

## ЛЕКЦІЯ 12

### Таблично-керований передбачаючий (прогнозуючий) розбір

На рис.1 зображена структура передбачаючого аналізатора, який визначає чергове правило з таблиці. Таку таблицю можна побудувати безпосередньо з граматики.

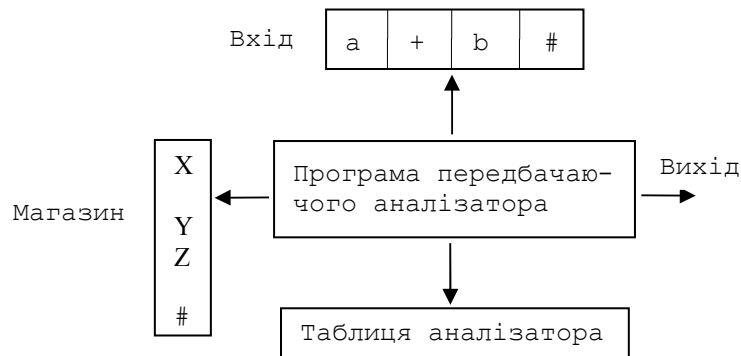


Рис 1.

Таблично-керований передбачаючий аналізатор має вхідний буфер, магазин, таблицю аналізу і вихід. Вхідний буфер містить рядок розпізнавання, за яким йде # – правий кінцевий маркер (ознаку кінця рядка). Магазин містить послідовність символів граматики з символом # на дні. Спочатку магазин містить початковий символ граматики на верхівці і символ # на дні. Таблиця аналізу – це двовимірний масив  $M[A,a]$ , де  $A$  – нетермінал, і  $a$  – термінал або символ #.

Аналізатор керується програмою, яка працює наступним чином. Програма розглядає  $X$  – символ на верхівці магазину і  $a$  – поточний вхідний символ. Ці два символи визначають дію аналізатора. Є три подальші варіанти:

1. Якщо  $X = a = \#$ , аналізатор зупиняється і повідомляє про успішне завершення розбору.
2. Якщо  $X = a \neq \#$ , аналізатор видаляє  $X$  з магазину і просуває покажчик входу на наступний вхідний символ.
3. Якщо  $X$  – нетермінал, програма бере з таблиці елемент  $M[X,a]$ . За цими координатами в таблиці зберігається або правило для нетерміналу  $X$ , або ознака помилки. Якщо, наприклад,  $M[X,a] = \{X \rightarrow UVW\}$ , аналізатор замінює  $X$  на верхівці магазину на  $WVU$  (в результаті на верхівці буде  $U$ ). Будемо вважати, що аналізатор в якості виходу просто друкує використані правила виводу. Якщо  $M[X,a] = \text{error}$ , аналізатор звертається до підпрограми аналізу помилок.

Поведінка аналізатора може бути описана в термінах конфігурацій автомата розбору.

Спочатку аналізатор знаходиться в конфігурації, в якій магазин містить  $S\#$  ( $S$  – початковий символ граматики) у вхідному буфері  $w\#$  ( $w$  – вхідний рядок), змінна  $InSym$  містить перший символ вхідного рядка. Алгоритм, який використовує таблицю аналізатора  $M$  для здійснення розбору показано нижче (Алгоритм 1).

## Алгоритм 1. Нерекурсивний передбачаючий аналіз.

```

repeat  X := верхній символ магазина;
        if X – термінал чи #
        then if X = InSym
                then
                        begin
                                видалити X з магазина;
                                InSym := черговий символ;
                        end
                else error()
        else /*X = нетермінал*/
                if M[X, InSym] = X → Y1Y2...Yk
                then
                        begin
                                видалити X з магазина;
                                помістити Yk, Yk-1, ..., Y1 в магазин (Y1 на верхівку);
                                вивести правило X → Y1Y2...Yk;
                        end
                else error() /*вхід таблиці M пустий*/
until x X = #; /*магазин пустий*/

