



Автоматы и формальные языки

Карпов Юрий Глебович
профессор, д.т.н., зав.кафедрой
“Распределенные вычисления и компьютерные сети”
Санкт-Петербургского Политехнического университета

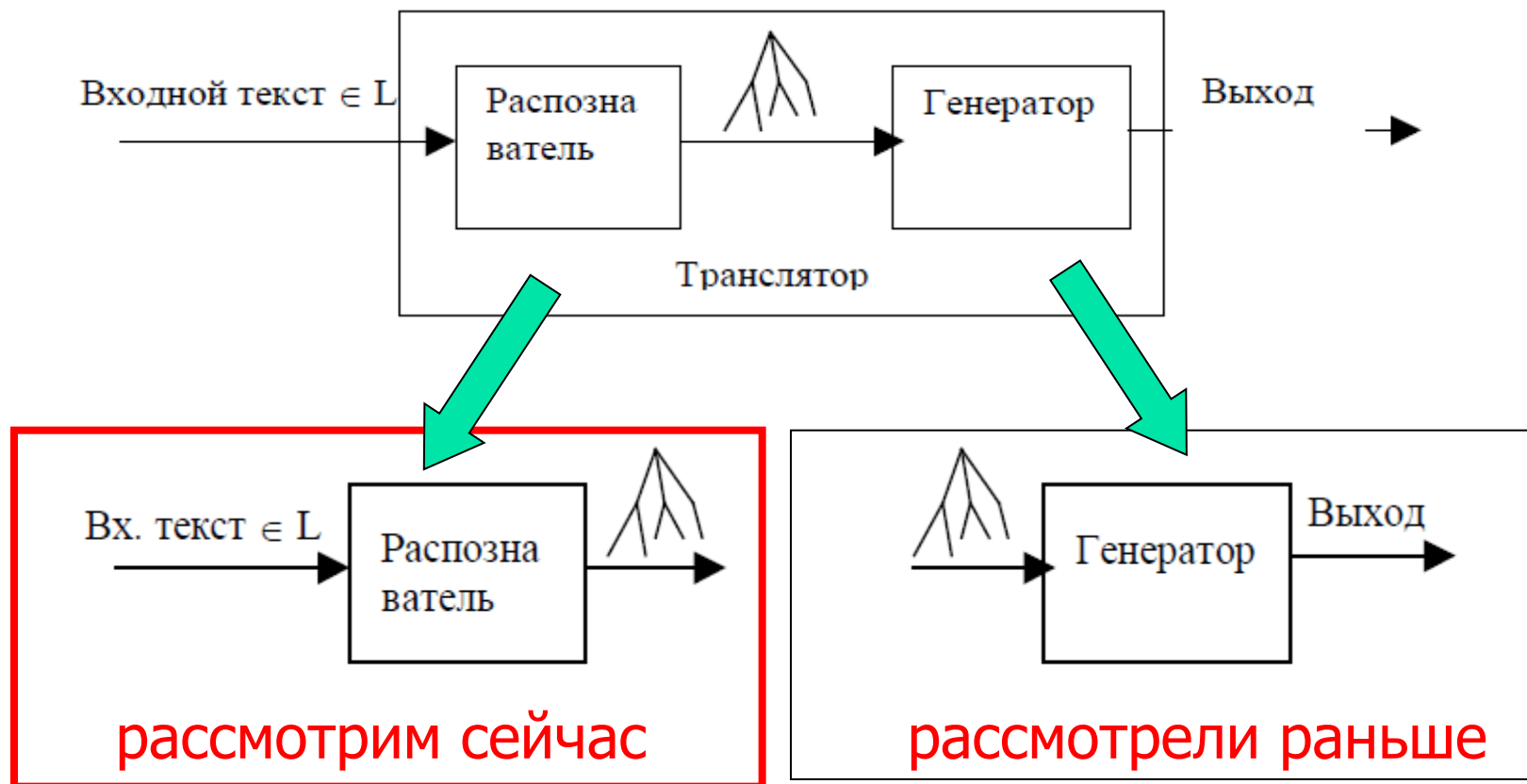
karpov@dcn.infos.ru



Структура курса

- Конечные автоматы-распознаватели – 4 л
- Порождающие грамматики Хомского – 3 л
- Атрибутные трансляции и двусмысленные КС-грамматики – 2 л
- Распознаватели КС-языков и трансляция – 6 л
 - Лекция 10. s-грамматики, LL(k)-грамматики, грамматики рекурсивного спуска
 - Лекция 11. Построение транслятора языка Милан методом рекурсивного спуска
 - Лекция 12. Грамматики предшествования, LR(k)-грамматики
 - Лекция 13. SLR(k) и LALR(k)-грамматики.
 - Лекция 14. Компиляторы компиляторов. Yacc и Bison.
 - Лекция 15. Грамматики Кока-Янгера-Касами и Эрли
- Дополнительные лекции - 2 л

Задачи распознавателя и генератора



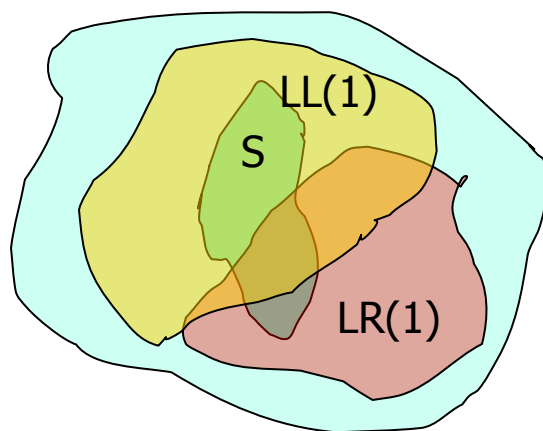
Сегодня рассмотрим, как построить дерево вывода любой входной цепочки языка, порождаемого заданной грамматикой

Синтаксический анализ - одна из наиболее глубоко проработанных и понятных ветвей Computer Science

- Для КС-грамматик синтаксический анализ - это восстановление дерева вывода данной цепочки в данной грамматике, фактически, это алгоритм проверки принадлежности входной цепочки языку, заданному КС-грамматикой



Существуют “универсальные” алгоритмы СА, работающие для любой КС-грамматики, **но они неэффективны**: они должны уметь строить несколько деревьев



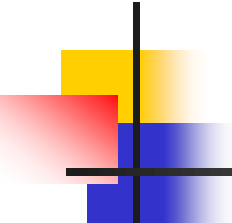
Подклассы КС-грамматик
s-грамматики
LL(k)-грамматики,
Грамматики рекурсивного спуска
Грамматики предшествования
LR(k)-грамматики,
LALR(k)-грамматики,
...

Возможность применения эффективных алгоритмов СА зависит от вида КС-грамматики



Функция грамматики языка

- Грамматика – конечный формализм определения бесконечного числа предложений языка
- Но функции грамматики не только в этом: на основе грамматики выполняются и синтаксический анализ, и семантические вычисления
- Формальный язык может иметь практическую ценность только в том случае, если для него можно построить эффективный транслятор, а эффективность транслятора зависит от вида грамматики
- Для языка нужно ПОДБИРАТЬ грамматику, позволяющую выполнить эффективную трансляцию (эффективный синтаксический анализ и эффективные семантические вычисления)
- При создании языка в грамматику языка (и в язык!) включаются правила с необходимыми разделителями, со структурой, которая позволяет эффективно выполнить синтаксический анализ и семантические вычисления



Преобразования грамматик, сохраняющие эквивалентность



Преобразования КС-грамматик, сохраняющие эквивалентность

- Две грамматики эквивалентны, если они порождают один и тот же язык
- Проблема проверки эквивалентности двух КС-грамматик в общем случае неразрешима. Однако часто удается преобразовать грамматику к виду (удобному для применения того или иного метода анализа), не изменяя язык, порождаемый этой грамматикой

Это, фактически, задача приведения КС-грамматики к “хорошему”, “чистому” виду: например:

выбрасывание неиспользуемых правил;

выбрасывание правил вида $A \rightarrow A$;

выбрасывание нетерминалов, которые вообще не встречаются в выводах

Применение этих алгоритмов к грамматике называют
“гигиеной” грамматик

Приведение КС-грамматик

(выбрасывание правил, не влияющих на выводы)

$G:: S \rightarrow aSb$

$S \rightarrow cA$

$S \rightarrow aD$

$A \rightarrow aAbB$

$A \rightarrow ad$

$B \rightarrow baa$

$B \rightarrow bA$

$C \rightarrow abS \longrightarrow$ нетерминал C не выводится из S

$A \rightarrow A \longrightarrow$ бессмысленное правило

$D \rightarrow bDcD \longrightarrow$ нетерминал D выводится из начального символа, но из него нельзя вывести терминальную строку

Нетерминал A полезный, iff

$$S \Rightarrow^* \alpha A \beta \Rightarrow^* \gamma \quad \gamma \in T^*$$

его можно “вывести” из S

из него выводится произвольная терминальная цепочка

Все правила, содержащие бесполезные нетерминалы, могут быть удалены из грамматики Порождаемый язык не изменится

- Нетерминал A *достижимый*, если существует вывод $S \Rightarrow^* \alpha A \beta$.
Нетерминал A *продуктивный*, если из него можно вывести терминальную цепочку, т.е. существует вывод $A \Rightarrow^* \gamma, \gamma \in T^*$
- Полезные нетерминалы: *достижимые* И *продуктивные*
- **Приведение грамматик:** найти все *бесполезные* нетерминалы и выбросить из грамматики все правила, в которых они встречаются

Построение множества достижимых нетерминалов

- Для построения множества **достижимых** нетерминалов последовательно строятся множества $M_0, M_1, \dots, M_k, \dots$
- Начальное множество M_0 достижимых нетерминалов равно $\{S\}$. В следующее множество M_{k+1} добавляются нетерминалы B , такие, что в грамматике существует правило $A \Rightarrow \alpha B \beta$ для какого-нибудь A из предыдущего множества M_k

$$M_0 = \{S\},$$

// любой вывод начинается с S

$$M_{k+1} = M_k \cup \{B \mid \exists A \in M_k: A \Rightarrow \alpha B \beta\}$$

// добавляем те, которые выводятся

Пример.

$S \rightarrow SbAc \mid dA \mid d$
 $A \rightarrow AcC \mid abc \mid dAE$
 $B \rightarrow ScC \mid dB$
 $C \rightarrow cC \mid DdAS$
 $D \rightarrow cE \mid CdA$
 $E \rightarrow Ac \mid Dd$
 $F \rightarrow cC \mid a \mid dAE$
 $G \rightarrow AcC \mid bc \mid dA$

$$M_0 = \{S\};$$

$$M_1 = \{S, A\};$$

$$M_2 = \{S, A, C, E\};$$

$$M_3 = \{S, A, C, E, D\};$$

$$M_4 = M_3$$

$$M_\infty = \{S, A, C, D, E\}$$

Пример.

