсказывающего синтаксического разбора (сверху-вниз):



Основная проблема предсказывающего разбора – определение правила вывода, которое нужно применить к нетерминалу.

Структура предсказывающего анализатора:



Предсказывающий анализатор для грамматики  $G = \langle N, T, P, S \rangle$  получает на вход цепочку  $x \in T^*$ , в конце которой обязательно стоит правый концевой маркер  $\$$ .

Магазин содержит последовательность символов  $y \in (N \cup T)^*$  с символом  $\$$  на дне. В начале магазин содержит аксиому грамматики  $S$  на верхушке и  $\$$  на дне.

Пока будем считать, что на выходе анализатора получается последовательность применённых правил вывода.

Таблица анализатора задаёт функцию  $\delta : N \times T \longrightarrow (N \cup T)^* \cup \{\text{ERROR}\}$ .



Программа предсказывающего анализатора работает следующим образом. Она рассматривает  $X$  – символ на верхушке магазина и  $a$  – текущий входной символ. При этом имеются три возможности.

1. Если  $X = a = \$$ , анализатор останавливается и сообщает об успешном окончании разбора.
2. Если  $X = a \neq \$$ , анализатор удаляет  $X$  из магазина и продвигает указатель входа на следующий входной символ.
3. Если  $X$  – нетерминал, программа заглядывает в таблицу  $\delta(X, a)$ , где хранится либо правило для  $X$ , либо ошибка. Если, например,  $\delta(X, a) = UVW$ , анализатор заменяет  $X$  на верхушке магазина на  $WVU$  ( $U$  оказывается на верхушке) и отправляет на выход правило  $X \rightarrow UVW$ . Если  $\delta(X, a) = \text{ERROR}$ , анализатор переходит в режим восстановления при синтаксических ошибках.

```

TopDownParse( $N$ ,  $T$  (содержит  $\$$ ),  $S$ ,  $\delta$ , цепочка  $x$ ) {
     $Result$  = пустая последовательность;
    поместить в магазин  $S\$$ ;  $a$  = первый символ из  $x$ ;
    do {
         $X$  = верхний символ магазина;
        if ( $X \in T$ ) {
            if ( $X = a$ ) {
                удалить  $X$  из магазина;
                 $a$  = очередной символ из  $x$ ;
            } else error();
        } else {
            if ( $\delta(X, a) = Y_1Y_2...Y_k$ ) {
                удалить  $X$  из магазина;
                поместить  $Y_k...Y_2Y_1$  в магазин ( $Y_1$  сверху);
                добавить правило  $X \rightarrow Y_1Y_2...Y_k$  в  $Result$ ;
            } else error();
        }
    }
    while ( $X \neq \$$ );
    return  $Result$ ;
}

```

## Пример. Объектные Рефал-выражения

$T = \{\text{symbol}, '(', ')', \$\}$ ,  $N = \{\text{expr}, \text{term}\}$ ,  $S = \text{expr}$ ,  
 $P =$

$\text{expr} ::= \text{term expr} \mid \varepsilon.$

$\text{term} ::= \text{symbol} \mid '(\text{expr})'.$

Функция  $\delta$ :

	symbol	'('	)'	\$
expr	term expr	term expr	$\varepsilon$	$\varepsilon$
term	symbol	'(' expr ')'	ERROR	ERROR

## §19. Множества FIRST и FOLLOW

Пусть дана грамматика  $G = \langle N, T, P, S \rangle$ . Будем обозначать нетерминальные символы грамматики буквами  $X$  и  $Y$ , терминальные символы – буквой  $a$ , а цепочки – буквами  $u$  и  $v$ .

Множества FIRST и FOLLOW связаны с грамматикой языка и позволяют построить таблицу предсказывающего разбора.

Определим  $\text{FIRST}(u)$  как множество терминалов, с которых начинаются цепочки, выводимые из  $u$ . Если  $u \Rightarrow^* \varepsilon$ , то  $\varepsilon$  также принадлежит  $\text{FIRST}(u)$ .

$$\text{FIRST}(u) = \{a \in T \mid u \Rightarrow^* av\} \cup \{\varepsilon \mid u \Rightarrow^* \varepsilon\}$$

Определим  $\text{FOLLOW}(X)$  как множество терминалов  $a$  таких, что существует вывод вида  $S \Rightarrow^* uXav$ . Если  $X$  может быть самым правым символом некоторой сентенциальной формы, то  $\$$  принадлежит  $\text{FOLLOW}(X)$ .

$$\text{FOLLOW}(X) = \{a \in T \mid S\$ \Rightarrow^* uXav\}$$

Определим функцию  $\mathcal{F} : (N \cup T)^* \longrightarrow 2^T$ , работающую в предположении, что для каждого нетерминального символа грамматики определено некоторое множество FIRST, и вычисляющую множество FIRST для любой цепочки, то есть  $\text{FIRST}(u) = \mathcal{F}[[u]]$ .

$$\begin{aligned}\mathcal{F}[[au]] &= \{a\} \\ \mathcal{F}[[Xu]] &= \begin{cases} \text{FIRST}(X), & \text{если } \varepsilon \notin \text{FIRST}(X) \\ (\text{FIRST}(X) \setminus \{\varepsilon\}) \cup \mathcal{F}[[u]], & \text{если } \varepsilon \in \text{FIRST}(X) \end{cases} \\ \mathcal{F}[[\varepsilon]] &= \{\varepsilon\}\end{aligned}$$

Запишем алгоритм построения множеств FIRST для нетерминальных символов:

Шаг 1. Пусть  $\text{FIRST}(X) = \emptyset$  для любого  $X \in N$ .

Шаг 2. Для каждого правила  $X \rightarrow u$  добавить  $\mathcal{F}[[u]]$  в  $\text{FIRST}(X)$ .

Шаг 3. Повторять шаг 2 до тех пор, пока множества FIRST не перестанут изменяться.

Запишем алгоритм построения множеств FOLLOW для всех нетерминальных символов:

Шаг 1. Положить  $\text{FOLLOW}(X) = \emptyset$  для любого  $X \in N$ .

Шаг 2. Поместить \$ в  $\text{FOLLOW}(S)$ , где  $S$  – аксиома.

Шаг 3. Если есть правило вывода  $X \rightarrow uYv$ , то все из  $\text{FIRST}(v)$ , за исключением  $\varepsilon$ , добавить к  $\text{FOLLOW}(Y)$ .

Шаг 4. Пока ничего нельзя будет добавить ни к какому множеству  $\text{FOLLOW}(X)$ : если есть правило вывода  $X \rightarrow uY$  или  $X \rightarrow uYv$ , где  $\text{FIRST}(v)$  содержит  $\varepsilon$  (т.е.  $v \Rightarrow^* \varepsilon$ ), то все символы из  $\text{FOLLOW}(X)$  добавить к  $\text{FOLLOW}(Y)$ .

## Пример. Арифметические выражения

$T = \{ '+', '*', n, '(', ')', \$ \}$ ,  $N = \{ E, E', T, T', F \}$ ,  $S = E$ ,  
 $P =$

$E ::= T E'.$

$E' ::= '+' T E' \mid \varepsilon.$

$T ::= F T'.$

$T' ::= '*' F T' \mid \varepsilon.$

$F ::= n \mid '(' E ')'$ .

Множества FIRST и FOLLOW:

$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ '(', n \}$

$\text{FIRST}(E') = \{ '+', \varepsilon \}$

$\text{FIRST}(T') = \{ '*', \varepsilon \}$

$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{ ')', \$ \}$

$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{ '+', ')', \$ \}$

$\text{FOLLOW}(F) = \{ '+', '*', ')', \$ \}$

## §20. LL(1)-грамматики

**Определение.** Грамматика  $G$  является LL(1) тогда и только тогда, когда для любых двух правил вида  $A \rightarrow u \mid v$  выполняется следующее:

1.  $\text{FIRST}(u) \cap \text{FIRST}(v) = \emptyset$ ;

*(Поясняющий пример 1. Предположим, что  $a \in \text{FIRST}(u) \cap \text{FIRST}(v)$ . Тогда невозможно решить, по какому правилу раскрывать  $A$ .)*

*(Поясняющий пример 2. Предположим, что  $\varepsilon \in \text{FIRST}(u) \cap \text{FIRST}(v)$  и  $a \in \text{FOLLOW}(A)$ . Тогда  $A$  нужно раскрыть, как  $\varepsilon$ , но непонятно, каким из двух правил при этом нужно воспользоваться.)*

2. если  $v \Rightarrow^* \varepsilon$ , то  $\text{FIRST}(u) \cap \text{FOLLOW}(A) = \emptyset$ .