```

## Приклад 1. Розглянемо граматику арифметичних виразів у такому вигляді:

- |                              |  |
|------------------------------|--|
| 1. $E \rightarrow T E'$      | $FIRST(E) = FIRST(T) = FIRST(F) = \{ (, id \}$ |
| 2. $E' \rightarrow + T E'$   | $FIRST(E') = \{ +, \epsilon \}$                |
| 3. $E' \rightarrow \epsilon$ | $FIRST(T') = \{ *, \epsilon \}$                |
| 4. $T \rightarrow F T'$      | $FOLLOW(E) = FOLLOW(E') = \{ ), \# \}$         |
| 5. $T' \rightarrow * F T'$   | $FOLLOW(T) = FOLLOW(T') = \{ +, ), \# \}$      |
| 6. $T' \rightarrow \epsilon$ | $FOLLOW(F) = \{ +, *, ), \# \}$                |
| 7. $F \rightarrow ( E )$     |  |
| 8. $F \rightarrow id$        |  |

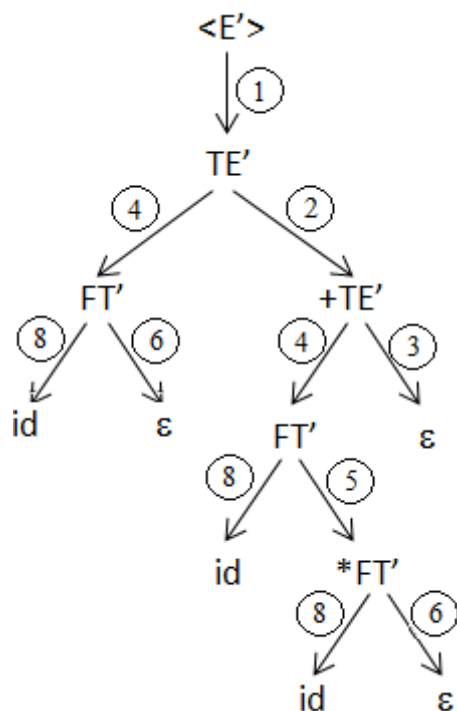
Таблиця передбачаючого аналізатора показана нижче (Таблиця 1). Тут пусті клітинки – входи помилок. Непусті клітинки містять правила, за якими виконується розгортка нетермінала.

Таблиця 1

Нетермінал	Вхідний символ					
	id	+	*	(	)	#
<b>E</b>	$E \rightarrow T E'$			$E \rightarrow T E'$		
<b>E'</b>		$E' \rightarrow + T E'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
<b>T</b>	$T \rightarrow F T'$			$T \rightarrow F T'$		
<b>T'</b>		$T' \rightarrow \epsilon$	$T' \rightarrow * F T'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
<b>F</b>	$F \rightarrow id$			$F \rightarrow ( E )$		

Таблиця 2

Вхід	Магазин	Вихід (застосоване правило)	№ викор. правила
id+id*id#	E#		
id+id*id#	TE'#	$E \rightarrow TE'$	1
id+id*id#	FT'E'#	$T \rightarrow FT'$	4
id+id*id#	idT'E'#	$F \rightarrow id$	8
+id*id#	T'E'#		
+id*id#	E'#	$T' \rightarrow \epsilon$	6
+id*id#	+TE'#	$E' \rightarrow +TE'$	2
id*id#	TE'#		
id*id#	FT'E'#	$T \rightarrow FT'$	4
id*id#	idE'T'#	$F \rightarrow id$	8
*id#	T'E'##		
*id#	F*T'E'##	$T' \rightarrow *FT'$	5
id#	FT'E'##		
id#	idT'E'##	$F \rightarrow id$	8
#	T'E'##		
#	E'##	$T' \rightarrow \epsilon$	6
#	#	$E' \rightarrow \epsilon$	3



Якщо на вхід подається рядок `id + id * id`, передбачаючий аналізатор виконує послідовність кроків, яка показана в таблиці 2. Вказівник вхідного рядка показує на найлівіший символ в колонці «Вхід». Якщо уважно проаналізувати дії аналізатора, то видно, що він виконує лівий вивід, тобто правила застосовуються в відповідності до лівого виводу. За вже переглянутими вхідними символами йдуть символи граматики в магазині (зверху вниз), що відповідає лівим сентенціальним формам виводу. Дерево розбору для рядка `id + id * id` зображено на рис. 2.