$S \rightarrow SbAc \mid dA \mid d$
 $A \rightarrow AcC \mid abc \mid dAE$
 ~~$B \rightarrow ScC \mid dB$~~
 $C \rightarrow cC \mid DdAS$
 $D \rightarrow cE \mid CdA$
 $E \rightarrow Ac \mid Dd$
 ~~$F \rightarrow cC \mid a \mid dAE$~~
 ~~$G \rightarrow AcC \mid bc \mid dA$~~

**Бесполезные
правила
выбрасываем**

Построение множества продуктивных нетерминалов

(из которых могут быть порождены терминальные цепочки)

- Для построения множества продуктивных нетерминальных символов последовательно строятся множества V_0, \dots, V_k, \dots ,
- V_0 включает такие нетерминалы, из которых за один шаг выводятся терминальные цепочки
- $V_{k+1} = V_k \cup \{B \mid (\exists \alpha \in V^*): B \Rightarrow \alpha\}$ - в следующее множество V_{k+1} добавляются нетерминалы B такие, что в грамматике существует правило $B \Rightarrow \alpha$, где α включает терминалы и нетерминалы только из предыдущего множества V_k

Грамматика G

$S \rightarrow SbAc \mid dA \mid d$
 $A \rightarrow AcC \mid abc \mid dAE$
 $B \rightarrow ScC \mid dB$
 $C \rightarrow cC \mid DdAS$
 $D \rightarrow cE \mid CdA$
 $E \rightarrow Ac \mid Dd$
 $F \rightarrow cC \mid a \mid dAE$
 $G \rightarrow AcC \mid bc \mid dA$

$$V_0 = \{S, A, F, G\};$$

$$V_1 = \{S, A, F, G, E\};$$

$$V_2 = \{S, A, F, G, E, D\}; V_3 = V_2$$

$$M_\infty = \{S, A, C, E, D\}$$

$$V_\infty = \{S, A, F, G, E, D\}$$

Берем пересечение M и V

$$M_\infty \cap V_\infty = \{S, A, D, E\}$$

Грамматика G'

$S \rightarrow SbAc \mid dA \mid d$
 $A \rightarrow abc \mid dAE$
 $D \rightarrow cE$
 $E \rightarrow Ac \mid Dd$

В грамматике оставляем только те правила, в которые входят нетерминалы только из пересечения

Построили приведенную грамматику G'

Пример

- Устранить бесполезные нетерминалы в грамматике:

$S \rightarrow aC \mid Ac \mid bDe$

$A \rightarrow cAB$

$B \rightarrow b \mid dB$

$C \rightarrow cCB \mid a$

$E \rightarrow cE \mid BdA \mid d$

- Решение

- Достижимые терминальные символы: $\{S, A, B, C, D\}$
- Продуктивные терминальные символы: $\{B, C, E, S\}$
- Полезные терминальные символы (пересечение множеств достижимых и продуктивных): $\{S, B, C\}$

- Приведенная грамматика:

$S \rightarrow aC$

$B \rightarrow b \mid dB$

$C \rightarrow cCB \mid a$

$M_0 = \{S\};$

$M_1 = \{S, C, A, D\};$

$M_2 = \{S, C, A, D, B\}; M_3 = M_2$

$V_0 = \{C, B, E\};$

$V_1 = \{S, C, B, E\}; V_2 = V_1$

Проблема пустоты КС-языка разрешима: язык, порождаемый КС-грамматикой НЕпуст, если и только если S является продуктивным символом



ϵ - свободные и неукорачивающие КС-грамматики

- **Определение.** КС-грамматика называется **неукорачивающей**, если она не включает продукций вида $A \rightarrow \epsilon$
- Любой шаг вывода в неукорачивающей грамматике не может уменьшить длину выводимой цепочки – отсюда и ее название.
- КС-грамматика называется **ϵ - свободной**, если она неукорачивающая или в ней существует ровно одна продукция вида $S \rightarrow \epsilon$, где S – начальный нетерминал, и S не встречается в правой части ни одной из продукций

Приведение к ϵ -свободной грамматике

- **Теорема.** По любой КС-грамматике может быть построена эквивалентная ϵ -свободная КС-грамматика. По любой КС-грамматике, порождающей язык, не включающий пустую цепочку, может быть построена эквивалентная неукорачивающая грамматика.
 - **Доказательство.** Для любого правила $A \rightarrow \epsilon$ скопируем (повторим) все продукции, включающие A в правой части, с выбрасыванием символа A из правых частей этих копий, причем если A входит несколько раз в правую часть продукции, то такое выбрасывание произведем во всех возможных сочетаниях. Далее, продукцию $A \rightarrow \epsilon$ выбросим из множества R продукций грамматики.
 - Если S - начальный нетерминал, и $S \rightarrow \epsilon$ – среди продукций, то выбираем новый начальный нетерминал S' , и к продукциям грамматики добавляем две новых: $S' \rightarrow \epsilon$ и $S' \rightarrow S$

- **Пример.** Привести к неукорачивающему виду грамматику:

$$\begin{aligned} S &\rightarrow cA \mid \epsilon \\ A &\rightarrow cA \mid bA \mid \epsilon \end{aligned}$$

- **Решение.**

$$\begin{aligned} S' &\rightarrow S \mid \epsilon \\ S &\rightarrow cA \mid c \\ A &\rightarrow cA \mid bA \mid c \mid b \end{aligned}$$

ϵ – свободные языки

- **Определение.** Язык называется ϵ - свободным, если он не включает пустой цепочки
- По любой КС-грамматике можно проверить, является ли порождаемый ею язык ϵ - свободным
 - Для этого нужно проверить, есть ли в эквивалентной ϵ - свободной грамматике продукция $S \rightarrow \epsilon$.
- Любую грамматику можно преобразовать в неукорачивающую:
 - добавим в грамматику новый нетерминал S' и продукцию $S' \rightarrow S\epsilon$
 - применим алгоритм приведения грамматики к ϵ - свободному виду

Следствие: Проблема пустоты (содержит ли порождаемый грамматикой язык пустую цепочку?) для КС грамматик разрешима

В некоторых случаях для упрощения удобно к грамматике добавлять правило $S' \rightarrow \# S \#$

Ясно, цепочки языка, порождаемого новой грамматикой будут иметь вид $\#\alpha\#$, где α - цепочка прежнего языка

Сингулярные продукции (вида $A \rightarrow B$)

- Эффективный алгоритм, позволяющий по любой ε -свободной КС-грамматике построить эквивалентную КС-грамматику без сингулярных продукции:
 - Включим в множество R продукции грамматики G вместо каждой сингулярной продукции вида $A \rightarrow B$ все продукции $A \rightarrow \beta$ такие, что продукция $B \rightarrow \beta \in R$ и несингулярна

- **Пример.** Для грамматики:

1. $S \rightarrow BA$
2. $A \rightarrow C \mid ac$
3. $B \rightarrow b$
4. $C \rightarrow A$

эквивалентная грамматика без сингулярных продукции имеет вид:

1. $S \rightarrow BA$
2. $A \rightarrow ac$
3. $B \rightarrow b$
4. $C \rightarrow ac$

В результате приведения этой грамматики последняя продукция будет выброшена. Эквивалентная приведенная грамматика:

1. $S \rightarrow BA$
2. $A \rightarrow ac$
3. $B \rightarrow b$



Сингулярные productions: осторожно!

- В языках программирования сингулярные productions используются достаточно часто, и их выбрасывание может нарушить очевидность структуры productions. Исключение составляет специальный случай сингулярных productions, которые приводят к циклу вида $A \Rightarrow^+ A$.
- **Определение.** КС-грамматика, в которой не существует вывода $A \Rightarrow^+ A$, называется ациклической.

Любую КС-грамматику можно привести к эквивалентному ациклическому виду.

Действительно, построим эквивалентную ε - свободную грамматику, и затем приведем ее к эквивалентному виду без сингулярных productions.

Левая рекурсия

- **Определение.** КС-грамматика называется грамматикой с левой рекурсией, если в ней существует вывод $A \Rightarrow^* A\beta$.

- Леворекурсивные грамматики не всегда удобны. В частности, нисходящие методы синтаксического анализа, восстанавливающие дерево вывода от его корня, не могут быть применены для грамматик с левой рекурсией. Существует алгоритм, позволяющий привести любую КС-грамматику к эквивалентному виду без левой рекурсии

- Рассмотрим случай "**прямой левой рекурсии**", когда грамматика имеет продукции вида $A \rightarrow A\alpha$
- Очевидно, что A будет "полезным" символом, только если в грамматике есть продукции вида $A \rightarrow \beta$, т.е. грамматика содержит, как минимум, пару продукций

$$A \rightarrow A\alpha \mid \beta$$

Многократное применение правила $A \rightarrow A\alpha$ даст $A \Rightarrow^* A\alpha^*$, и последняя замена A на β даст $A \Rightarrow^* \beta\alpha^*$

$$A \Rightarrow A\alpha \Rightarrow A\alpha\alpha \Rightarrow A\alpha\alpha\alpha \Rightarrow^* A\alpha^* \Rightarrow \beta\alpha^*$$

- Вывод: пара продукций

$$A \rightarrow A\alpha \mid \beta$$

// из A выводятся только цепочки вида $\beta\alpha\alpha\alpha\alpha \dots$

может быть заменена следующими нелеворекурсивными продукциями:

$$A \rightarrow \beta X$$

$$X \rightarrow \alpha X \mid \varepsilon$$

Пример

- Классическая грамматика арифметических выражений:
$$E \rightarrow E+T \mid T$$
 - Используя алгоритм, преобразуем эти правила с прямой левой рекурсией к следующему эквивалентному виду:
$$E \rightarrow TE'$$
$$E' \rightarrow +TE' \mid \varepsilon$$
- Этот метод легко обобщается на случай, когда нетерминал имеет несколько альтернатив с прямой левой рекурсией. Более сложный общий случай избавления от непрямой левой рекурсии требует чуть более сложного алгоритма

В БНФ-нотации грамматика арифметических выражений без левой рекурсии:

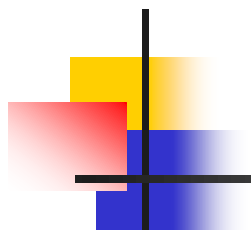
$$\begin{aligned} E &\rightarrow E+T \mid T \\ T &\rightarrow T*R \mid R \\ R &\rightarrow i \end{aligned}$$

$$\begin{aligned} E &::= T \{+T\} \\ T &::= R \{*R\} \\ R &::= i \end{aligned}$$

<выражение> - это **<терм>** с идущими за ним произвольным числом **+<терм>**

<терм> - это **<первичное>** с идущими за ним произвольным числом ***<первичное>**

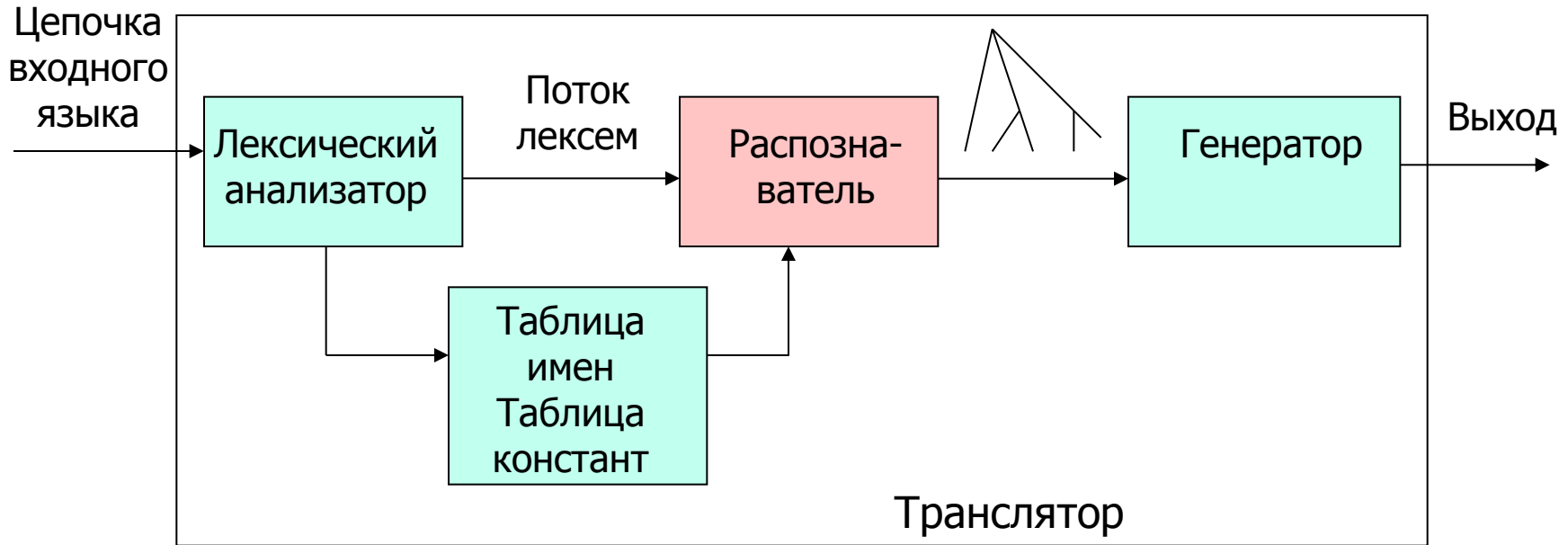
'+' можно понимать, как addop, '*' – как mulop



Постановка задачи синтаксического анализа

Синтаксический анализ КС-языков

- Для КС-грамматик синтаксический анализ - это восстановление дерева вывода данной цепочки в данной грамматике



Наша проблема – проблема построения алгоритма распознавания:
как воссоздать вывод произвольной цепочки в заданной грамматике?