*(Поясняющий пример. Предположим, что  $a \in \text{FIRST}(u) \cap \text{FOLLOW}(A)$ . Тогда невозможно определить, то ли  $A$  надо раскрыть как  $u$ , то ли  $A$  раскрывается как  $\varepsilon$ , и символ  $a$  расположен после  $A$ .)*

*(Примечание. В поясняющих примерах  $A$  находится на вершущке магазина,  $a$  – текущий входной символ.)*



Напомним, что программа предсказывающего разбора работает на основе таблицы, задаваемой функцией  $\delta : N \times T \longrightarrow (N \cup T)^* \cup \{\text{ERROR}\}$ . Такая таблица может быть построена только для LL(1)-грамматик.

Основная трудность при использовании предсказывающего анализа – это составление для входного языка LL(1)-грамматики.

Грамматики языков программирования часто бывают «почти LL(1)», то есть отличаются от LL(1)-грамматики тем, что содержат левую рекурсию и нуждаются в левой факторизации.

## Удаление левой рекурсии

**Определение.** КС-грамматика *леворекурсивна*, если в ней имеется нетерминал  $A$  такой, что существует вывод  $A \Rightarrow^+ Au$ , где  $u$  – некоторая цепочка.

Мы не будем рассматривать алгоритм удаления левой рекурсии, ограничившись рассмотрением простейшего случая.

Пусть в грамматике есть правила  $A \rightarrow Av_1 \mid Av_2 \mid \dots \mid Av_n \mid u_1 \mid u_2 \mid \dots \mid u_m$ .

Тогда для удаления левой рекурсии выполним преобразование:

$$A \rightarrow u_1 A' \mid u_2 A' \mid \dots \mid u_m A'.$$

$$A' \rightarrow v_1 A' \mid v_2 A' \mid \dots \mid v_n A' \mid \varepsilon.$$

## Левая факторизация

Пусть в грамматике имеются два правила для нетерминала  $A$ :

$$A \rightarrow uv_1 \mid uv_2.$$

Основная идея левой факторизации заключается в том, что, когда неясно, какую из двух альтернатив надо использовать для развертки нетерминала  $A$ , нужно переделать правила для  $A$  так, чтобы отложить решение до тех пор, пока не будет достаточно информации, чтобы принять правильное решение.

Для левой факторизации надо преобразовать грамматику:

$$A \rightarrow uA'.$$

$$A' \rightarrow v_1 \mid v_2.$$

То есть если входная цепочка начинается с непустой цепочки, выводимой из  $u$ , то в исходной грамматике мы не знаем, разворачивать ли  $A$  по  $uv_1$  или по  $uv_2$ . Однако в факторизованной грамматике можно отложить решение, развернув  $A \rightarrow uA'$ . Тогда после анализа того, что выводимо из  $u$ , можно развернуть  $A' \rightarrow v_1$  или  $A' \rightarrow v_2$ .

## §21. Построение таблиц предсказывающего анализатора

Алгоритм построения таблиц работает для любой КС-грамматики, так как реально строит функцию  $\delta' : N \times T \longrightarrow 2^{(N \cup T)^*} \cup \{\text{ERROR}\}$ . То есть в ячейке таблицы в общем случае могут находиться несколько правил.

### Алгоритм построения таблиц предсказывающего анализатора

Пусть для начала  $\delta'(X, a) = \text{ERROR}$  для любых  $X \in N$  и  $a \in T$ .

Затем проходим по всем правилам грамматики, и для каждого правила  $X \rightarrow u$  выполняем следующее:

- а) перебираем все  $a \in \text{FIRST}(u)$  и добавляем  $u$  к  $\delta'(X, a)$ .
- б) если  $\varepsilon \in \text{FIRST}(u)$ , перебираем все  $b \in \text{FOLLOW}(X)$  и добавляем  $u$  к  $\delta'(X, b)$ .

**Замечание.** Добавление  $u$  к  $\delta'(X, a)$  заключается в том, что если  $\delta'(X, a) = \text{ERROR}$ , то  $\delta'(X, a)$  становится равным  $\{u\}$ , а если  $\delta'(X, a) = A$ , то  $\delta'(X, a)$  становится равным  $A \cup \{u\}$ .

## Пример. Грамматика арифметических выражений

$E ::= T E'.$

$E' ::= '+' T E' \mid \varepsilon.$

$T ::= F T'.$

$T' ::= '*' F T' \mid \varepsilon.$

$F ::= n \mid '(' E ')'$

Инициализируем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
E'	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
T	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
T'	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
F	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR

Рассматриваем правило

$$E ::= T E'$$

Имеем

$$\text{FIRST}(T E') = \{'(', n\},$$

$$\text{FOLLOW}(E) = \{')', \$\}.$$

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	<u>T E'</u>	<u>T E'</u>	ERROR	ERROR
E'	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
T	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
T'	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
F	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR

Рассматриваем правило

$$E' ::= '+' T E'$$

Имеем

$$\text{FIRST}('+' T E') = \{ '+' \},$$

$$\text{FOLLOW}(E') = \{ ')', \$ \}.$$

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	T E'	T E'	ERROR	ERROR
E'	<u>'+' T E'</u>	ERROR	ERROR	ERROR	ERROR	ERROR
T	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
T'	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
F	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR

Рассматриваем правило

$$E' ::= \varepsilon$$

Имеем

$$\text{FIRST}(\varepsilon) = \{\varepsilon\},$$

$$\text{FOLLOW}(E') = \{')', \$\}.$$

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	T E'	T E'	ERROR	ERROR
E'	'+' T E'	ERROR	ERROR	ERROR	$\underline{\varepsilon}$	$\underline{\varepsilon}$
T	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
T'	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
F	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR



Рассматриваем правило

$T ::= F T'$

Имеем

$\text{FIRST}(F T') = \{'(', n\}$ ,

$\text{FOLLOW}(T) = \{'+', ')', \$\}$ .

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	$T E'$	$T E'$	ERROR	ERROR
E'	$'+' T E'$	ERROR	ERROR	ERROR	$\epsilon$	$\epsilon$
T	ERROR	ERROR	<u><math>F T'</math></u>	<u><math>F T'</math></u>	ERROR	ERROR
T'	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR
F	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR

Рассматриваем правило

$T' ::= '*' F T'$

Имеем

$FIRST('*' F T') = \{'*\}$ ,

$FOLLOW(T') = \{'+', ')', '\$'\}$ .

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	$T E'$	$T E'$	ERROR	ERROR
$E'$	$'+' T E'$	ERROR	ERROR	ERROR	$\varepsilon$	$\varepsilon$
T	ERROR	ERROR	$F T'$	$F T'$	ERROR	ERROR
$T'$	ERROR	<u><math>'*' F T'</math></u>	ERROR	ERROR	ERROR	ERROR
F	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR

Рассматриваем правило

$$T' ::= \varepsilon$$

Имеем

$$\text{FIRST}(\varepsilon) = \{\varepsilon\},$$

$$\text{FOLLOW}(T') = \{'+', ')', \$\}.$$

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	$T E'$	$T E'$	ERROR	ERROR
$E'$	$'+' T E'$	ERROR	ERROR	ERROR	$\varepsilon$	$\varepsilon$
T	ERROR	ERROR	$F T'$	$F T'$	ERROR	ERROR
$T'$	$\underline{\varepsilon}$	$'*' F T'$	ERROR	ERROR	$\underline{\varepsilon}$	$\underline{\varepsilon}$
F	ERROR	ERROR	ERROR	ERROR	ERROR	ERROR

Рассматриваем правило

$F ::= n$

Имеем

$\text{FIRST}(n) = \{n\},$

$\text{FOLLOW}(F) = \{' + ', '*' , '\)' , '\$'\}.$

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	T E'	T E'	ERROR	ERROR
E'	'+' T E'	ERROR	ERROR	ERROR	$\varepsilon$	$\varepsilon$
T	ERROR	ERROR	F T'	F T'	ERROR	ERROR
T'	$\varepsilon$	'*' F T'	ERROR	ERROR	$\varepsilon$	$\varepsilon$
F	ERROR	ERROR	<u>n</u>	ERROR	ERROR	ERROR

Рассматриваем правило

$F ::= '(' E ')'$

Имеем

$\text{FIRST}('(' E ')') = \{ '(' \},$

$\text{FOLLOW}(F) = \{ '+', '*', ')', \$ \}.$

Заполняем таблицу:

	'+'	'*'	n	'('	')'	\$
E	ERROR	ERROR	T E'	T E'	ERROR	ERROR
E'	'+' T E'	ERROR	ERROR	ERROR	$\varepsilon$	$\varepsilon$
T	ERROR	ERROR	F T'	F T'	ERROR	ERROR
T'	$\varepsilon$	'*' F T'	ERROR	ERROR	$\varepsilon$	$\varepsilon$
F	ERROR	ERROR	n	<u>'(' E ')'</u>	ERROR	ERROR

## §22. Алгоритм Эрли

---

Earley, J. *An Efficient Context-Free Parsing Algorithm* [PhD Thesis]. – Carnegie-Mellon University, 1968.