## Конструювання таблиць передбачаючого аналізатора

Для конструювання таблиць передбачаючого аналізатора за граматикою  $G$  може бути використано алгоритм, який базується на наступній ідеї. Припустимо, що  $A \rightarrow \alpha$  – правило виводу граматика і  $a \in \text{FIRST}(\alpha)$ . Тоді аналізатор виконує розгортку  $A$  по  $\alpha$ , якщо вхідним символом є  $a$ . Складність виникає, коли  $\alpha = \epsilon$  чи  $\alpha \Rightarrow^* \epsilon$ . В цьому випадку потрібно розгорнути  $A$  в  $\alpha$ , якщо поточний вхідний символ належить  $\text{FOLLOW}(A)$  або дорівнює  $\#$ , причому  $\# \in \text{FOLLOW}(A)$ .

**Алгоритм 2.** Побудова таблиць передбачаючого аналізатора.

Для кожного правила виводу  $A \rightarrow \alpha$  граматика виконати кроки 1 і 2:

Крок 1. Для кожного терміналу  $a$  з  $\text{FIRST}(u)$  додати  $A \rightarrow u$  до  $M[A, a]$ .

Крок 2. Якщо  $\epsilon \in \text{FIRST}(\alpha)$ , додати  $A \rightarrow \alpha$  до  $M[A, b]$  для кожного терміналу  $b$  з  $\text{FOLLOW}(A)$ . Якщо  $\epsilon \in \text{FIRST}(\alpha)$  і  $\# \in \text{FOLLOW}(A)$ , додати  $A \rightarrow \alpha$  до  $M[A, \#]$ .

Крок 3. Всі невизначені входи прирівняти до  $\epsilon$ тог.

**Приклад 2.** Застосуємо алгоритм 2 до граматика, яка була приведена в прикладі 1.

Оскільки  $\text{FIRST}(TE') = \text{FIRST}(T) = \{ (, id \}$ , у відповідності з правилом виводу  $E \rightarrow TE'$  входи  $M[E, (]$  і  $M[E, id]$  стають рівними  $E \rightarrow TE'$ .

В відповідності з правилом виводу  $E' \rightarrow + TE'$  вхід  $M[E', +]$  буде дорівнювати  $E' \rightarrow + TE'$ .

В відповідності з правилом виводу  $E' \rightarrow \epsilon$  входи  $M[E', )]$  і  $M[E', \#]$  будуть дорівнювати  $E' \rightarrow \epsilon$ , оскільки  $\text{FOLLOW}(E') = \{ ), \# \}$ .

Таблиця аналізу, побудована алгоритмом 2, приведена в таблиці 1.

## Умови безповоротного LL(1) синтаксичного аналізу

Якщо для **LL(1)**-граматика можна побудувати синтаксичний аналізатор, що працює без повернень, то отримаємо дві наступні переваги:

- 1) загальний час синтаксичного аналізу буде обмежений величиною, що пропорційна довжині вхідного рядка;
- 2) програма синтаксичного аналізатора може читати вхідні символи по одному, не зберігаючи раніше прочитані символи.

Відомо, що правила будь-якої КВ-граматика можна подати у вигляді правил, що належать лише наступним п'яти типам:

Тип 1.  $X_p \rightarrow X_q | X_r$ , де  $q < p$  і  $r < p$

Тип 2.  $X_p \rightarrow X_q X_r$ , де  $q < p$  і якщо  $X_q \rightarrow^+ \epsilon$ , то і  $r < p$

Тип 3.  $X_p \rightarrow X_q$ , де  $q < p$

Тип 4.  $X_p \rightarrow a$

Тип 5.  $X_p \rightarrow \epsilon$

Для всіх типів правил  $X_p, X_q, X_r \in N$ ,  $a \in T$ .