Синтаксический анализ – одна из наиболее проработанных и понятных ветвей
Computer Science, и это - очень красивая ветвь

Задача блока синтаксического анализа (СА)

- Рассмотрим пример КС-грамматики.

$G_{4.5}$:

1. $S \rightarrow abScB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bBc$
6. $B \rightarrow c$

Пусть в этой грамматике есть вывод некоторой терминальной цепочки:

$$S \Rightarrow abScB \Rightarrow abbAcB \Rightarrow abbcBAcB \Rightarrow abbccAcB \Rightarrow abbccabcbB \Rightarrow abbccabcc$$

- Задача блока синтаксического анализа обратна: по заданной терминальной цепочке восстановить вывод ее из начального символа грамматики:

$$S \Rightarrow \dots \quad ? \quad ? \quad ? \quad ? \quad \dots \quad \Rightarrow abbccabcc$$

- Синтаксический анализ можно понимать и как восстановление синтаксического дерева заданной входной цепочки. Для этого дерева известны листья – (входная цепочка) и корень (начальный символ грамматики)

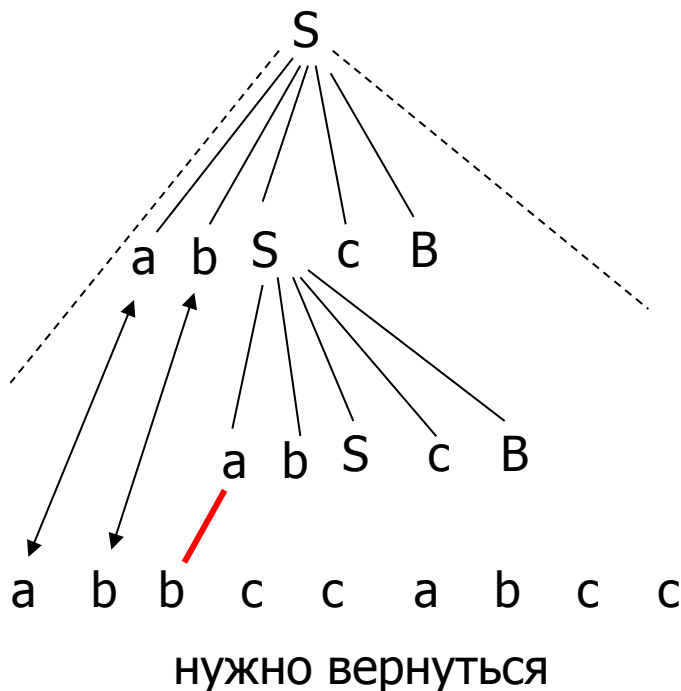
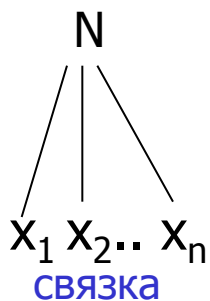
Нисходящие алгоритмы синтаксического анализа

$S \Rightarrow_1 abScB \Rightarrow_2 abbAcB \Rightarrow_4 abbcBAcB \Rightarrow_6 abbccAcB \Rightarrow_3 abbccabcB \Rightarrow_6 abbccabcc$

левый вывод

1. $S \rightarrow abScB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bBc$
6. $B \rightarrow c$

$N \rightarrow X_1 X_2 \dots X_n$



Нисходящий синтаксический анализ работает так:

Начиная от начального символа дерева пытаемся найти, какой альтернативой (какой правой частью продукции) заменить очередной нетерминал, чтобы раскрыть новый узел синтаксического дерева.

Восстанавливает **ЛЕВЫЙ** канонический вывод

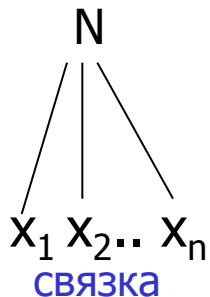
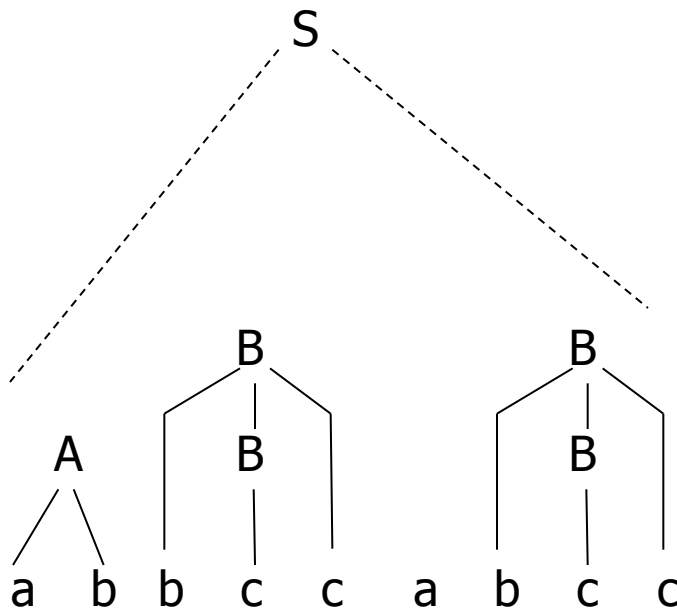
Повторяющимся вопросом здесь является следующий: какую альтернативу для текущего нетерминала нужно выбрать, чтобы из начального символа по текущей сентенциальной форме построить следующую сентенциальную форму вывода, чтобы получить входную цепочку?

Восходящие алгоритмы синтаксического анализа

$S \Rightarrow_1 abScB \Rightarrow_2 abScc \Rightarrow_4 abbAcc \Rightarrow_6 abbcBAcc \Rightarrow_3 abbcBabcc \Rightarrow_6 abbccabcc$

правый вывод

1. $S \rightarrow abScB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bBc$
6. $B \rightarrow c$



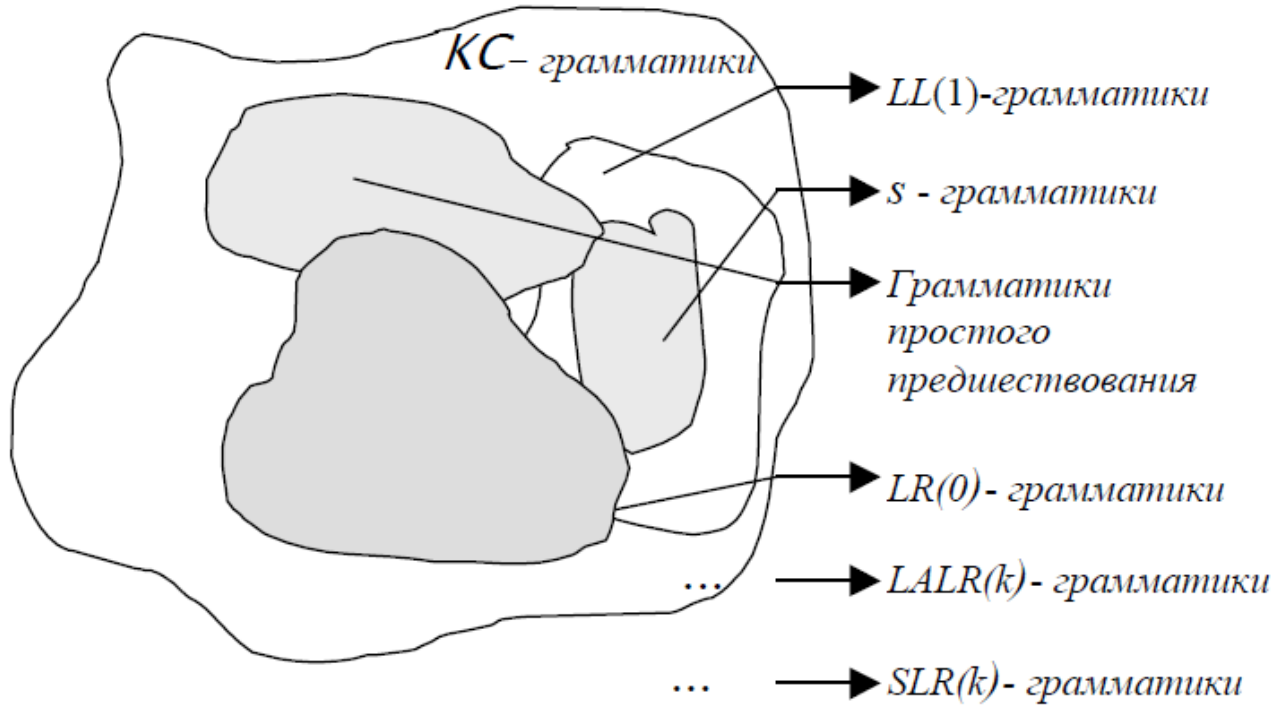
Восходящий синтаксический анализ работает так.

Начиная от терминальной строки листьев дерева пытаемся найти "связку" – правую часть продукции грамматики, которую нужно заменить нетерминалом (левой частью продукции), чтобы получить новый узел синтаксического дерева.

Восстанавливает ПРАВЫЙ канонический вывод

Основным повторяющимся вопросом здесь является следующий: какую подстроку в сентенциальной форме нужно выбрать в качестве связки, и каким нетерминалом ее заменить с получением предыдущей сентенциальной формы вывода, чтобы свернуть всю цепочку к S ?

Подклассы КС-грамматик, допускающие эффективные методы синтаксического анализа



Для всех КС-грамматик существуют общие неэффективные алгоритмы как нисходящие, так и восходящего синтаксического анализа (сложность n^3). Выделенные собственные подклассы – это те КС-грамматики, для которых были найдены различные по простоте и эффективности специфические алгоритмы нисходящего либо восходящего синтаксического анализа



Нисходящие методы синтаксического анализа

Нисходящий синтаксический анализ

$S \Rightarrow_1 abScB \Rightarrow_2 abbAcB \Rightarrow_4 abbcBAcB \Rightarrow_6 abbccAcB \Rightarrow_3 abbccabcbB \Rightarrow_6 abbccabcc$

1. $S \rightarrow abScB$ $S \Rightarrow \dots ??? \Rightarrow \dots abbccabcc$ КАК ВОССТАНОВИТЬ?

2. $S \rightarrow bA$

3. $A \rightarrow ab$

4. $A \rightarrow cBA$

5. $B \rightarrow bBc$

6. $B \rightarrow c$

1. По какому правилу заменен S на I шаге, по 1 или 2?

Можем гадать, а можем смотреть на входную цепочку $abbccabcc$

$S \Rightarrow_1 abSBc$ Что выберем? По первому символу анализируемой цепочки!! Выбираем \Rightarrow_1
 $S \Rightarrow_2 bA$

$S \Rightarrow_1 abSBc \Rightarrow \dots \Rightarrow abbccabcc$ КАК ВОССТАНОВИТЬ ДАЛЬШЕ?

S	\Rightarrow_1	$abSBc$	\Rightarrow_1	$ababScBBc$	$??$	\Rightarrow^*	$abbccabcc$	эта подстановка для S ?
S	\Rightarrow_1	$abSBc$	\Rightarrow_2	$abbABc$	$??$	\Rightarrow^*	$abbccabcc$	эта подстановка для S ?