---

Алгоритм Эрли подходит для любой КС-грамматики и для заданного предложения языка позволяет найти все левые выводы. Алгоритм Эрли относится к алгоритмам, работающих «сверху-вниз».

Оценка времени:  $O(n)$  – для LL(1)-грамматик,  $O(n^2)$  – для непротиворечивых грамматик,  $O(n^3)$  – для всех остальных КС-грамматик.

Работа алгоритма основана на построении последовательности множеств, называемых *множествами Эрли*. Для входной цепочки  $a_1a_2\dots a_n$  конструируются  $n+1$  множеств: начальное множество  $S_0$  и по одному множеству  $S_i$  для каждого входного символа  $a_i$ .

**Определение.** Элемент множества Эрли (Earley item), принадлежащий множеству  $S_k$  – это тройка  $\langle p, i, j \rangle$ , в которой:  
 $p$  – это правило грамматики;  
 $i$  – позиция в правой части правила, показывающая, какая часть правила уже обработана;  
 $j \leq k$  – номер множества Эрли, в котором началось распознавание по правилу  $p$ .

Элемент  $\langle p, i, j \rangle$  принято записывать в виде  $[A \rightarrow u \bullet v, j]$ , где  $A \rightarrow uv$  – это правило  $p$ , а знак  $\bullet$  задаёт позицию  $i$ .

В начале работы алгоритма Эрли в множество  $S_0$  записывается единственный элемент  $[S' \rightarrow \bullet S, 0]$ , в котором  $S$  – это аксиома грамматики, а  $S' \rightarrow S$  – вспомогательное правило, добавляемое в грамматику для удобства описания алгоритма.

В процессе работы алгоритма Эрли одни элементы порождают другие. Для того чтобы можно было восстановить последовательность применяемых правил, необходимо строить *дерево потомков*, в узлах которого находятся элементы Эрли, а дуги отражают тот факт, что один элемент породил другой элемент.

Алгоритм Эрли работает за  $n + 1$  шагов.

На каждом шаге происходит вычисление множества  $S_i$  и инициализация множества  $S_{i+1}$ :

$S_{i+1} = \emptyset;$

$Q = S_i;$

**loop** {

$Q' = \emptyset;$

**for** (каждый элемент  $x \in Q$ ) {

        SCANNER( $x$ );                   /\* Может добавлять в  $S_{i+1}$  \*/

        PREDICTOR( $x, Q'$ );           /\* Может добавлять в  $Q'$  \*/

        COMPLETER( $x, Q'$ );          /\* Может добавлять в  $Q'$  \*/

    }

**if** ( $Q' \subseteq S_i$ ) **break**;

$S_i = S_i \cup Q';$

$Q = Q';$

}



Пусть  $a_1a_2\dots a_n$  – входная цепочка,  $A$  и  $B$  обозначают нетерминальные символы,  $b$  обозначает терминальный символ,  $u$  обозначает цепочку, составленную из терминальных и нетерминальных символов.

SCANNER( $x$ ). Если  $x = [A \rightarrow \dots \bullet b \dots, j]$  и  $b = a_{i+1}$ ,  
добавить  $[A \rightarrow \dots b \bullet \dots, j]$  в  $S_{i+1}$ .

PREDICTOR( $x, Q$ ). Если  $x = [A \rightarrow \dots \bullet B \dots, j]$ , то для каждого правила вида  $B \rightarrow u$  добавить в  $Q$  элемент  $[B \rightarrow \bullet u, i]$ .  
При этом, если  $B \Rightarrow^* \varepsilon$ , то в  $Q$  также нужно добавить элемент  $[A \rightarrow \dots B \bullet \dots, j]$ .

COMPLETER( $x, Q$ ). Если  $x = [A \rightarrow \dots \bullet, j]$ , то для каждого элемента  $[B \rightarrow \dots \bullet A \dots, k] \in S_j$  добавить в  $Q$  элемент  $[B \rightarrow \dots A \bullet \dots, k]$ .

**Пример.** Грамматика  $E \rightarrow E '+' E \mid n$  и входная цепочка  $n '+' n$ .

$$S_0$$

$Q_1$	$\left\{ \begin{array}{l} S' \rightarrow \bullet E \\ E \rightarrow \bullet E '+' E \\ E \rightarrow \bullet n \end{array} \right.$	$, 0$
-------	---	-------

$n$

$$S_1$$

$Q_1$	$\left\{ \begin{array}{l} E \rightarrow n \bullet \end{array} \right.$	$, 0$
$Q_2$	$\left\{ \begin{array}{l} S' \rightarrow E \bullet \\ E \rightarrow E \bullet '+' E \end{array} \right.$	$, 0$

$'+'$

$$S_2$$

$Q_1$	$\left\{ \begin{array}{l} E \rightarrow E '+' \bullet E \end{array} \right.$	$, 0$
$Q_2$	$\left\{ \begin{array}{l} E \rightarrow \bullet E '+' E \\ E \rightarrow \bullet n \end{array} \right.$	$, 2$

$n$

$$S_3$$

$Q_1$	$\left\{ \begin{array}{l} E \rightarrow n \bullet \end{array} \right.$	$, 2$
$Q_2$	$\left\{ \begin{array}{l} E \rightarrow E '+' E \bullet \\ E \rightarrow E \bullet '+' E \end{array} \right.$	$, 0$
$Q_3$	$\left\{ \begin{array}{l} S' \rightarrow E \bullet \\ E \rightarrow E \bullet '+' E \end{array} \right.$	$, 0$

## §23. Расширенная БНФ

**Определение.** Говорят, что КС-грамматика  $G = \langle N, T, E, S \rangle$  записана в расширенной форме Бэкуса-Наура (РБНФ), если её правила имеют вид  $X \rightarrow R$ , где  $X \in N$  – нетерминальный символ, а  $R$  – дерево, синтаксис которого задаётся грамматикой

$R ::=$	$\varepsilon$	; пустое дерево
	$a$	; терминальный символ
	$X$	; нетерминальный символ
	$RR$	; конкатенация поддеревьев
	$R'   R$	; альтернатива
	$R'^*$	; вхождение 0 и более раз
	$R'^+$	; вхождение 1 и более раз
	$R'^?$	; вхождение 0 или 1 раз

Множество таких синтаксических деревьев мы будем обозначать как RHS.

При текстовой записи правил наивысший приоритет имеют операции '\*', '+', '-', за ними следует конкатенация, а наименьший приоритет имеет альтернатива. При этом круглые скобки служат для изменения приоритета.

**Пример.** Арифметические выражения.

$$\text{Expr} ::= ('+' | '-')^? \text{Term} (('+' | '-') \text{Term})^*.$$
$$\text{Term} ::= \text{Factor} (('*' | '/') \text{Factor})^*.$$
$$\text{Factor} ::= \text{Number} | '(' \text{Expr} ')'.$$

В книгах Н. Вирта используется альтернативная запись РБНФ, называемая синтаксической нотацией Вирта:

Метасимвол	Значение
=	Разделяет левую и правую части правила.
.	Конец правила.
	Альтернатива
[ ]	Вхождение 0 или 1 раз.
{ }	Вхождение 0 или более раз.
( )	Изменение приоритета операций.

Приоритет операций – такой же, как и в РБНФ.

**Пример.** Арифметические выражения.

```

Expr    = [ "+" | " − " ] Term { ( "+" | " − ") Term } .
Term     = Factor { ( "*" | "/" ) Factor } .
Factor  = Number | "(" Expr ")".

```

## §24. Множества FIRST и FOLLOW для РБНФ

Пусть дана грамматика  $G = \langle N, T, E, S \rangle$ , правила которой заданы в РБНФ. Будем обозначать нетерминальные символы грамматики буквами  $X$  и  $Y$ , терминальные символы – буквой  $a$ , а выражения РБНФ – буквами  $u$  и  $v$ .