Відношення  $q < p$  в даному випадку означає таку властивість: якщо граматика не має ліворекурсивних нетерміналів, то всі нетермінали  $X_1, X_2, \dots, X_n$  можна впорядкувати так, що

$X_p \Rightarrow^L X_q$  лише тоді, коли  $q < p$ .

Наприклад, для грамматики

$S \rightarrow P \#$

$P \rightarrow LE$

$E \rightarrow +LE \mid \epsilon$

$L \rightarrow MT$

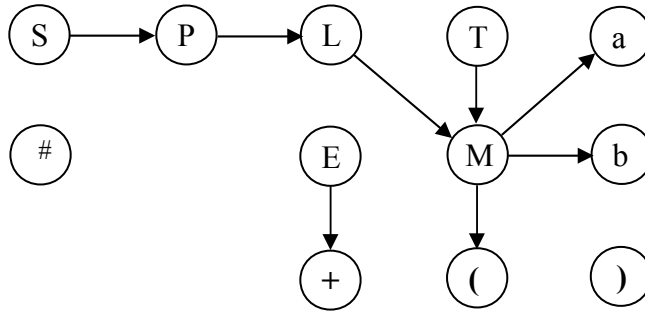
$T \rightarrow MT \mid \epsilon$

$M \rightarrow a \mid b \mid (P),$

де  $S, P, L, E, M, T \in N$ ;

$a, b \in T$

таким впорядкуванням буде  $E, M, T, L, P, S$ .



**Умови безповоротного низхідного розбору формулюються наступним чином:**

1. Жоден нетермінальний символ не є ліворекурсивним.
2. Для жодного правила типу 1 із  $X_q$  і  $X_r$  не виводяться рядки, що починаються з одного і того ж термінального символу, тобто  $FIRST(X_q) \cap FIRST(X_r) = \emptyset$ .
3. Для будь-якого правила типу 1, де  $X_r$  може породжувати порожній рядок  $X_r \Rightarrow^* \epsilon$ , із іншого нетерміналу  $X_q$  не може виводитись рядок, що починається з терміналу, який належить множині FOLLOW нетерміналу лівої частини правила  $X_p$ , тобто завжди виконується  $FIRST(X_q) \cap FOLLOW(X_p) = \emptyset$ . Те ж саме повинне виконуватись, якщо  $X_q$  і  $X_r$  міняються ролями.
4. В жодному правилі типу 1 символ  $X_q$  або  $X_r$  не повинні бути «безвідмовним», іншими словами йому повинна відповідати процедура аналізуючої машини Кнута, яка хоча б в одному випадку повертає значення «false».

**Визначення.** Нетермінальний символ  $X_p$  є «безвідмовним» тоді і тільки тоді, коли відповідне йому правило:

- 1) або належить типу 5 ( $X_p \rightarrow \epsilon$ );
- 2) або належить типу 3 і  $X_q$  є безвідмовним;
- 3) або належить типу 2 і  $X_q$  є безвідмовним;
- 4) або належить типу 1 і  $X_q$  або  $X_r$  є безвідмовним.

Перевірити умову 4 можна шляхом послідовного дослідження символів  $X_p$  на безвідмовність в порядку  $X_1, X_2, \dots, X_p$ , тобто від меншого до більшого.

**Теорема.** Довільна  $LL(1)$  – мова може бути описана граматикою, всі правила якої належать одному з двох типів:

Тип 1.  $A \rightarrow a_1 \beta_1 \mid a_2 \beta_2 \mid \dots \mid a_m \beta_m$

Тип 2.  $A \rightarrow a_1 \beta_1 \mid a_2 \beta_2 \mid \dots \mid a_m \beta_m \mid \epsilon$

де  $a_1, a_2, \dots, a_m$  – різні термінальні символи; в правилах типу 2 жоден із символів  $a_i$  не належить множині FOLLOW(A);  $\beta_i$  – довільні сентенціальні форми.