- Рассуждения о восстановлении вывода цепочки из начального символа эквивалентны рассуждениям о восстановлении синтаксического дерева: восстановление каждого шага вывода соответствует добавлению очередного узла в дерево, и обратно

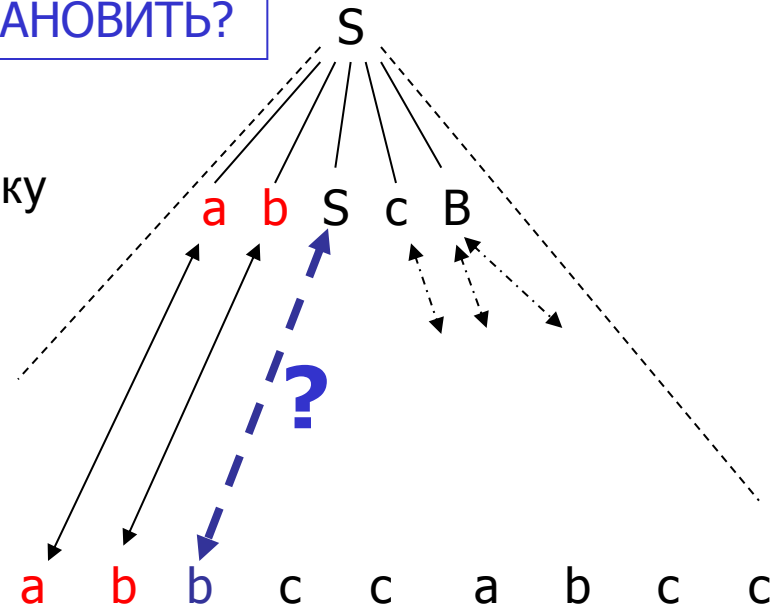
Нисходящий синтаксический анализ

$S \Rightarrow_1 abScB \Rightarrow_2 abbAcB \Rightarrow_4 abbcBAcB \Rightarrow_6 abbccAcB \Rightarrow_3 abbccabcB \Rightarrow_6 abbccabcc$

1. $S \rightarrow abScB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bBc$
6. $B \rightarrow c$

$S \Rightarrow^* abbccabcc$ КАК ВОССТАНОВИТЬ?

Пусть первую подстановку
для S сделали



$S \Rightarrow_1 abScB$

$S \Rightarrow_2 bA$

1. По какому правилу теперь заменять S , по 1 или 2?

Можем гадать, а можем смотреть на остаток входной цепочки $\dots bccabcc$

Что выберем? По очередному терминальному символу входной цепочки!!

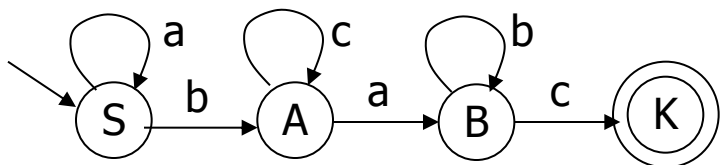
При выводе цепочки в КС-грамматиках если терминальный символ попал в начало сентенциальной формы, то он там и останется в процессе всего вывода



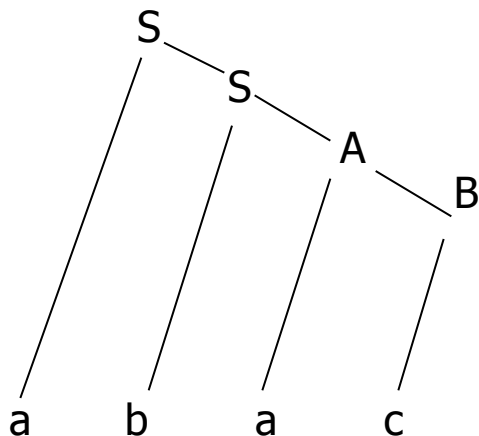
s-грамматики (simple grammars)

Автоматные грамматики и s-грамматики

1. $S \rightarrow abScB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bBc$
6. $B \rightarrow c$



1. $S \rightarrow aS$
2. $S \rightarrow bA$
3. $A \rightarrow aB$
4. $A \rightarrow cA$
5. $B \rightarrow bB$
6. $B \rightarrow c$



$S \Rightarrow aS \Rightarrow abA \Rightarrow abaB \Rightarrow abac$

	a	b	c
S	S	A	ош
A	B	ош	A
B	ош	B	!

Таблица переходов

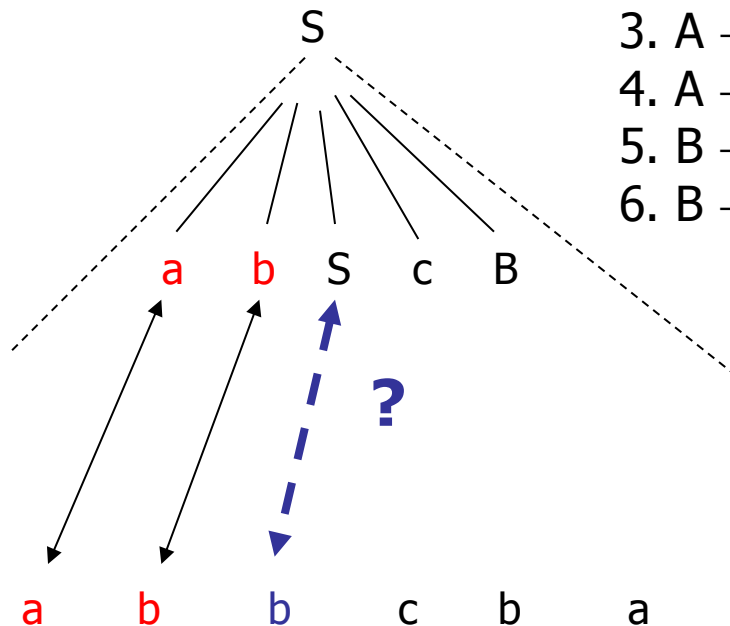


Таблица
принятия
решений

	a	b	c
S	abScB	bA	ош
A	ab	ош	cBA
B	ош	bBc	c
a	X	ош	ош
b	ош	X	ош
c	ош	ош	X

s-грамматики: пары <нетерминал, терминал>

- Автоматные языки удобны тем, что для них может быть построен детерминированный конечный автомат, реализующий эффективный (линейный) нисходящий алгоритм распознавания входной цепочки.
- По каждой паре: <очередной (самый левый) нетерминал sentенциальной формы Q , очередной терминал входной строки> распознаватель автоматного языка (если соответствующий конечный автомат детерминированный) однозначно выбирает подстановку, т.е. правую часть правила $A \rightarrow aB$ либо $A \rightarrow b$, заменяющую A в sentенциальной форме.
- Однозначность выбора достигается тем, что в детерминированном автомате в множестве альтернатив каждого нетерминала $A \rightarrow a_1B_1 \mid a_2B_2 \mid \dots \mid a_nB_n$ все терминальные символы a_i различны
- Назовем первый терминал правой части правила ключом этого правила. Для любого автоматного языка может быть построена такая порождающая его автоматная грамматика, что все ключи в различных альтернативах каждого нетерминала различны
- Для любого неавтоматного языка если альтернативы одного и того же нетерминала начинаются РАЗЛИЧНЫМИ терминалами, то выбор правил подстановки тоже однозначен. Такие грамматики без ϵ -правил называются s-грамматиками



S-грамматики как обобщение автоматных грамматик

- s-грамматики (здесь 's' от слова *simple*, простой) являются простым обобщением автоматных грамматик. s-грамматики имеют вид правил $A \rightarrow a\beta$, где a – терминал (ключ правила), а β – произвольная цепочка из терминалов и нетерминалов, причем у всех альтернатив одного и того же нетерминала ключи различны.
- Это позволяет построить простой детерминированный нисходящий алгоритм распознавания для s-грамматик

Алгоритм нисходящего синтаксического анализа для s-грамматик линейный

- Алгоритм нисходящего синтаксического анализа для s-грамматик реализуется простым алгоритмом со стеком - МП-автоматом с одним состоянием. Алфавит магазина этого МП-автомата включает все терминальные и нетерминальные символы грамматики
 - Вначале в магазин помещается начальный символ грамматики.
 - Очередной терминал входной цепочки и верхний символ магазина определяют, какая цепочка должна быть помещена в магазин
 - Если верхний символ магазина – нетерминал, то очередной входной терминал – это ключ, определяющий правую часть правила грамматики, которой должен быть заменен нетерминал в восстанавливаемом выводе
 - Если верхний символ магазина – терминал, он должен совпадать с очередным терминалом входной цепочки. При таком совпадении терминал выбрасывается из магазина и выполняется сдвиг по входной цепочке – анализатор переходит к следующему терминалу входной цепочки
 - Цепочка распознана, если по исчерпанию цепочки магазин пуст
 - Несовпадение терминала в верхушке магазина с очередным терминалом входной цепочки свидетельствует об ошибке во входной строке
 - Ошибка также обнаруживается, если для нетерминала, находящегося в верхушке магазина, очередной входной терминал не совпадает ни с одним ключом правил грамматики для этого нетерминала

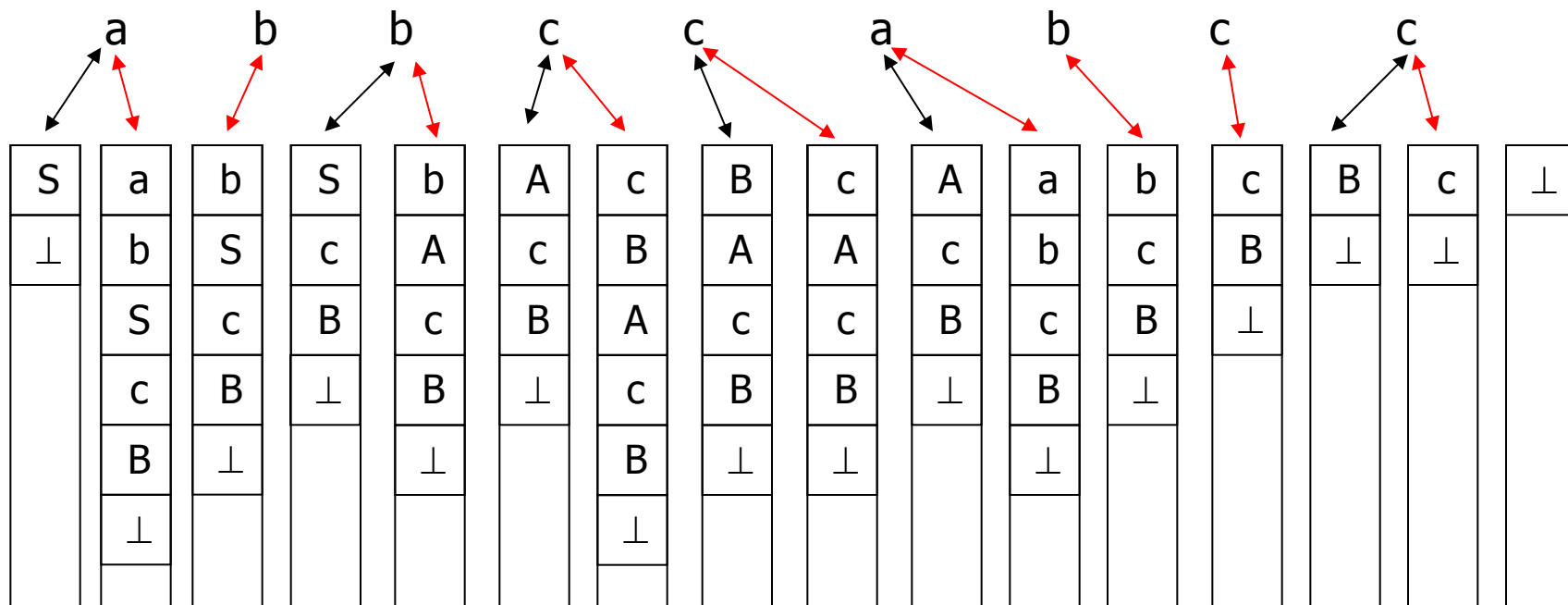
Пример: синтаксический анализ цепочки, порождаемой s-грамматикой

1. $S \rightarrow abScB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bBc$
6. $B \rightarrow c$

$S \Rightarrow abScB \Rightarrow abbAcB \Rightarrow abbcBAcB \Rightarrow$
 $abbccAcB \Rightarrow abbccabcB \Rightarrow abbccabcc$

В решающей таблице могут быть
номера правил, а не правые части
правил

	a	b	c
S	abScB	bA	ош
A	ab	ош	cBA
B	ош	bBc	c
a	X	ош	ош
b	ош	X	ош
c	ош	ош	X



Вопрос: алгоритм распознавания проработал – где дерево вывода???