Определим функцию  $\mathcal{F} : \text{RHS} \longrightarrow 2^T$ , работающую в предположении, что для каждого нетерминального символа грамматики определено некоторое множество FIRST, и вычисляющую множество FIRST для любого выражения РБНФ, то есть  $\text{FIRST}(u) = \mathcal{F}[u]$ .

$$\begin{aligned}\mathcal{F}[a] &= \{a\} \\ \mathcal{F}[X] &= \text{FIRST}(X) \\ \mathcal{F}[uv] &= \begin{cases} \mathcal{F}[u], & \text{если } \varepsilon \notin \mathcal{F}[u] \\ (\mathcal{F}[u] \setminus \{\varepsilon\}) \cup \mathcal{F}[v], & \text{если } \varepsilon \in \mathcal{F}[u] \end{cases} \\ \mathcal{F}[u \mid v] &= \mathcal{F}[u] \cup \mathcal{F}[v] \\ \mathcal{F}[u^*] &= \mathcal{F}[u] \cup \{\varepsilon\} \\ \mathcal{F}[u^+] &= \mathcal{F}[u] \\ \mathcal{F}[u^?] &= \mathcal{F}[u] \cup \{\varepsilon\} \\ \mathcal{F}[\varepsilon] &= \{\varepsilon\}\end{aligned}$$

Алгоритм построения множеств FIRST  
для всех нетерминальных символов грамматики:

Шаг 1. Пусть  $\text{FIRST}(X) = \emptyset$  для любого  $X \in N$ .

Шаг 2. Для каждого правила  $X \rightarrow u$  добавить  $\mathcal{F}[[u]]$  в  $\text{FIRST}(X)$ .

Шаг 3. Повторять шаг 2 до тех пор, пока множества FIRST не перестанут изменяться.

Определим функцию  $\mathcal{H} : N \longrightarrow 2^T$ , работающую в предположении, что для каждого нетерминального символа грамматики вычислено множество FIRST и определено некоторое приближение множества FOLLOW, и вычисляющую новое приближение множества FOLLOW для некоторого  $X$ .

$$\mathcal{H}(X) = \bigcup_{(Y \rightarrow u) \in E} \mathcal{H}'(X, \text{FOLLOW}(Y), \llbracket u \rrbracket)$$

Вспомогательная функция  $\mathcal{H}' : N \times 2^T \times \text{RHS} \longrightarrow 2^T$  определяется как:

$$\begin{aligned} \mathcal{H}'(X, s, \llbracket a \rrbracket) &= \emptyset \\ \mathcal{H}'(X, s, \llbracket Y \rrbracket) &= \begin{cases} s, & \text{если } X = Y \\ \emptyset, & \text{если } X \neq Y \end{cases} \\ \mathcal{H}'(X, s, \llbracket uv \rrbracket) &= \begin{cases} \mathcal{H}'(X, \mathcal{F}(v), \llbracket u \rrbracket) \cup \mathcal{H}'(X, s, \llbracket v \rrbracket), & \text{если } \varepsilon \notin \mathcal{F}(v) \\ \mathcal{H}'(X, (\mathcal{F}(v) \setminus \{\varepsilon\}) \cup s, \llbracket u \rrbracket) \cup \mathcal{H}'(X, s, \llbracket v \rrbracket), & \text{если } \varepsilon \in \mathcal{F}(v) \end{cases} \\ \mathcal{H}'(X, s, \llbracket u | v \rrbracket) &= \mathcal{H}'(X, s, \llbracket u \rrbracket) \cup \mathcal{H}'(X, s, \llbracket v \rrbracket) \\ \mathcal{H}'(X, s, \llbracket u^* \rrbracket) &= \mathcal{H}'(X, (\mathcal{F}(u) \setminus \{\varepsilon\}) \cup s, \llbracket u \rrbracket) \\ \mathcal{H}'(X, s, \llbracket u^+ \rrbracket) &= \mathcal{H}'(X, (\mathcal{F}(u) \setminus \{\varepsilon\}) \cup s, \llbracket u \rrbracket) \\ \mathcal{H}'(X, s, \llbracket u^? \rrbracket) &= \mathcal{H}'(X, s, \llbracket u \rrbracket) \\ \mathcal{H}'(X, s, \llbracket \varepsilon \rrbracket) &= \emptyset \end{aligned}$$



Алгоритм построения множеств FOLLOW  
для всех нетерминальных символов грамматики:

Шаг 1. Пусть  $\text{FOLLOW}(S) = \{\$ \}$  и  $\text{FOLLOW}(X) = \emptyset$  для любого  $X \neq S$ .

Шаг 2.  $\text{FOLLOW}(X) = \mathcal{H}(X)$  для каждого  $X \in N$ .

Шаг 3. Повторять шаг 2 до тех пор, пока множества FOLLOW не перестанут изменяться.

## §25. Рекурсивный спуск

Мы будем рассматривать метод рекурсивного спуска применительно к LL(1)-грамматикам, записанным в РБНФ.

В синтаксическом анализаторе, написанном методом рекурсивного спуска, каждому нетерминалу грамматики соответствует отдельная функция, в которой закодирован эффект применения соответствующего нетерминалу правила (мы будем считать, что такое правило единственно).

Цель этой функции – анализ последовательности токенов, которые по её запросу выдаёт лексический анализатор, и проверка соответствия этой последовательности правилу грамматики:

```
X() {  
    PARSE( правило );  
}
```

Функция, соответствующая нетерминалу  $X$  грамматики, составляется с учётом того, что:

- существует глобальная переменная  $Sym$ , в которую ДО вызова функции помещается токен, соответствующий первому символу цепочки, в которую раскрывается нетерминал  $X$ ;
- функция должна либо полностью потребить последовательность токенов, в которую раскрывается нетерминал  $X$ , либо сгенерировать сообщение об ошибке;
- после завершения функции переменная  $Sym$  должна содержать токен, непосредственно следующий за раскрытым нетерминалом  $X$ .

Код функции, соответствующая нетерминалу  $X$  грамматики, порождается из правила грамматики согласно следующим шаблонам:

1. Пустое дерево

```
PARSE( $\varepsilon$ )  $\Rightarrow$   
    ; /* Ничего не делать. */
```

2. Терминальный символ

```
PARSE( $a$ )  $\Rightarrow$   
    {  
        if ( $Sym == a$ )  
             $Sym = \text{NextToken}()$ ;  
        else  
            ReportError();  
    }
```

### 3. Нетерминальный символ

$$\text{PARSE}(X) \Rightarrow$$
$$X();$$

### 4. Конкатенация поддеревьев

$$\text{PARSE}(R_1 R_2) \Rightarrow$$
$$\{ \text{PARSE}(R_1); \quad \text{PARSE}(R_2); \}$$

### 5. Альтернатива ( $\varepsilon \notin \text{FIRST}(R_1)$ )

$$\text{PARSE}(R_1 | R_2) \Rightarrow$$
$$\{$$
$$\quad \textbf{if } (Sym \in \text{FIRST}(R_1))$$
$$\quad \quad \text{PARSE}(R_1);$$
$$\quad \textbf{else}$$
$$\quad \quad \text{PARSE}(R_2);$$
$$\}$$

6. Вхождение 0 или более раз

$$\text{PARSE}(R^*) \Rightarrow$$
$$\begin{cases} \text{while } (Sym \in \text{FIRST}(R)) \\ \quad \text{PARSE}(R); \end{cases}$$

7. Вхождение 1 или более раз

$$\text{PARSE}(R^+) \Rightarrow$$
$$\begin{cases} \text{do} \\ \quad \text{PARSE}(R); \\ \text{while } (Sym \in \text{FIRST}(R)) \end{cases}$$

8. Вхождение 0 или 1 раз

```
PARSE( $R^?$ )  $\Rightarrow$   
  {  
    if ( $Sym \in \text{FIRST}(R)$ )  
      PARSE( $R$ );  
  }
```

**Пример.** Арифметические выражения.

1. Имеем правило для нетерминала Expr:

$$\text{Expr} ::= ('+' | '-')^? \text{Term} (('+' | '-') \text{Term})^*.$$

Создаём функцию Expr:

```
Expr() {  
    PARSE( ('+' | '-')^? Term (('+' | '-') Term)^* );  
}
```

2. Правило для Expr представляет собой конкатенацию трёх поддеревьев:

```
Expr() {  
    PARSE( ('+' | '-')^? );  
    PARSE( Term );  
    PARSE( (('+' | '-') Term)^* );  
}
```