## Ще одне визначення LL(1)-граматики

Алгоритм 2 для побудови таблиці аналізу  $M$  може бути застосовано до будь-якої граматички. Проте для деяких граматик  $M$  може мати неоднозначно визначені входи. Наприклад, якщо граматика ліворекурсивна чи неоднозначна,  $M$  буде мати хоча б один неоднозначно-визначений вхід. Граматики, для яких таблиці аналізу не мають неоднозначно-визначених входів, повинні бути  $LL(k)$ -граматиками, особливо  $LL(1)$ -граматиками.

Можна показати, що алгоритм 2 для кожної  $LL(1)$ -граматички  $G$  будує таблиці, за якими розпізнаються всі ланцюжки з  $L(G)$ .  $LL(1)$ -граматики мають декілька відмінних властивостей. Неоднозначна чи ліворекурсивна граматика не може бути  $LL(1)$ . Можна також показати, що граматика  $G$  являється  $LL(1)$  тоді і тільки тоді, коли для двох правил вигляду  $A \rightarrow \alpha \mid \beta$  виконується наступне:

- 1) ні для якого терміналу  $a$  одночасно з  $\alpha$  і  $\beta$  не виводяться рядки, які починаються з терміналу  $a$ ;
- 2) тільки з одного із рядків  $\alpha$  чи  $\beta$  може виводитися пустий рядок;
- 3) якщо  $\beta \Rightarrow^* \epsilon$ , то з  $\alpha$  не виводиться ніякий рядок, який починається з терміналу, що належать множині  $FOLLOW(A)$ .

Еквівалентним є наступне визначення:

КВ-граматика називається  $LL(1)$ -граматикою, якщо з існування двох лівих виводів

$$(1) S \Rightarrow^* \gamma A \alpha \Rightarrow \gamma \beta \alpha \Rightarrow^* \gamma x,$$

$$(2) S \Rightarrow^* \gamma A \alpha \Rightarrow \gamma \delta \alpha \Rightarrow^* \gamma y,$$

для яких  $FIRST(x) = FIRST(y)$ , впливає, що  $\beta = \delta$ . Це означає, що для даного ланцюжка  $\gamma A \alpha$  і першого символу, який виводиться з  $A \alpha$  (чи  $\#$ ), існує не більше одного правила, яке може бути застосовано до  $A$ , щоб отримати вивід деякого термінального ланцюжка, який починається з  $\gamma$  і продовжується цим першим символом. Мова, для якої можна побудувати  $LL(1)$ -граматику, називають  $LL(1)$ -мовою. Якщо таблиця аналізу має неоднозначно-визначені входи, то граматика не являє собою  $LL(1)$ . Прикладом може служити наступна граматика:

```
St  $\rightarrow$  if Ex then St
      | if Ex then St else St
      | Cont
Ex  $\rightarrow$  ...
```

Ця граматика неоднозначна, що ілюструється на рис.3. Оскільки граматика неоднозначна, то вона не є  $LL(1)$ -граматикою. Проблема, чи породжує граматика  $LL$ -мову, не має алгоритмічного розв'язку.

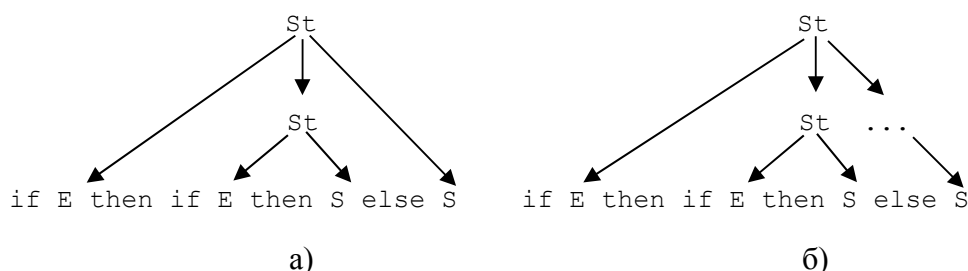


Рис.3.