Семантические вычисления для s-грамматик: операционная семантика

Семантики могут быть включены в любом месте правил

1. $S \rightarrow abScB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bBc$
6. $B \rightarrow c$

$S \Rightarrow abScB \Rightarrow abbAcB \Rightarrow abbcBAcB \Rightarrow$
 $abbccAcB \Rightarrow abbccabcB \Rightarrow abbccabcc$

Выполненные семантики для
цепочки abbccabcc:

$y_1 y_3 y_4 y_3$

1. $S \rightarrow ab y_1 ScB$
2. $S \rightarrow b A$
3. $A \rightarrow a y_4 b$
4. $A \rightarrow c B A$
5. $B \rightarrow bB y_5 c$
6. $B \rightarrow c y_3$

S	a	b	y_1	S	b	A	c	B	c	y_3	A	a	y_4	b	c	B	c	y_3	\perp
\perp	b	y_1	S	c	A	c	B	A	y_3	A	c	y_4	b	c	B	\perp	y_3	\perp	
	y_1	S	c	B	c	B	A	c	A	c	B	b	c	B	\perp		\perp		
	S	c	B	\perp	B	\perp	c	B	c	B	\perp	c	B	\perp					
	c	B	\perp		\perp		B	\perp	B	\perp		B	\perp						
	B	\perp					\perp		\perp			\perp							
	\perp																		

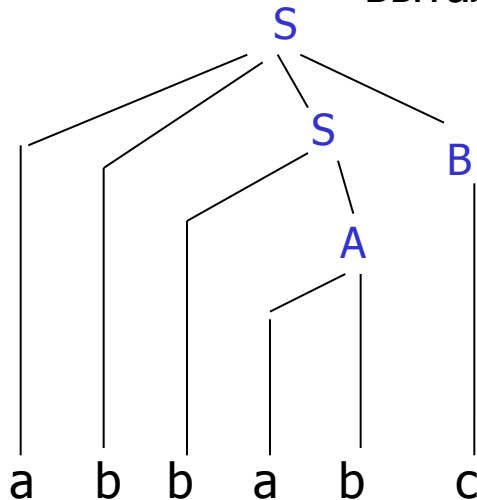
Представления структуры дерева вывода скобочной формой

Грамматика

1. $S \rightarrow abSB$
2. $S \rightarrow bA$
3. $A \rightarrow ab$
4. $A \rightarrow cBA$
5. $B \rightarrow bB$
6. $B \rightarrow c$

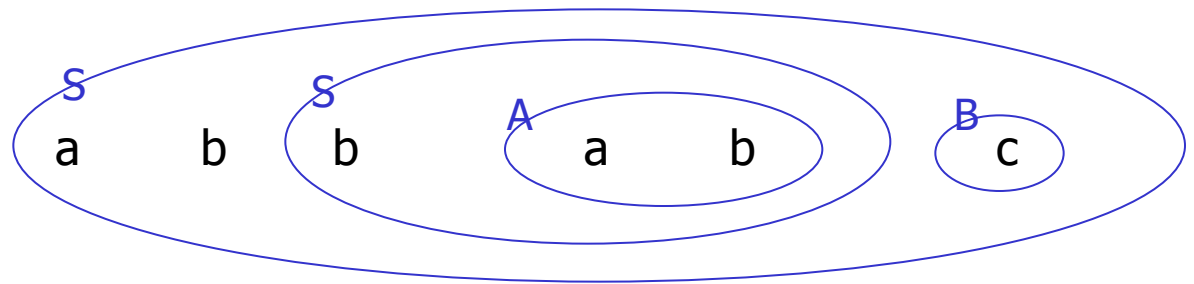
$S \Rightarrow abSB \Rightarrow abbAB \Rightarrow abbabB \Rightarrow abbabc$

При анализе выводить все
выталкиваемые из стека символы



Грамматика с семантическими действиями

1. $S \rightarrow ('_S' abScB '_S')$
2. $S \rightarrow ('_S' b A '_S')$
3. $A \rightarrow ('_A' a b '_A')$
4. $A \rightarrow ('_A' cBA '_A')$
5. $B \rightarrow ('_B' bB '_B')$
6. $B \rightarrow ('_B' c '_B')$

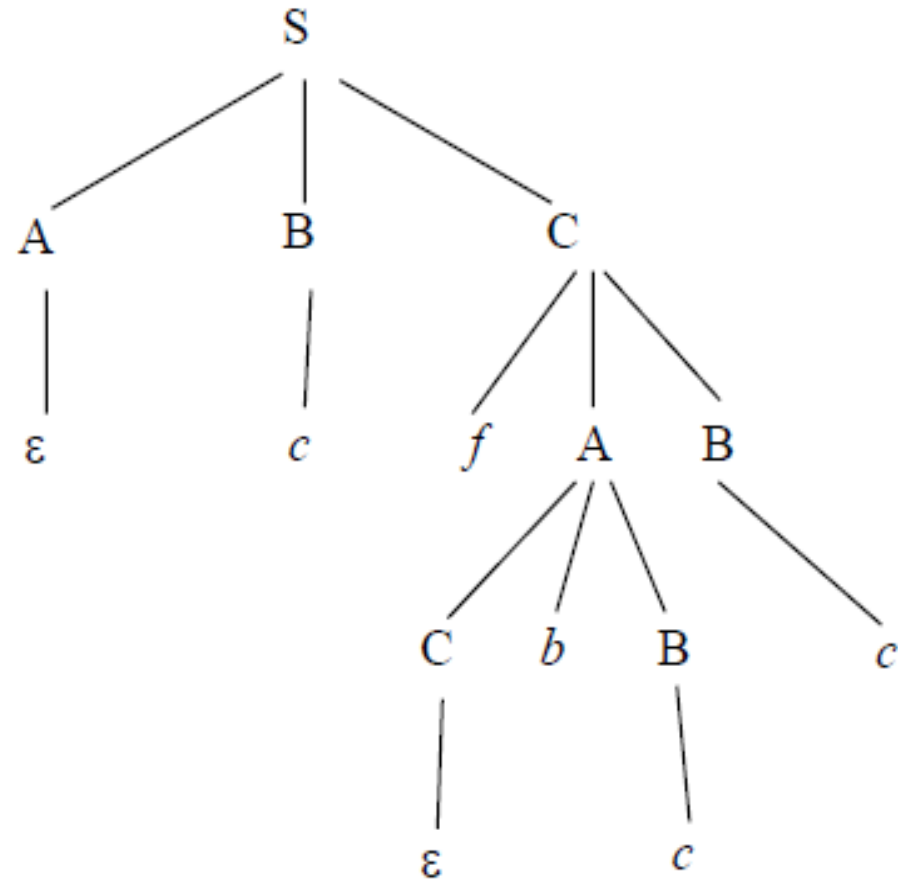


$S(a\ b\ S(b\ A(a\ b)\)\ B(c)\)$

Представление цепочки в виде дерева

■ $S(A()B(c)C(fA(C())bB(c))B(c)))$

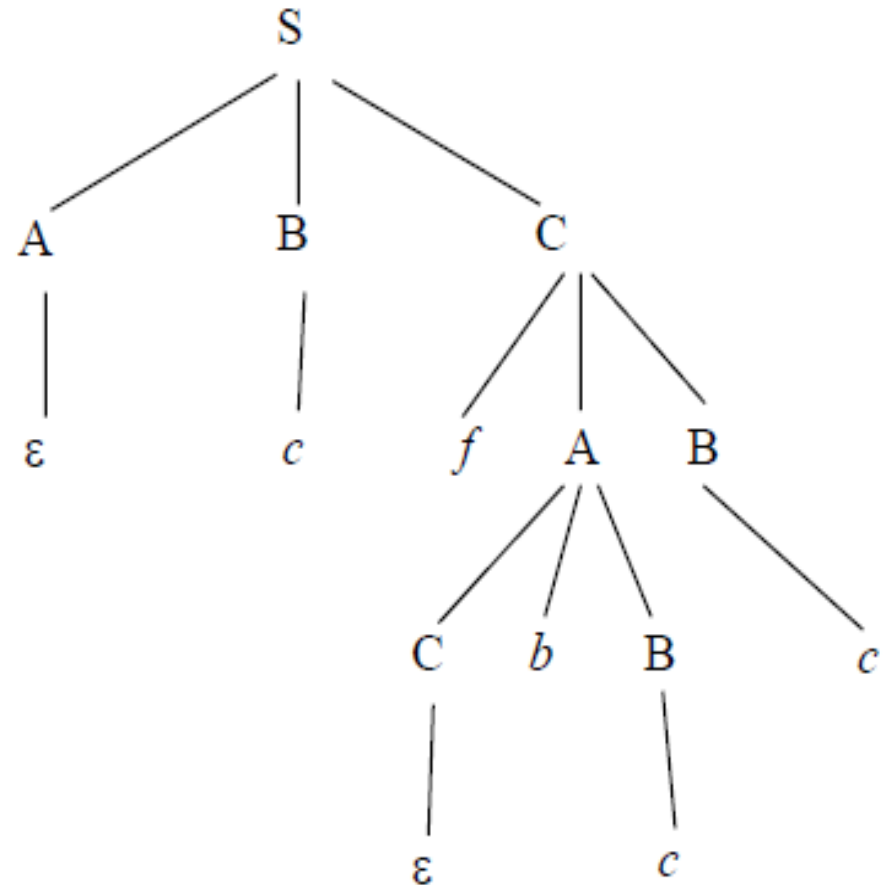
- Полученная на выходе при трансляции последовательность символов представляет структуру входной цепочки – дерево синтаксического анализа в скобочной форме
- При необходимости в качестве семантических действий могут быть выбраны такие, которые представляют дерево в списочной форме, в форме рисунка и т.п.
- Таким образом, само дерево синтаксического разбора ВИРТУАЛЬНО всегда присутствует как результат последовательности применения правил при синтаксическом анализе входной цепочки



Семантические вычисления при нисходящем разборе

■ $S(A()B(c)C(fA(C())bB(c))B(c)))$

- Трансляция при нисходящем синтаксическом анализе проводится семантическими процедурами, которые обычно вставляются непосредственно в правые части правил грамматики
- При синтаксическом анализе правые части грамматики заталкиваются в магазин вместе с семантическими процедурами на основании таблицы решений. Семантические процедуры выполняются при их выталкивании из магазина
 - Фактически, входная цепочка просто управляет последовательностью выполнения этих семантических процедур при синтаксическом анализаторе, и для каждой входной цепочки такая последовательность будет своей



Обобщение s-грамматик (для $k=1$)

a lookup table

Пример

1. $S \rightarrow bSbAc$
2. $S \rightarrow AbBc$
3. $A \rightarrow aSc$
4. $A \rightarrow cABc$
5. $B \rightarrow a$
6. $B \rightarrow dD$
7. $D \rightarrow Ba$

Правая часть правила 7 начинается не с терминала. Но других альтернатив для D нет!

	a	b	c	d
S	2	1	2	ош
D	7	ош	ош	7

Правая часть правила 2 начинается не с терминала
Но все цепочки, которые можно вывести из цепочки $AbBc$ начинаются с терминалов $\{a, c\}$. Следовательно, если очередной нетерминал S , то можно однозначно сделать выбор правила, по которому нужно S заменить, анализируя только один очередной терминал на входе. Множество выбора для правила 1 есть $\{b\}$, множество выбора для правила 2 есть $\{a, c\}$. Они не пересекаются

Пример

0. $S' \rightarrow S\$$
1. $S \rightarrow bSaAc$
2. $S \rightarrow \varepsilon$
3. $A \rightarrow bSc$
4. $A \rightarrow cABc$
5. $B \rightarrow aBc$
6. $B \rightarrow bc$

Правая часть правила 2 – пустая цепочка.
Можно ли однозначно определить, следует ли S заменять пустой цепочкой?
Для этого определим множество выбора правила $N \rightarrow \alpha$:

множество выбора правила $N \rightarrow \alpha$ – это терминалы, с которых в sententialных формах могут продолжаться цепочки, если нетерминал N заменим на α .