3. Раскрываем  $\text{PARSE} ( ('+' | '-')^? )$ :

```
Expr ( ) {  
    if ( Sym ∈ FIRST ( '+' | '-' ) )  
        PARSE ( '+' | '-' );  
    PARSE ( Term );  
    PARSE ( ( ('+' | '-') Term )* );  
}
```

4. Раскрываем  $\text{PARSE} ( \text{Term} )$ :

```
Expr ( ) {  
    if ( Sym ∈ FIRST ( '+' | '-' ) )  
        PARSE ( '+' | '-' );  
    Term ( );  
    PARSE ( ( ('+' | '-') Term )* );  
}
```

5. Раскрываем  $\text{PARSE}((\text{'+'} \mid \text{'-'}) \text{Term})^*$  ):

```
Expr ( ) {  
    if ( Sym  $\in$  FIRST(  $\text{'+'} \mid \text{'-'} \text{'}$  ) )  
        PARSE(  $\text{'+'} \mid \text{'-'} \text{'}$  );  
    Term ( );  
    while ( Sym  $\in$  FIRST(  $(\text{'+'} \mid \text{'-'}) \text{Term}$  ) )  
        PARSE(  $(\text{'+'} \mid \text{'-'}) \text{Term}$  );  
}
```

6. Раскрываем  $\text{PARSE}( '+' | '-' )$ :

```
Expr ( ) {  
    if (  $Sym \in \text{FIRST}( '+' | '-' )$  ) {  
        if (  $Sym \in \text{FIRST}( '+' )$  )  
            PARSE( '+' );  
        else  
            PARSE( '-' );  
    }  
    Term ( );  
    while (  $Sym \in \text{FIRST}( ('+' | '-') \text{Term} )$  )  
        PARSE( ('+' | '-') Term );  
}
```

7. Раскрываем  $\text{PARSE}( '+' )$  и  $\text{PARSE}( '-' )$ :

```
Expr ( ) {  
    if ( Sym  $\in$  FIRST( '+' | '-' ) ) {  
        if ( Sym  $\in$  FIRST( '+' ) ) {  
            if ( Sym == '+' ) Sym = NextToken ( );  
            else ReportError ( );  
        } else {  
            if ( Sym == '-' ) Sym = NextToken ( );  
            else ReportError ( );  
        }  
    }  
    Term ( );  
    while ( Sym  $\in$  FIRST( ( '+' | '-' ) Term ) )  
        PARSE( ( '+' | '-' ) Term );  
}
```

8. Упрощаем функцию:

```
Expr ( ) {  
    if ( Sym == '+' || Sym == '-' )  
        Sym = NextToken ( ) ;  
    Term ( ) ;  
    while ( Sym ∈ FIRST( ('+' | '-') Term ) )  
        PARSE( ('+' | '-') Term ) ;  
}
```

9. Раскрываем  $\text{PARSE} ( '+' | '-' ) \text{Term}$  :

```
Expr ( ) {  
    if ( Sym == '+' || Sym == '-' )  
        Sym = NextToken ( );  
    Term ( );  
    while ( Sym ∈ FIRST ( '+' | '-' ) Term ) {  
        PARSE ( '+' | '-' );  
        PARSE ( Term );  
    }  
}
```

10. Раскрываем `PARSE( '+' | '-' )`:

```
Expr() {  
    if (Sym == '+' || Sym == '-')  
        Sym = NextToken();  
    Term();  
    while (Sym ∈ FIRST( '+' | '-' ) Term) {  
        if (Sym ∈ FIRST( '+' | '-' )) {  
            if (Sym ∈ FIRST( '+' )) {  
                if (Sym == '+') Sym = NextToken();  
                else ReportError();  
            } else {  
                if (Sym == '-') Sym = NextToken();  
                else ReportError();  
            }  
        }  
        PARSE( Term );  
    }  
}
```

11. Упрощаем функцию и раскрываем PARSE( Term ):

```
Expr ( ) {  
    if ( Sym == '+' || Sym == '-' )  
        Sym = NextToken ( );  
    Term ( );  
    while ( Sym == '+' || Sym == '-' ) {  
        Sym = NextToken ( );  
        Term ( );  
    }  
}
```



12. Имеем правило для нетерминала Term:

$\text{Term} ::= \text{Factor} (('*' | '/') \text{Factor})^*.$

Создаём функцию Term:

```
Term() {  
    PARSE( Factor (('*' | '/') Factor)* );  
}
```

13. После раскрытия  $\text{PARSE}( \text{Factor} (('*' | '/') \text{Factor})^* )$  и упрощения получаем:

```
Term() {  
    Factor();  
    while (Sym == '*' || Sym == '/') {  
        Sym = NextToken();  
        Factor();  
    }  
}
```

14. Имеем правило для нетерминала Factor:

Factor ::= Number | '(' Expr ')'

Создаём функцию Factor:

```
Factor() {  
    PARSE( Number | '(' Expr ')' );  
}
```

15. Раскрываем PARSE( Number | '(' Expr ')' ):

```
Factor() {  
    if (Sym ∈ FIRST( Number ))  
        PARSE( Number );  
    else  
        PARSE( '(' Expr ')' );  
}
```

15. Раскрываем `PARSE( Number )`:

```
Factor() {  
    if (Sym ∈ FIRST( Number )) {  
        if (Sym == Number)  
            Sym = NextToken();  
        else  
            ReportError();  
    } else  
        PARSE( '(' Expr ')' );  
}
```

16. Упрощаем функцию:

```
Factor() {  
    if (Sym == Number)  
        Sym = NextToken();  
    else  
        PARSE( '(' Expr ')' );  
}
```

17. Раскрываем PARSE( '(' Expr ')' ):

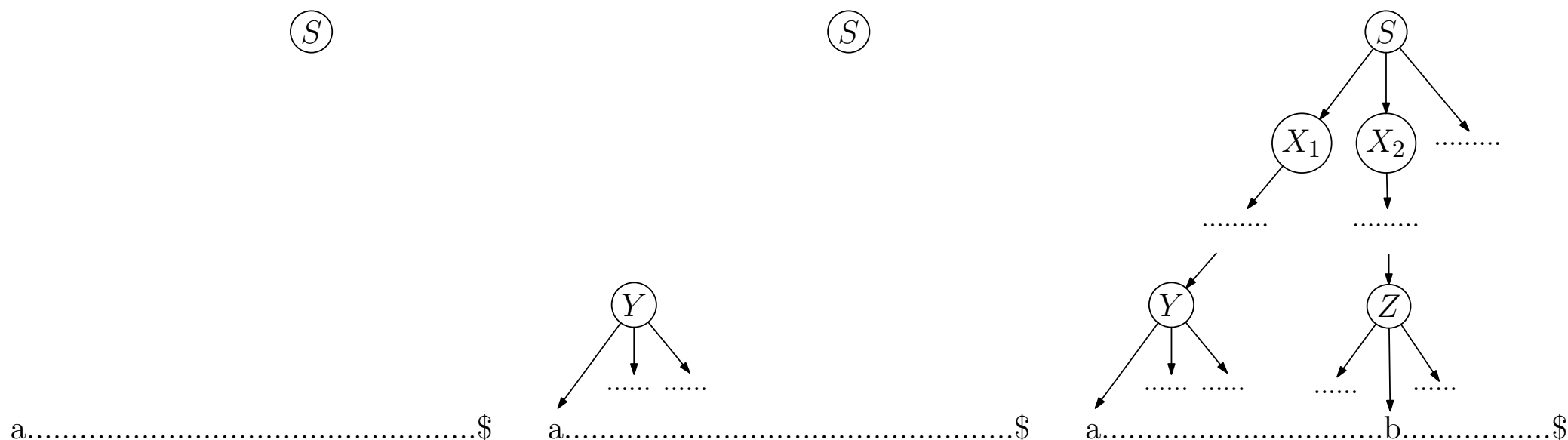
```
Factor() {  
    if (Sym == Number)  
        Sym = NextToken();  
    else {  
        PARSE( '(' );  
        PARSE( Expr );  
        PARSE( ')' );  
    }  
}
```

19. Раскрывая все вхождения PARSE, в итоге получаем:

```
Factor() {  
    if (Sym == Number)  
        Sym = NextToken();  
    else {  
        if (Sym == '(')  
            Sym = NextToken();  
        else  
            ReportError();  
  
        Expr();  
  
        if (Sym == ')')  
            Sym = NextToken();  
        else  
            ReportError();  
    }  
}
```

## §26. Синтаксический разбор типа «перенос–свёртка»

Последовательность формирования дерева вывода в процессе синтаксического разбора типа «перенос–свёртка» (снизу-вверх):



ПС-алгоритм (алгоритм синтаксического разбора типа «перенос–свёртка») строит правосторонний вывод предложения языка в обратном порядке.