## Видалення лівої рекурсії

Основна складність при використанні передбачаючого аналізу – це написання такої граматики для вхідної мови, щоб за нею можна було побудувати передбачаючий аналізатор. Іноді за допомогою деяких простих перетворень не LL(1)-граматику можна привести до LL(1)-вигляду. Серед цих перетворень найбільш очевидні являються ліва факторизація і видалення лівої рекурсії. Тут необхідно зробити два зауваження: по-перше, не всяка граматика після цих перетворень стає LL(1), по-друге, після видалення лівої рекурсії і лівої факторизації отримана граматика може стати важкою для розуміння. Граматика ліворекурсивна, якщо в ній міститься нетермінал  $A$  такий, що існує вивід  $A \Rightarrow^+ Au$  для деякого рядка  $u$ . Ліворекурсивні граматики не можуть аналізуватися методами зверху-вниз, тому необхідно видалення лівої рекурсії.

Безпосередню ліву рекурсію (рекурсію вигляду  $A \rightarrow Au$ ) можна видалити наступним чином. Спочатку групуємо  $A$ -правила:

$$A \rightarrow Au_1 \mid Au_2 \mid \dots \mid Au_m \mid v_1 \mid v_2 \mid \dots \mid v_n$$

де ніякий із рядків  $v_i$  не починається з  $A$ . Потім заміняємо  $A$ -правила на

$$A \rightarrow v_1A' \mid v_2A' \mid \dots \mid v_nA' \\ A' \rightarrow u_1A' \mid u_2A' \mid \dots \mid u_mA' \mid \epsilon$$

Нетермінал  $A$  породжує ті ж рядки, що і раніше, але тепер немає лівої рекурсії. За допомогою цієї процедури видаляються всі безпосередні ліві рекурсії, але не видаляється ліва рекурсія, яка включає два чи більше кроки. Приведений нижче алгоритм 3 дозволяє видалити всі ліві рекурсії з граматики.

### Алгоритм 3. Видалення лівої рекурсії.

Крок 1. Впорядковуємо нетермінали в довільному порядку.

Крок 2. for  $i := 1$  to  $n$  do

    for  $j := 1$  to  $i-1$  do

        нехай  $A_j \rightarrow v_1 \mid v_2 \mid \dots \mid v_k$  – всі поточні правила для  $A_j$ ;

        замінити всі правила вигляду  $A_i \rightarrow A_j u$  на правила

$A_i \rightarrow v_1 u \mid v_2 u \mid \dots \mid v_k u$ ;

    end;

    видалити безпосередню ліву рекурсію в правилах для  $A_i$ ;

end

Після  $(i-1)$ -ї ітерації зовнішнього циклу на кроці 2 для любого правила вигляду  $A_k \rightarrow A_l u$ , де  $kk$ . В результаті на наступній ітерації (по  $i$ ) внутрішній цикл (по  $j$ ) послідовно збільшує нижню межу по  $m$  в любому правилі  $A_i \rightarrow A_m u$ , поки не буде виконано  $m \geq i$ . Потім, викидаючи безпосередню ліву рекурсію для  $A_i$ -правил, робимо  $m$  більше  $i$ .

Алгоритм 3 застосовується, якщо граматика не має циклів (виводів вигляду  $A \Rightarrow^+ A$ ) і  $\epsilon$ -правил (правил вигляду  $A \rightarrow \epsilon$ ). Як цикли, так і  $\epsilon$ -правила можуть бути видалені попередньо. Отримувана граматика без лівої рекурсії може мати  $\epsilon$ -правила.

## Ліва факторизація

Основна ідея лівої факторизації в тому, коли незрозуміло яку з двох альтернатив потрібно використовувати для розгортки нетерміналу  $A$ , потрібно переробити  $A$ -правила так, щоб відкласти рішення до того часу, коли не буде достатньо інформації, щоб прийняти правильне рішення.

Якщо  $A \rightarrow uv1 \mid uv2$  – два  $A$ -правила і вхідний рядок починається з непустого рядка, який виводиться з  $u$ , ми не знаємо, чи розгортати по  $uv1$  чи по  $uv2$ . Проте можна відкласти рішення, розгорнувши  $A \rightarrow uA'$ . Тоді після аналізу того, що виводимо з  $u$ , можна розгорнути  $A' \rightarrow v1$  або  $A' \rightarrow v2$ .