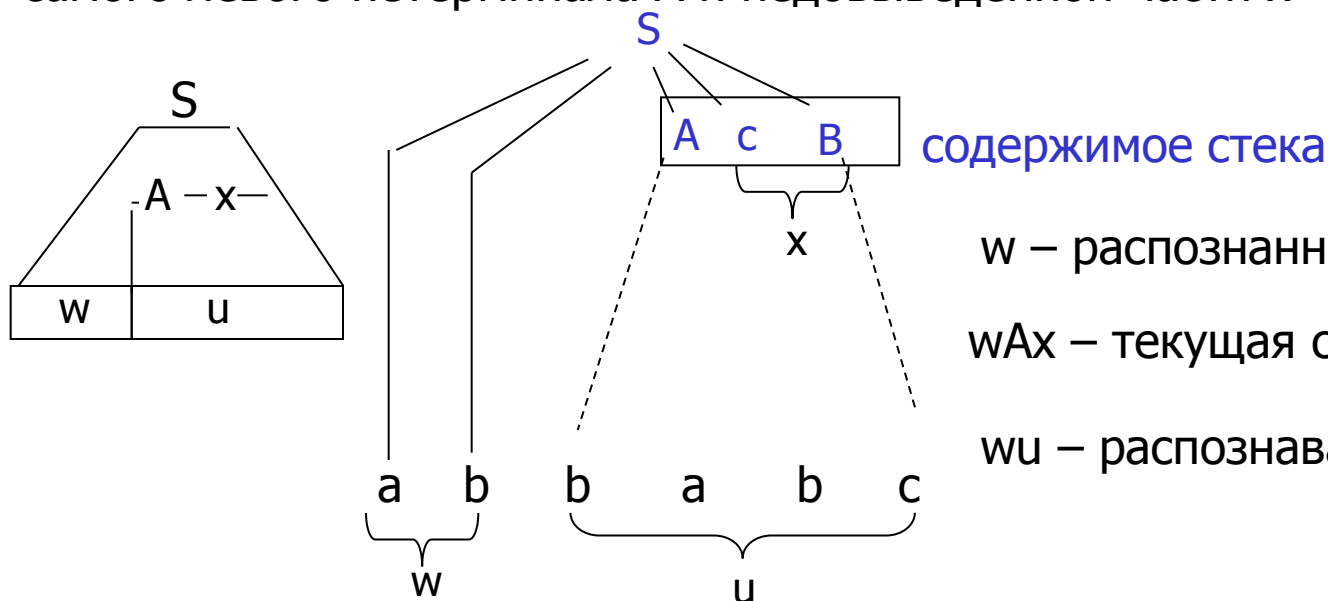
Множество выбора правила $S \rightarrow \varepsilon$ – это $\{\$, a, c\}$. Оно не пересекается с множеством выбора для правила 1, которое есть $\{b\}$

a lookup table

	a	b	c	\$
S	2	1	2	2

LL(k) языки и грамматики

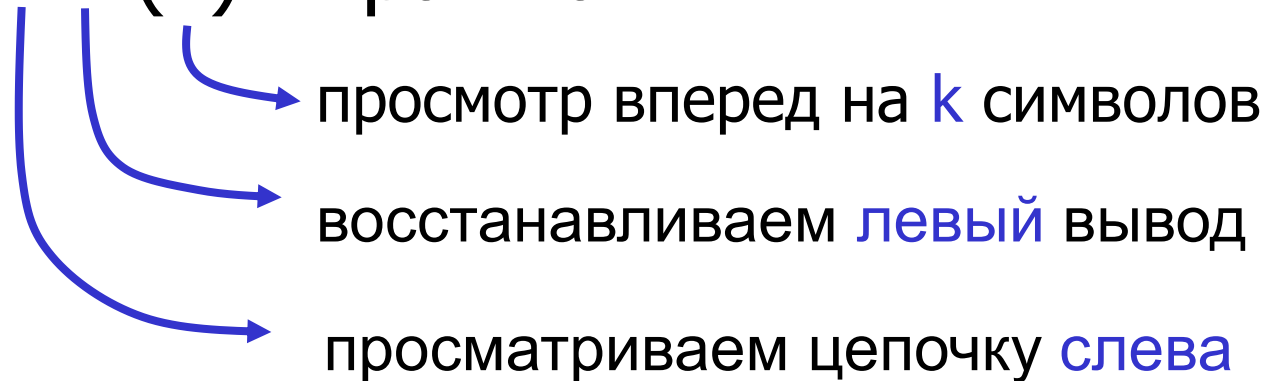
Построение дерева вывода в процессе восстановления левого вывода цепочки. Промежуточная цепочка в процессе вывода состоит из цепочки терминалов wu , самого левого нетерминала A и недовыведенной части x



Для продолжения разбора требуется заменить нетерминал A по одному из правил вида $A \rightarrow \alpha$. Чтобы разбор был детерминированным (без возвратов), правило требуется выбирать однозначно на основе анализа оставшейся части входной цепочки

Грамматика имеет свойство LL(k), если для однозначного выбора правила достаточно рассмотреть только Ax и первые k символов непросмотренной цепочки u . Обычно используется $k=1$, т.е. выполнить анализ одного очередного терминального символа оставшейся части входной цепочки u

L L (k) - грамматики



- первая буква L в названии LL(k): входная цепочка читается **слева** направо
- вторая L: восстанавливается **левый** вывод входной цепочки
- k означает, что на каждом шаге для принятия решения используется k следующих символов непрочитанной части входной цепочки (**look ahead**, заглядывание вперед по цепочке)

Теория LL(1)-грамматик

Функция **FIRST**(α) – множество терминальных символов, с которых могут начинаться цепочки, выводимые из цепочки α (цепочка α может содержать и нетерминальные символы). Если ε выводится из α , то ε включается в **FIRST**(α).

Если $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$, то выбор между правилами $A \rightarrow \alpha$ и $A \rightarrow \beta$ однозначен

Функция **FOLLOW**(A) – множество терминальных символов, с которых могут начинаться цепочки, которые **следуют** за A в любых sentential формах. Если A является последним символом sentential формы, то в **FOLLOW**(A) включается \$.

Если $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \emptyset$, то выбор между правилами $A \rightarrow \alpha$ и $A \rightarrow \varepsilon$ однозначен

- Грамматика имеет свойство LL(k), если из существования двух цепочек левых выводов:
 - $S \Rightarrow^* wAx \Rightarrow w\alpha x \Rightarrow^* wu$ (если A заменяется по правилу $A \rightarrow \alpha$)
 - $S \Rightarrow^* wAx \Rightarrow w\beta x \Rightarrow^* wv$ (если A заменяется по правилу $A \rightarrow \beta$)
 - из условия $\text{FIRST}_k(u) = \text{FIRST}_k(v)$ следует $\alpha = \beta$.
 - $\text{FIRST}_k(u)$ – терминальные цепочки длиной k, обобщение функции **FIRST**(u)
- В случае k=1, для выбора альтернативы для A достаточно знать только нетерминал A и терминал a - очередной символ входной цепочки (первый символ цепочки u):
 - следует выбрать правило $A \Rightarrow \alpha$, если a входит в $\text{FIRST}_1(\alpha)$
 - следует выбрать правило $A \Rightarrow \varepsilon$, если a входит в $\text{FOLLOW}_1(A)$

Функция $\text{FIRST}(\alpha)$

■ $\text{FIRST}(\alpha)$

- α - произвольная строка из $(T \cup N)^*$ (т.е. строка из терминальных и нетерминальных символов)
- Результат – множество терминалов
- Если $\alpha \Rightarrow^* a\beta$, где a – терминал, то $a \in \text{FIRST}(\alpha)$

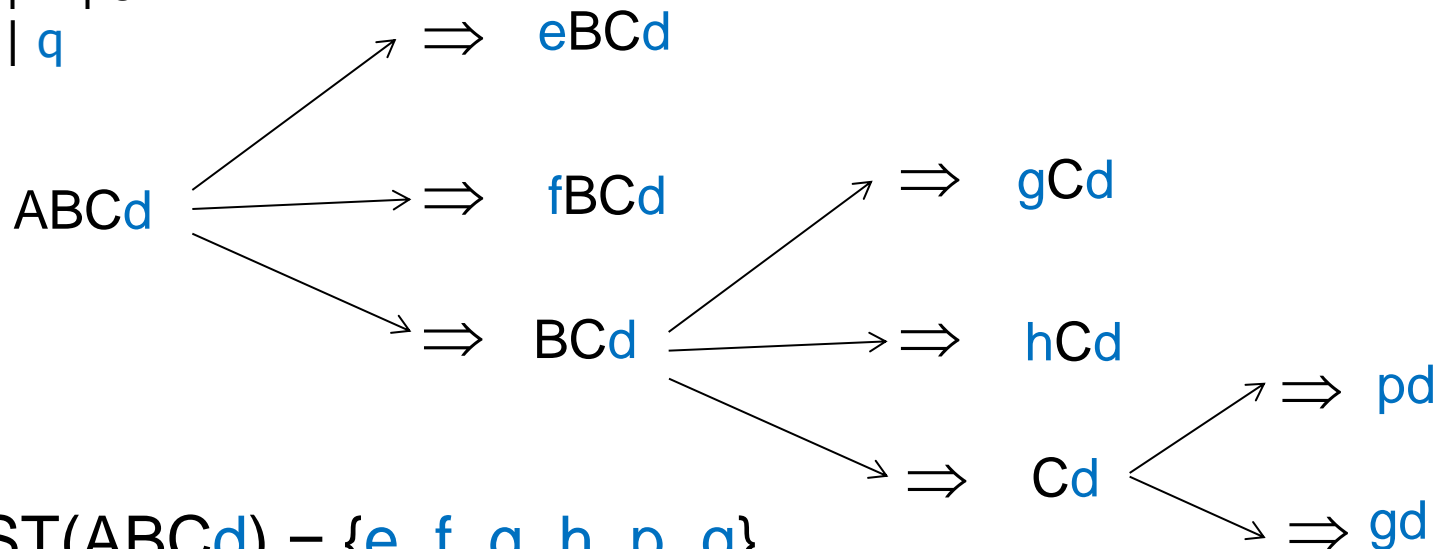
$S \rightarrow ABCd$

$A \rightarrow e \mid f \mid \varepsilon$

$B \rightarrow g \mid h \mid \varepsilon$

$C \rightarrow p \mid q$

$\text{FIRST}(ABCd)$



$\text{FIRST}(ABCd) = \{e, f, g, h, p, q\}$

Множество LOOKAHEAD

LL(1)-анализаторы просматривают поток только на один символ вперед при принятии решения о том, какое правило грамматики необходимо применить

Иногда строят LOOKAHEAD таблицу – множество выбора конкретной альтернативы для нетерминала (правила) по очередному терминальному символу

$\text{LOOKAHEAD}(A \rightarrow \alpha) = \text{FIRST}(\alpha)$, если неверно, что $\alpha \Rightarrow^* \varepsilon$

$\text{LOOKAHEAD}(A \rightarrow \alpha) = \text{FIRST}(\alpha) \cup \text{FOLLOW}(A)$, если $\alpha \Rightarrow^* \varepsilon$

Пример:

0. $S' \rightarrow S\$$
1. $S \rightarrow AbB$
2. $S \rightarrow d$
3. $A \rightarrow CAb$
4. $A \rightarrow B$
5. $B \rightarrow cSd$
6. $B \rightarrow \varepsilon$
7. $C \rightarrow a$
8. $C \rightarrow ed$

Вычисление множеств выбора для $k=1$

1. $\text{LOOKAHEAD}(S \rightarrow AbB) = \text{FIRST}(AbB) = \{a, b, c, e\}$
2. $\text{LOOKAHEAD}(S \rightarrow d) = \text{FIRST}(d) = \{d\}$
3. $\text{LOOKAHEAD}(A \rightarrow CAb) = \text{FIRST}(CAb) = \{a, e\}$
4. $\text{LOOKAHEAD}(A \rightarrow B) = \text{FIRST}(B) \cup \text{FOLLOW}(A) = \{c, b\}$
5. $\text{LOOKAHEAD}(B \rightarrow cSd) = \text{FIRST}(cSd) = \{c\}$
6. $\text{LOOKAHEAD}(B \rightarrow \varepsilon) = \text{FOLLOW}(B) = \{b, d, \$\}$
7. $\text{LOOKAHEAD}(C \rightarrow a) = \text{FIRST}(a) = \{a\}$
8. $\text{LOOKAHEAD}(C \rightarrow ed) = \text{FIRST}(ed) = \{e\}$

Lookahead table

- LOOKAHEAD ($A \rightarrow X_1 X_2 \dots X_n$) = $\cup \{ \text{FIRST}(X_i) \mid \text{EMPTY}(X_1 X_2 \dots X_{i-1}) \}$
 \cup if $\text{EMPTY}(X_1 X_2 \dots X_n)$ then FOLLOW(A) else \emptyset

Пример.