**Пример.** Арифметические выражения.

$$T = \{ '+', '*', n, '(', ')', \$ \}, N = \{ E, T, F \}, S = E,$$

$P =$

$$E ::= E \text{ ' + ' } T \mid T.$$

$$T ::= T \text{ ' * ' } F \mid F.$$

$$F ::= n \mid \text{ ' ( ' } E \text{ ' ) ' }.$$

Для цепочки  $n * (n + n)$  ПС-алгоритм строит вывод:

$$\underline{n} * (n + n) \Leftarrow \underline{E} * (n + n) \Leftarrow T * (\underline{n} + n) \Leftarrow T * (\underline{E} + n) \Leftarrow T * (\underline{T} + n) \Leftarrow \\ \Leftarrow T * (E + \underline{n}) \Leftarrow T * (E + \underline{E}) \Leftarrow T * (\underline{E} + \underline{T}) \Leftarrow T * (\underline{E}) \Leftarrow \underline{T} * \underline{F} \Leftarrow \underline{T} \Leftarrow E.$$

**Определение.** Основа правой сентенциальной формы  $u$  – это сочетание правила  $X \rightarrow v$  и позиции цепочки  $v$  в  $u$ , такое, что вхождение  $v$  в  $u$  может быть заменено нетерминалом  $X$  для получения предыдущей правой сентенциальной формы в правостороннем выводе  $u$ .

После основы в правой сентенциальной форме идут только терминальные символы. В примере<sup>↑</sup> все основы подчёркнуты.

ПС-алгоритм работает с двумя структурами данных:

- стек, в который помещаются символы грамматики  $(N \cup T)$ , и дно которого помечено символом \$;
- входной буфер, содержащий нераспознанный суффикс входной цепочки и оканчивающийся символом \$.

В начале работы ПС-алгоритма стек содержит \$, а входной буффер –  $u\$$ , где  $u$  – распознаваемая цепочка.

ПС-алгоритм выполняет четыре действия:

1. *перенос* – символ в начале входного буфера переносится в стек;
2. *свёртка* – основа на вершине стека заменяется нетерминалом;
3. *допуск* – успешное окончание разбора (во входном буфере – \$, в стеке –  $\$S$ , где  $S$  – аксиома грамматики);
4. *ошибка* – вызов подпрограммы реакции на ошибку во входном потоке.

Основные проблемы ПС-алгоритма – поиск основы в стеке и выбор правила, по которому нужно осуществлять свёртку. В общем случае ПС-алгоритм осуществляет возвраты (см. [Ахо, Ульман]).



**Пример.** Разбор цепочки  $n * (n + n)$ .

№	Стек	Вход	Действие
1	\$	$n*(n+n)$$	Перенос
2	$\$n$	$*(n+n)$$	Свёртка по правилу $F \rightarrow n$
3	$\$F$	$*(n+n)$$	Свёртка по правилу $T \rightarrow F$
4!	$\$T$	$*(n+n)$$	Перенос $\times 3$
7	$\$T*(n$	$+n)$$	Свёртка по правилу $F \rightarrow n$
8	$\$T*(F$	$+n)$$	Свёртка по правилу $T \rightarrow F$
9	$\$T*(T$	$+n)$$	Свёртка по правилу $E \rightarrow T$
10	$\$T*(E$	$+n)$$	Перенос $\times 2$
12	$\$T*(E+n$	$)$$	Свёртка по правилу $F \rightarrow n$
13	$\$T*(E+F$	$)$$	Свёртка по правилу $T \rightarrow F$
14!	$\$T*(E+T$	$)$$	Свёртка по правилу $E \rightarrow E '+' T$
15	$\$T*(E$	$)$$	Перенос
16	$\$T*(E)$	$$$	Свёртка по правилу $F \rightarrow '(' E ')'$
17!	$\$T*F$	$$$	Свёртка по правилу $T \rightarrow T '*' F$
18	$\$T$	$$$	Свёртка по правилу $E \rightarrow T$
19	$\$E$	$$$	Допуск

## §27. Недетерминированные SLR-распознаватели

**Определение.** *Активный префикс* (viable prefix) – это префикс право-сентенциальной формы, который не выходит за правый конец основы этой сентенциальной формы.

Если ПС-алгоритм работает правильно (то есть не пошёл по ложному пути), то его стек содержит знак \$, выше которого располагается активный префикс.

**Пример.** На 14-ом шаге разбора цепочки  $n * (n + n)$  (см. пример в §26) в стеке содержится цепочка  
 $\$T^*(\underline{E+T}$   
в которой  $T^*(\underline{E+T})$  является активным префиксом.

**Определение.** *SLR-ситуация* (SLR item) – это пара  $\langle p, i \rangle$ , в которой:

- $p$  – правило грамматики;
- $i$  – позиция в правой части правила, показывающая, какая часть правила уже обработана.

Далее мы будем называть SLR-ситуации просто *ситуациями*.

Ситуацию  $\langle p, i \rangle$  принято записывать в виде  $X \rightarrow u \bullet v$  или  $[X \rightarrow u \bullet v]$ , где  $X \rightarrow uv$  – это правило  $p$ , а знак  $\bullet$  задаёт позицию  $i$ .

Для каждого правила  $X \rightarrow x_1 x_2 \dots x_n$  можно построить  $(n + 1)$  ситуаций:

$X \rightarrow \bullet x_1 x_2 \dots x_n,$   
 $X \rightarrow x_1 \bullet x_2 \dots x_n,$   
 $X \rightarrow x_1 x_2 \bullet \dots x_n,$   
 $\dots,$   
 $X \rightarrow x_1 x_2 \dots \bullet x_n,$   
 $X \rightarrow x_1 x_2 \dots x_n \bullet.$

Для правила вида  $X \rightarrow \varepsilon$  возможна единственная ситуация  $X \rightarrow \bullet$ .

**Определение.** *Расширенная грамматика* для грамматики  $G = \langle N, T, P, S \rangle$  – это грамматика  $G' = \langle N', T, P', S' \rangle$ , в которой  $N' = N \cup \{S'\}$ ,  $P' = P \cup \{S' \rightarrow S\}$ .

Основные проблемы ПС-алгоритма – поиск основы в стеке и выбор правила, по которому нужно осуществлять свёртку – для некоторых КС-грамматик можно решить с помощью распознавателей активных префиксов.

**Определение.** *Недетерминированный SLR-распознаватель активных префиксов* для расширенной грамматики  $G' = \langle N', T, P', S' \rangle$  – это конечный автомат, состояниями которого являются все возможные SLR-ситуации для  $P'$ , а переходы между ними задаются условиями:

- 1)  $[X \rightarrow u \bullet xv] \xrightarrow{x} [X \rightarrow ux \bullet v]$ , где  $x \in N' \cup T$ ;
- 2)  $[X \rightarrow u \bullet Yv] \xrightarrow{\varepsilon} [Y \rightarrow \bullet w]$ , где  $Y \in N'$ .