Лівофакторизовані правила приймають вигляд:

$$\begin{aligned} A &\rightarrow u A' \\ A' &\rightarrow v1 \mid v2 \end{aligned}$$

**Алгоритм 4.** Ліва факторизація граматики.

Для кожного нетерміналу  $A$  шукаємо найдовший префікс  $u$ , спільний для двох чи більше його альтернатив. Якщо  $u \neq \epsilon$ , тобто існує нетривіальний спільний префікс, замінюємо всі  $A$ -правила

$$A \rightarrow uv1 \mid uv2 \mid \dots \mid uvn \mid z,$$

де  $z$  – всі альтернативи, які не починаються з  $u$ ,

на

$$\begin{aligned} A &\rightarrow uA' \mid z \\ A' &\rightarrow v1 \mid v2 \mid \dots \mid vn \end{aligned}$$

Тут  $A'$  – новий нетермінал. Повторно застосовуємо це перетворення, поки дві альтернативи не будуть мати спільного префікса.

**Приклад 3.** Розглянемо знову граматику умовних операторів:

$$\begin{aligned} St &\rightarrow \text{if } Ex \text{ then } St \\ &\quad \mid \text{if } Ex \text{ then } St \text{ else } St \\ &\quad \mid \text{Cont} \\ Ex &\rightarrow \dots \end{aligned}$$

Після лівої факторизації граматика приймає вигляд

$$\begin{aligned} St &\rightarrow \text{if } Ex \text{ then } St \ St' \\ &\quad \mid \text{Cont} \\ St' &\rightarrow \text{else } St \mid \epsilon \\ Ex &\rightarrow \dots \end{aligned}$$

На жаль, граматика залишається неоднозначною, і значить не є  $LL(1)$ , що ілюструється на рис.4.

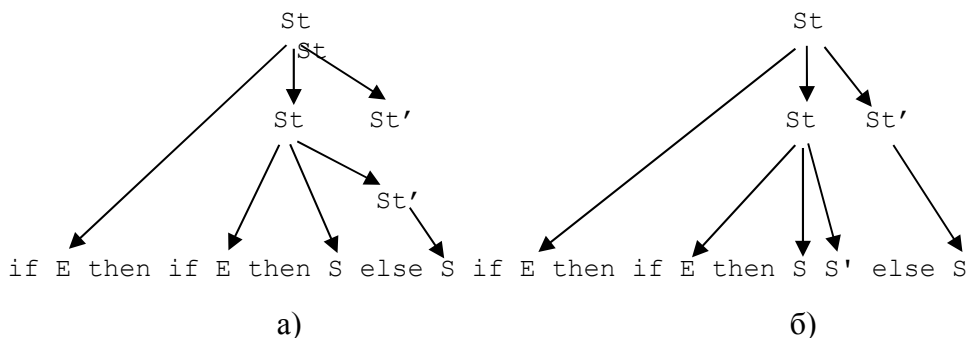


Рис.4.



## **Відновлення після синтаксичних помилок**

В наведених програмах використовувалась процедура реакції на синтаксичні помилки `error()`. В найпростішому випадку ця процедура видає діагностику и завершує роботу аналізатора. Але можна спробувати деяким чином продовжити роботу. Для розбору зверху вниз можна запропонувати наступний простий алгоритм.

Якщо в момент виявлення помилки на верхівці магазину опинився нетермінальний символ  $N$  і для нього нема правила, яке відповідає вхідному символу, то скануємо вхід доти, доки не зустрінемо символ або з `FIRST(N)`, або з `FOLLOW(N)`. В першому разі розгортаємо  $N$  за відповідним правилом, в другому – видаляємо  $N$  з магазину.

Якщо на верхівці магазину термінальний символ, то можна викинути всі термінальні символи з верхівки магазину аж до першого (зверху) нетермінального символу і продовжувати так, як це було описано вище.