$G:: S' \rightarrow S\$$
 $S \rightarrow BA \mid AAd$
 $A \rightarrow a \mid \varepsilon$
 $B \rightarrow bA \mid cB$

N	FIRST(N)	FOLLOW(N)
S'	{ a, b, c, d }	{ }
S	{ a, b, c, d }	{ \$ }
A	{ a }	{ a, d, \$ }
B	{ b, c }	{ a, \$ }

Правила	LOOKAHEAD
$S' \rightarrow S\$$	{ a, b, c, d, \$ }
$S \rightarrow BA$	{ b, c }
$S \rightarrow AAd$	{ a, d }
$A \rightarrow a$	{ a }
$A \rightarrow \varepsilon$	{ a, d, \$ }
$B \rightarrow bA$	{ b }
$B \rightarrow cB$	{ c }

Теорема. КС-грамматика G является LL(1)-грамматикой, если и только если для любых двух альтернатив $A \rightarrow \alpha$ и $A \rightarrow \beta$ любого нетерминала A грамматики G :
 $\text{LOOKAHEAD}(A \rightarrow \alpha) \cap \text{LOOKAHEAD}(A \rightarrow \beta) = \emptyset$

Примеры построения таблицы решений для LL(1)

1. $S' \rightarrow S\$$
2. $S \rightarrow aSb$
3. $S \rightarrow \varepsilon$

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

Грамматика LL(1)

Таблица решений

	a	b	\$
S'	1	ош	1
S	2	3	3

1. $S' \rightarrow S\$$
2. $S \rightarrow F$
3. $S \rightarrow (S + F)$
4. $F \rightarrow i$

$\text{FIRST}(F) = \{i\}$
 $\text{FIRST}(S+F) = \{(\}$

Грамматика LL(1)

Пустые клетки – это ошибки, их
обычно не пишут в таблице

	()	+	i	\$
S'	1			1	
S	3			2	
F				4	

Грамматика, порождающая подмножество типов языка Pascal:

1. $\langle \text{type} \rangle \rightarrow \langle \text{simple} \rangle$
2. $\langle \text{type} \rangle \rightarrow \wedge \text{id}$
3. $\langle \text{type} \rangle \rightarrow \text{array} [\langle \text{simple} \rangle] \text{ of } \langle \text{type} \rangle$
4. $\langle \text{simple} \rangle \rightarrow \text{integer}$
5. $\langle \text{simple} \rangle \rightarrow \text{char}$
6. $\langle \text{simple} \rangle \rightarrow \text{num} \text{ .. num}$

	\wedge	<u>array</u>	<u>id</u>	<u>ch</u>	<u>num</u>
$\langle \text{type} \rangle$	2	3	1	1	1
$\langle \text{simple} \rangle$	ош	ош	4	5	6

Здесь ошибки указаны явно

JFLAP: <http://www.jflap.org/> учебник по LL(1) синтаксическому анализу

Пример: прямая левая рекурсия

1. $E \rightarrow E+T$
2. $E \rightarrow T$
3. $T \rightarrow T * P$
4. $T \rightarrow P$
5. $P \rightarrow i$
6. $P \rightarrow (E)$

$A \rightarrow \alpha$	$FIRST(\alpha)$
1. $E \rightarrow E+T$	$\{ i, (\}$
2. $E \rightarrow T$	$\{ i, (\}$
3. $T \rightarrow T * P$	$\{ i, (\}$
4. $T \rightarrow P$	$\{ i, (\}$
5. $P \rightarrow i$	$\{ i \}$
6. $P \rightarrow (E)$	$\{ (\}$

Грамматика не является LL(1), т.к. правила для E и T содержат прямую левую рекурсию

$$E \rightarrow E+T \mid T$$

заменяем:

$$E \rightarrow TE_1 \quad E_1 \rightarrow +TE_1 \mid \varepsilon$$

или в БНФ:

$$E ::= T \{+T\}$$

0. $S \rightarrow E \$$
1. $E \rightarrow TE_1$
2. $E_1 \rightarrow +TE_1$
- 2a. $E_1 \rightarrow \varepsilon$
3. $T \rightarrow PT_1$
4. $T_1 \rightarrow *T_1$
- 4a. $T_1 \rightarrow \varepsilon$
5. $P \rightarrow i$
6. $P \rightarrow (E)$

Правило	LOOKAHEAD
1. $E \rightarrow TE_1$	$\{ i, (\}$
2. $E_1 \rightarrow +TE_1$	$\{ + \}$
2a. $E_1 \rightarrow \varepsilon$	$\{), \$ \}$
4. $T_1 \rightarrow *T_1$	$\{ * \}$
4a. $T_1 \rightarrow \varepsilon$	$\{ +,), \$ \}$
5. $P \rightarrow i$	$\{ i \}$
6. $P \rightarrow (E)$	$\{ (\}$

таблица решений

	+	*	i	()	\$
S			0	0		
E			1	1		
E ₁	2				2a	2a
T ₁	4a	4			4a	4a
P			5	6		

По 2 альтернативы - три символа, P, E₁ и T₁

Пустые клетки – это ошибки

Пример: трансляция оператора присваивания

0. $S \rightarrow i:=E \$$ $y_0(i)$
1. $E \rightarrow TE_1$
2. $E_1 \rightarrow +TE_1$ y_2
- 2a. $E_1 \rightarrow \varepsilon$
3. $T \rightarrow PT_1$
4. $T_1 \rightarrow *T_1$ y_4
- 4a. $T_1 \rightarrow \varepsilon$
5. $P \rightarrow i$ $y_5(i)$
6. $P \rightarrow (E)$

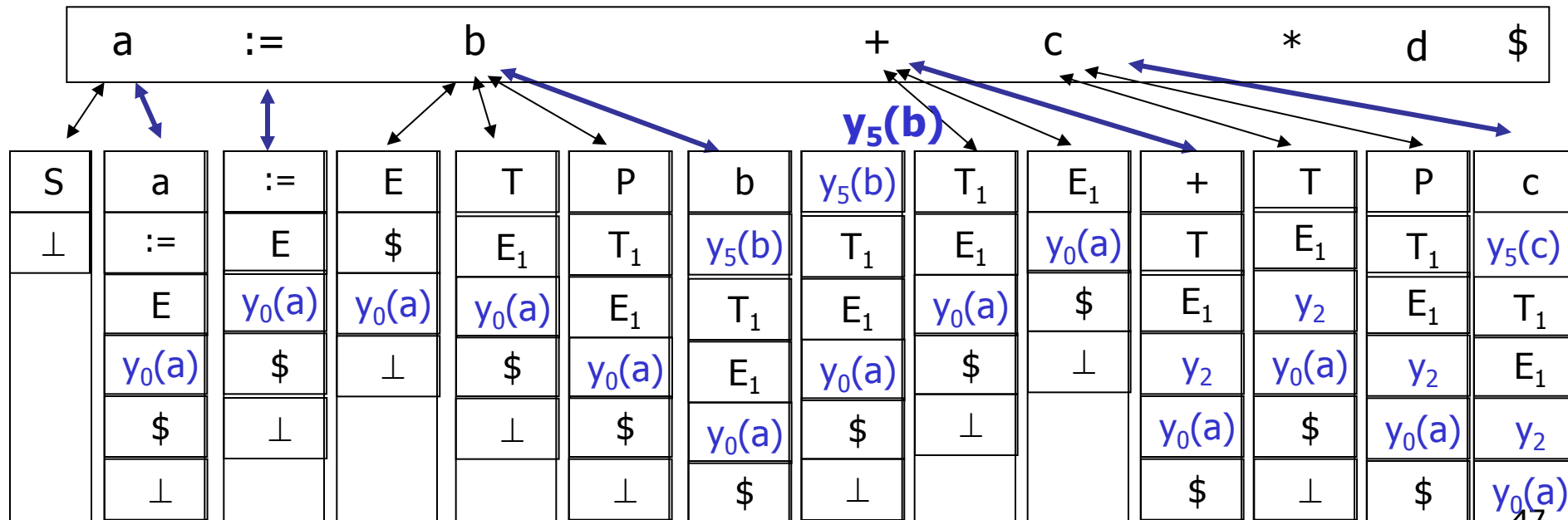
	+	*	i	()	\$
E_1	2				2a	2a
T			3	3		
T_1	4a	4			4a	4a
P			5	6		

Пустые клетки – это ошибки

Трансляция
 $a:=b+c*d \$$
 должна дать:
 Выберем
 семантики:

LOAD addr b
 LOAD addr c
 LOAD addr d
 MUL
 ADD
 STORE addr a

$y_0(i)$: Gen STORE addr i
 y_2 : Gen ADD
 y_4 : Gen MUL
 $y_5(i)$: Gen LOAD addr i

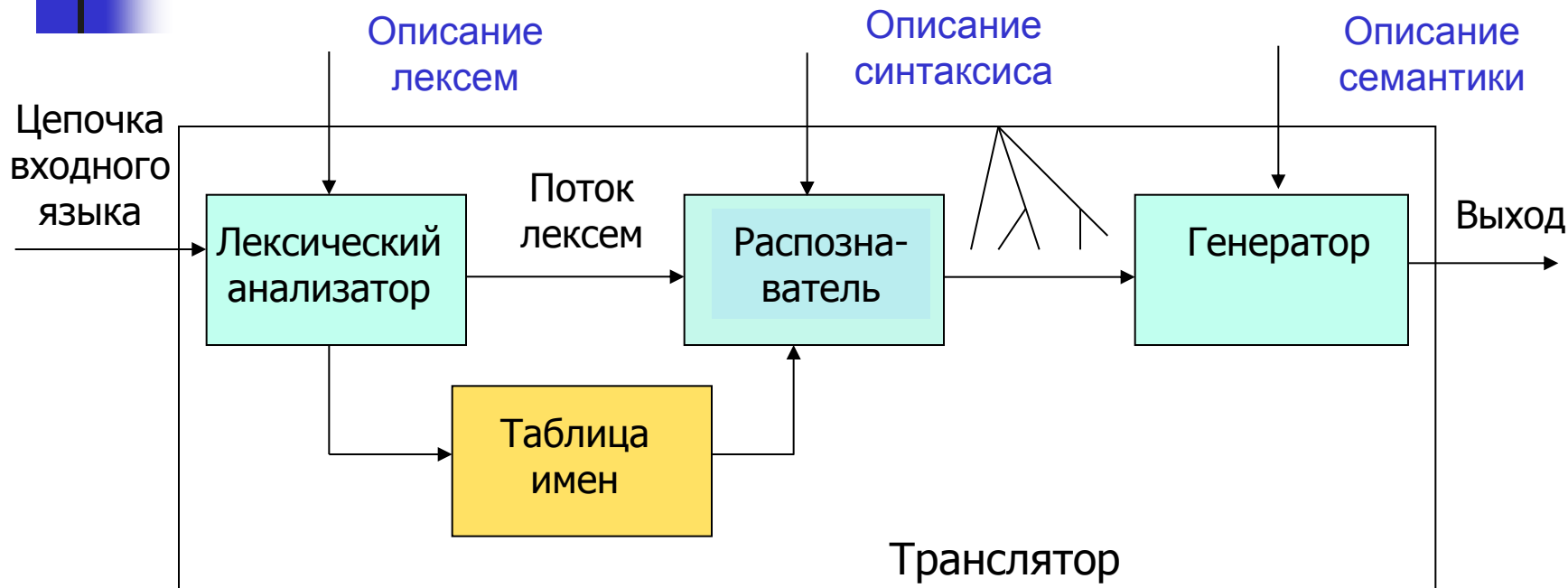


- Семантические вычисления при нисходящем синтаксическом анализе (LL-анализатор) выполняются просто
- Наследуемые атрибуты символов правой части правила сразу вычисляются при замене нетерминала по решающей таблице