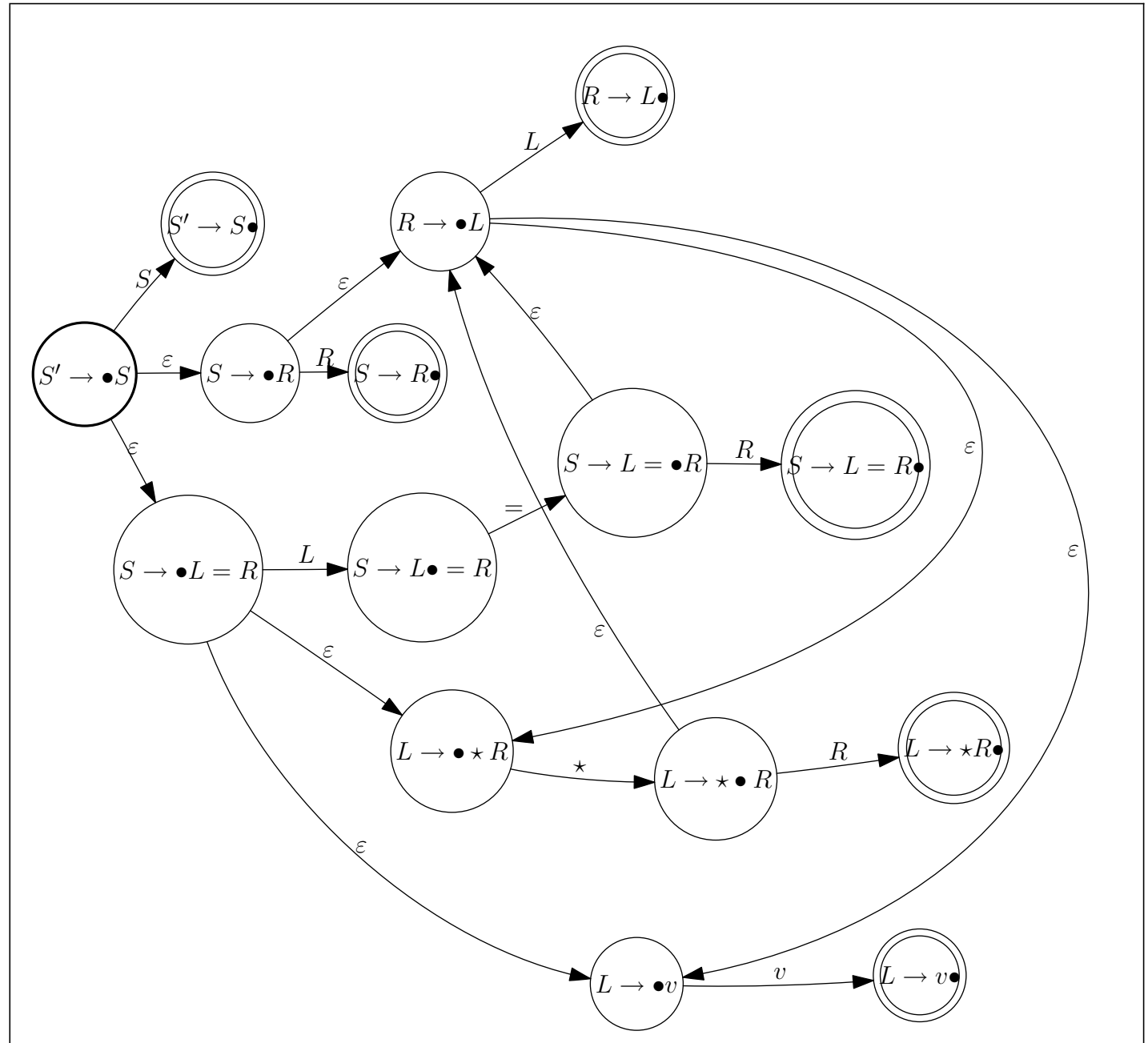
При этом начальным состоянием автомата является ситуация  $[S' \rightarrow S]$ , а заключительными состояниями являются все ситуации вида  $[X \rightarrow u \bullet]$ .

## Пример.

Правила  
грамматики:

$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow L = R \\ S &\rightarrow R \\ L &\rightarrow \star R \\ L &\rightarrow v \\ R &\rightarrow L \end{aligned}$$

Примеры  
выводимых  
предложений:

$$\begin{aligned} &v \\ &\star v \\ &\star \star v \\ &v = v \\ &\star \star v = \star v \end{aligned}$$


Если «натравить» SLR-распознаватель на активный префикс, лежащий в стеке ПС-алгоритма, то в результате мы можем получить некоторую ситуацию, которая является подсказкой для ПС-алгоритма:

- заключительная ситуация вида  $[X \rightarrow u\bullet]$  говорит о том, что нужно выполнить свёртку по правилу  $X \rightarrow u$ ;
- ситуация вида  $[X \rightarrow u \bullet xv]$  означает, что нужно выполнить перенос.

**Пример.** Пусть выполняется синтаксический анализ цепочки  $\star \star v = \star v$ , и на некотором шаге работы ПС-алгоритма на стеке лежит активный префикс  $L = \star R$ .

Распознавание активного префикса:

$$[S' \rightarrow \bullet S] \xrightarrow{\varepsilon} [S \rightarrow \bullet L = R] \xrightarrow{L} [S \rightarrow L\bullet = R] \xrightarrow{=} [S \rightarrow L = \bullet R] \xrightarrow{\varepsilon} [R \rightarrow \bullet L] \xrightarrow{\varepsilon} [L \rightarrow \bullet \star R] \xrightarrow{\star} [L \rightarrow \star \bullet R] \xrightarrow{R} [L \rightarrow \star R\bullet].$$

$[L \rightarrow \star R\bullet]$  – заключительное состояние, означающее свёртку основы  $\star R$  по правилу  $L \rightarrow \star R$ .

Если распознавание активного префикса оказалось невозможным, это означает, что обнаружена ошибка.

**Пример.** Пусть выполняется синтаксический анализ цепочки  $\star = \star v$ , и на некотором шаге работы ПС-алгоритма на стеке лежит активный префикс  $\star =$ .

Распознавание активного префикса:

$$[S' \rightarrow \bullet S] \xrightarrow{\varepsilon} [S \rightarrow \bullet L = R] \xrightarrow{\varepsilon} [L \rightarrow \bullet \star R] \xrightarrow{\star} [L \rightarrow \star \bullet R] \mapsto ?$$

Активный префикс не может быть распознан – обнаружена ошибка.

Если в результате распознавания активного префикса мы получили более одной ситуации, это означает, что КС-грамматика слишком сложна для SLR-распознавателя.

**Пример.** Пусть выполняется синтаксический анализ цепочки  $v = v$ , и на некотором шаге работы ПС-алгоритма на стеке лежит активный префикс  $L$ .

Распознавание активного префикса:

1.  $[S' \rightarrow \bullet S] \xrightarrow{\varepsilon} [S \rightarrow \bullet L = R] \xrightarrow{L} [S \rightarrow L \bullet = R]$  – перенос;
2.  $[S' \rightarrow \bullet S] \xrightarrow{\varepsilon} [S \rightarrow \bullet R] \xrightarrow{\varepsilon} [R \rightarrow \bullet L] \xrightarrow{L} [R \rightarrow L \bullet]$  – свёртка по правилу  $R \rightarrow L$ .

Имеем так называемый конфликт «перенос/свёртка».

Кроме конфликтов «перенос/свёртка» бывают ещё конфликты «перенос/перенос».



## §28. Детерминированные SLR-распознаватели

Если выполнить детерминизацию недетерминированного SLR-распознавателя, то получится *детерминированный SLR-распознаватель активного префикса*, состояниями которого являются не отдельные ситуации, а множества ситуаций.

При этом начальным состоянием детерминированного SLR-распознавателя является состояние, содержащее ситуацию  $[S' \rightarrow \bullet S]$ , а заключительными состояниями являются состояния, содержащие заключительные состояния соответствующего недетерминированного SLR-распознавателя.

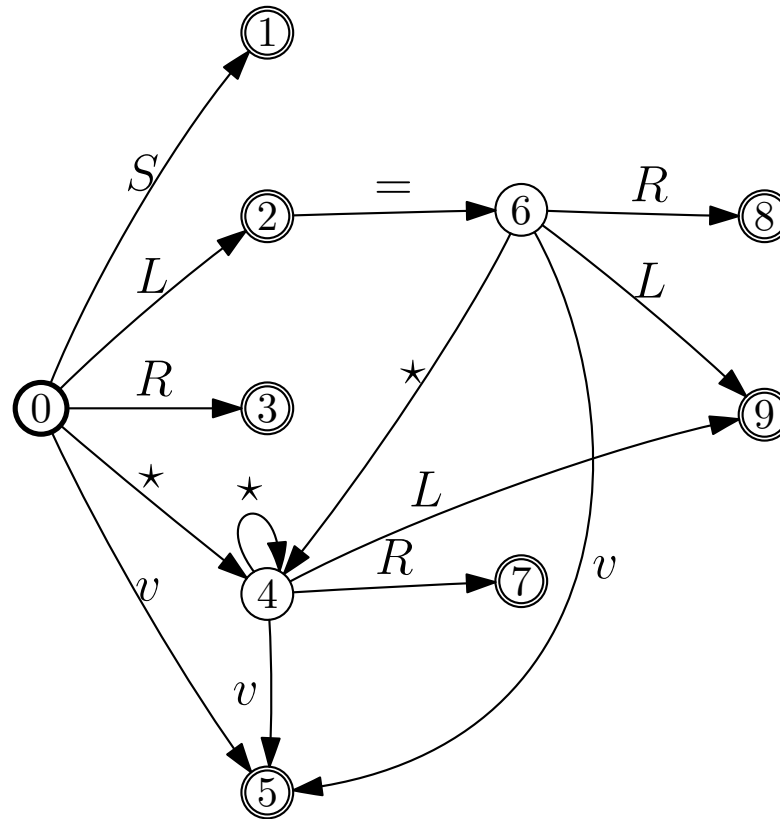
На практике используются только детерминированные SLR-распознаватели, потому что они позволяют распознавать активные префиксы за время, пропорциональное длине префикса.

## Пример.

Правила  
грамматики:

$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow L = R \\ S &\rightarrow R \\ L &\rightarrow \star R \\ L &\rightarrow v \\ R &\rightarrow L \end{aligned}$$

(В состоянии 2 –  
конфликт «перенос/свёртка».)



$$\begin{aligned} 0 : \quad & S' \rightarrow \bullet S \\ & S \rightarrow \bullet L = R \\ & S \rightarrow \bullet R \\ & L \rightarrow \bullet \star R \\ & L \rightarrow \bullet v \\ & R \rightarrow \bullet L \end{aligned}$$

$$\underline{1} : S' \rightarrow S \bullet$$

$$\begin{aligned} \underline{2} : \quad & S \rightarrow L \bullet = R \\ & R \rightarrow L \bullet \end{aligned}$$

$$\underline{3} : S \rightarrow R \bullet$$

$$\begin{aligned} 4 : \quad & L \rightarrow \star \bullet R \\ & R \rightarrow \bullet L \\ & L \rightarrow \bullet \star R \\ & L \rightarrow \bullet v \end{aligned}$$

$$\underline{5} : L \rightarrow v \bullet$$

$$\begin{aligned} 6 : \quad & S \rightarrow L = \bullet R \\ & R \rightarrow \bullet L \\ & L \rightarrow \bullet \star R \\ & L \rightarrow \bullet v \end{aligned}$$

$$\underline{7} : L \rightarrow \star R \bullet$$

$$\underline{8} : S \rightarrow L = R \bullet$$

$$\underline{9} : R \rightarrow L \bullet$$

Детерминированные SLR-распознаватели представляют в программах в виде двух таблиц: action и goto.

Таблица  $\text{action}[s, a]$ , где  $s$  – номер состояния SLR-распознавателя, а  $a \in T \cup \{\$\}$ , задаётся следующим образом:

- $\text{action}[s, a] \ni$  «перенос  $s'$ », если  $s \xrightarrow{a} s'$ ;
- $\text{action}[s, a] \ni$  «свёртка  $X \rightarrow u$ », если  $[X \rightarrow u\bullet] \in s$  и  $a \in \text{FOLLOW}(X)$ ;
- $\text{action}[s, \$] \ni$  «допуск», если  $[S' \rightarrow S\bullet] \in s$ ;
- $\text{action}[s, a] \ni$  «ошибка», если ни одна из вышеприведённых альтернатив не сработала.

Таблица  $\text{goto}[s, X]$ , где  $s$  – номер состояния SLR-распознавателя, а  $x \in N'$ , задаётся следующим образом:

- $\text{goto}[s, X] = s'$ , если  $s \xrightarrow{X} s'$ ;
- $\text{goto}[s, X] =$  «ошибка» в противном случае.

## Пример.

$\text{FOLLOW}(S') = \{\$, \}$ ,

$\text{FOLLOW}(S) = \{\$, \}$ ,

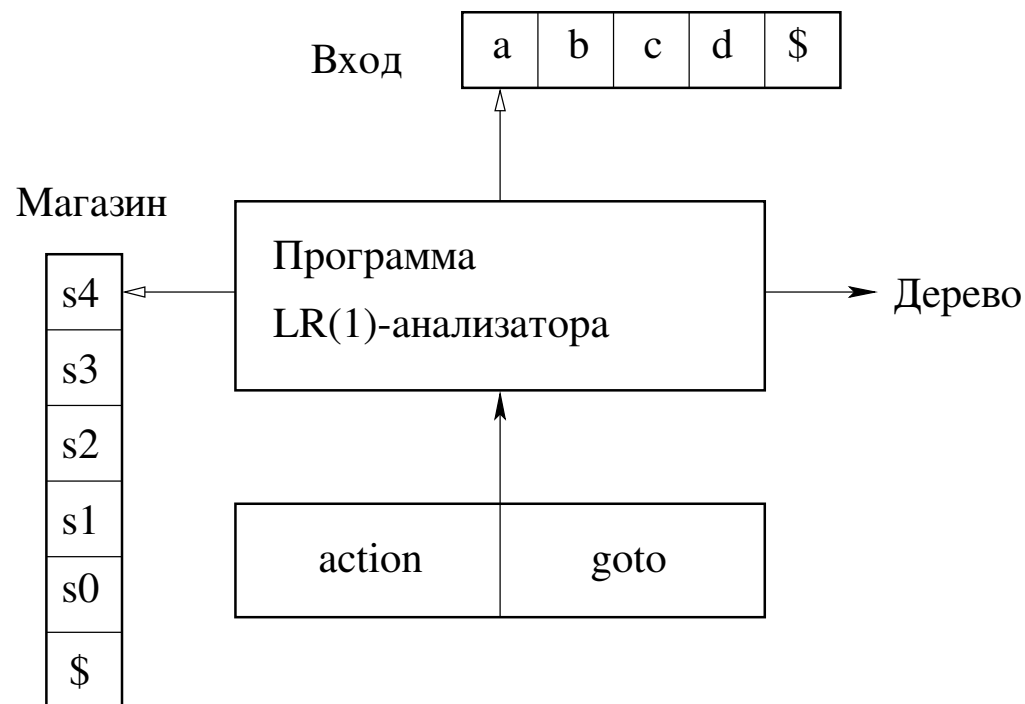
$\text{FOLLOW}(L) = \{=, \$\}$ ,

$\text{FOLLOW}(R) = \{=, \$\}$

№	action				goto			
	=	*	v	\$	$S'$	$S$	$L$	$R$
0	ош.	пер.4	пер.5	ош.	ош.	1	2	3
1	ош.	ош.	ош.	доп.	ош.	ош.	ош.	ош.
2	пер.6, св. $R \rightarrow L$	ош.	ош.	св. $R \rightarrow L$	ош.	ош.	ош.	ош.
3	ош.	ош.	ош.	св. $S \rightarrow R$	ош.	ош.	ош.	ош.
4	ош.	пер.4	пер.5	ош.	ош.	ош.	9	7
5	св. $L \rightarrow v$	ош.	ош.	св. $L \rightarrow v$	ош.	ош.	ош.	ош.
6	ош.	пер.4	пер.5	ош.	ош.	ош.	9	8
7	св. $L \rightarrow *R$	ош.	ош.	св. $L \rightarrow *R$	ош.	ош.	ош.	ош.
8	ош.	ош.	ош.	св. $S \rightarrow L = R$	ош.	ош.	ош.	ош.
9	св. $R \rightarrow L$	ош.	ош.	св. $R \rightarrow L$	ош.	ош.	ош.	ош.

## §29. Алгоритм LR(1)-анализа

Алгоритм LR(1)-анализа является вариантом ПС-алгоритма, работающим за линейное время. В стеке LR(1)-анализатора содержится последовательность состояний, через которые прошёл SLR-распознаватель (или более сложный LR-распознаватель) в процессе распознавания активного префикса. LR(1)-анализатор управляется таблицами распознавателя, не содержащими конфликтов.



```

LR1Parse( $N$ ,  $T$  (содержит  $\$$ ),  $S$ ,
         $s_0$ , action, goto, цепочка  $x$ ) {
    Result = пустая последовательность;
    поместить в стек  $\$$ , а затем  $s_0$ ;
     $a$  = первый символ из  $x$ ;
    loop {
         $s$  = состояние на вершине стека;
        if (action[ $s, a$ ] == "перенос  $s'$ ") {
            поместить в стек  $s'$ ;
             $a$  = следующий символ из  $x$ ;
        } else if (action[ $s, a$ ] == "свёртка  $X \rightarrow u$ ") {
            снять со стека  $|u|$  символов.
             $s'$  = состояние на вершине стека;
            поместить в стек goto[ $s', X$ ];
            добавить правило  $X \rightarrow u$  в Result;
        } else if (action[ $s, a$ ] == "допуск")
            return Result;
        else error();
    }
}

```