Действия СА при обнаружении ошибки

- Что, если в таблице решений синтаксический анализатор попадает на “ошибку”? Это означает, что входная строка НЕ выводится в данной грамматике, иными словами **во входной строке синтаксическая ошибка**
- При обнаружении ошибки синтаксический анализатор может просто остановиться и указать место возникновения ошибки. **Эта стратегия неудобна, поскольку при длинной программе количество повторных запусков на трансляцию может быть большим**
- Другая стратегия:
 - 1) определяем множество “маяков” (синхронизирующих лексем), которые ограничивают операторы, блоки, процедуры и т.п.
 - 2) при обнаружении синтаксической ошибки СА:
 - а) пропускает лексемы входной строки по одной, пока не будет найден маяк, и
 - б) управление в анализаторе передается в ближайшую точку анализатора, после которой встречается эта лексема
- Обычно в множество синхронизирующих лексем включают разделители, например, **;;**, **)**, **end** или **}**.. . Набор таких лексем определяется разработчиком анализируемого языка, для которого строится компилятор
- При такой стратегии восстановления может оказаться, что много лексем будет пропущено без проверки на наличие дополнительных ошибок
- Такая стратегия восстановления наиболее проста в реализации
- Существуют и другие стратегии восстановления СА после обнаружения ошибки во входном тексте
- **Задача для курсовой работы: почитать о стратегиях восстановления анализа после обнаружения синтаксической ошибки и реализовать ее в компиляторе для Милана**

Генератор компиляторов

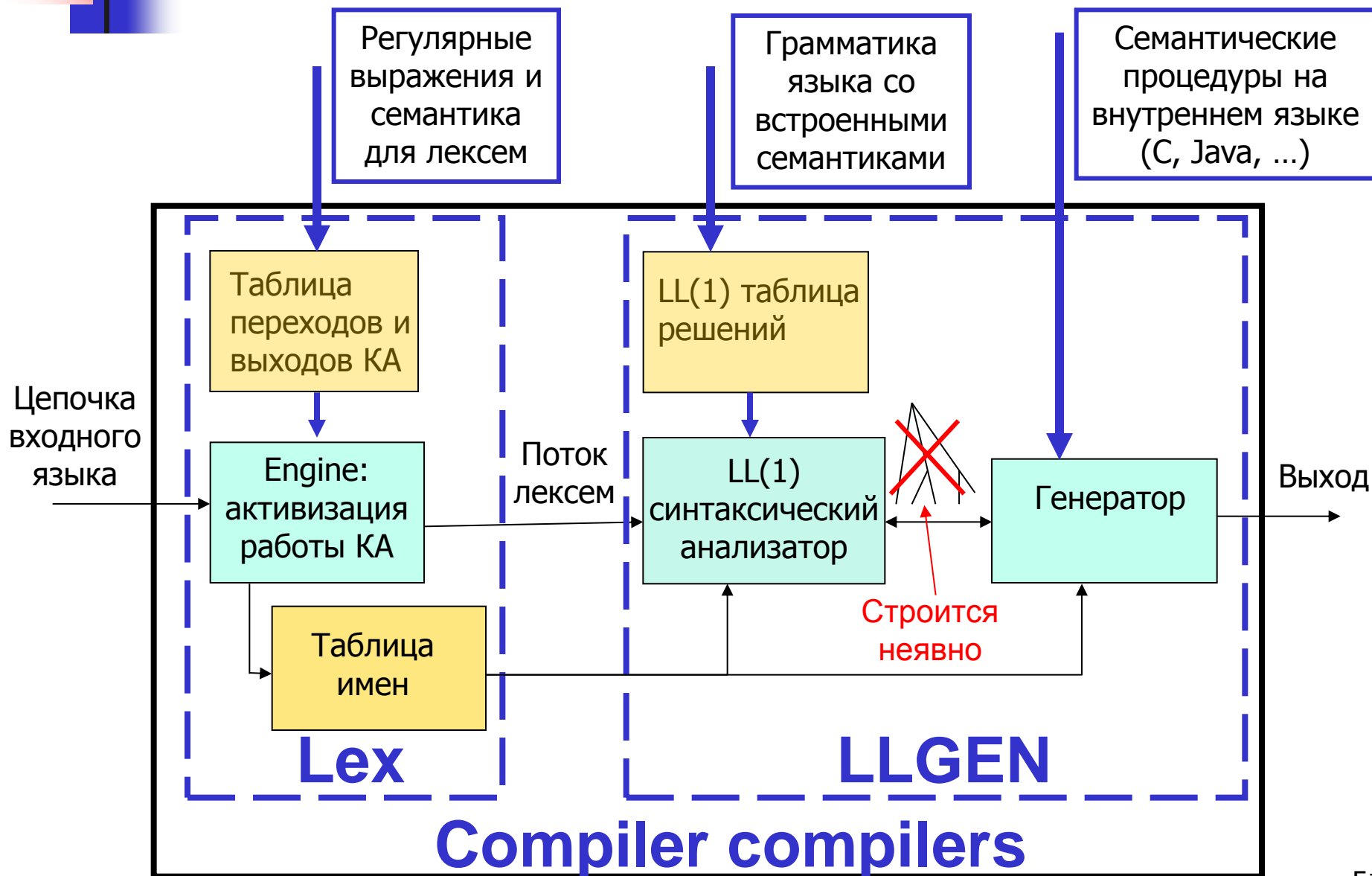


Можно ли построить такую программную систему для любой грамматики?

■ Да, если:

- задать описание лексических единиц (лексем)
- задать описание КС-грамматики исходного языка (ограничиться классом КС-грамматик, для которого есть алгоритм СА)
- задать описание семантики, связанной с правилами грамматики

Компилятор компиляторов для LL-грамматик



Генераторы лексических и синтаксических анализаторов

<http://www.kulichki.com/kit/tools/lexparse.html>

- **LEX** - классический генератор лексических анализаторов AT&T, поставляемый с Unix
- **YACC (Yet Another Compiler Compiler)** - классический генератор синтаксических анализаторов AT&T, поставляемый с Unix
- **FLEX** - GNU версия генератора сканеров Lex
- **BISON** - GNU версия Yacc
- **ACCENT** компилятор компиляторов, не накладывающий никаких ограничений на грамматики: никакой адаптации к специфическим методам синтаксического анализа, таким как LL(k) или LALR(k), не требуется. Поддерживает расширенную БНФ
- **AFLEX & AYACC** - аналогичны юниксовым инструментам Lex и Yacc, но написаны на Ада и генерируют на выходе компилятор на языке Ада
- **ANAGRAM** - генератор синтаксических анализаторов LALR с возможностью автоматического восстановления после синтаксических ошибок
- **BTYACC** - модифицированная версия Yacc, поддерживающая автоматический откат и семантическое устранение неоднозначности для разбора неоднозначных грамматик
- **BYACC** (Berkeley Yacc) - генератор синтаксических анализаторов LALR(1)
- **COCO/R** генерирует анализаторы, работающие по методу рекурсивного спуска и связанные с ними сканеры по атрибутным грамматикам
- **SCANGEN, LLGEN, LALRGEN** - генераторы лексических, LL(1) и LALR(1) генераторов
- **TP LEX AND YACC** - генератор лексических и синтаксических анализаторов для Turbo Pascal

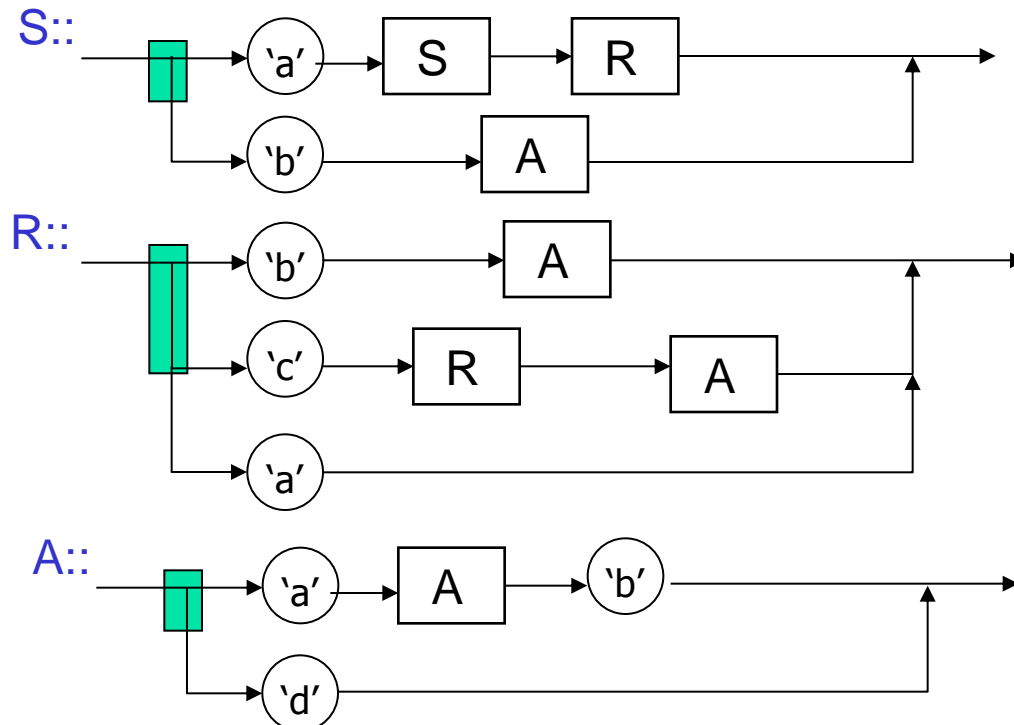
... .. Подобных генераторов десятки (если не сотни ...)



Грамматики рекурсивного спуска

Определение

- Анализатором рекурсивного спуска называется вариант предсказывающего LL(1) анализа, в котором каждому нетерминалу сопоставляется процедура (вообще говоря, рекурсивная), и стек образуется неявно при вызовах таких процедур
- Грамматикой рекурсивного спуска называется LL(1)-грамматика, которая представляется синтаксическими диаграммами для каждого нетерминала такими, что в любом разветвлении каждой синтаксической диаграммы выбор альтернативного пути может быть сделан на основе анализа одного очередного терминального символа входной цепочки
- Очевидно, что любая s-грамматика является грамматикой рекурсивного спуска

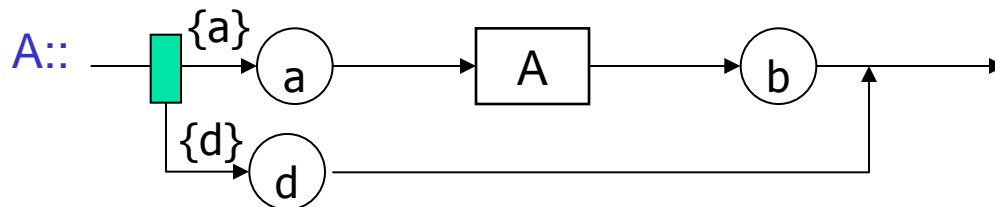
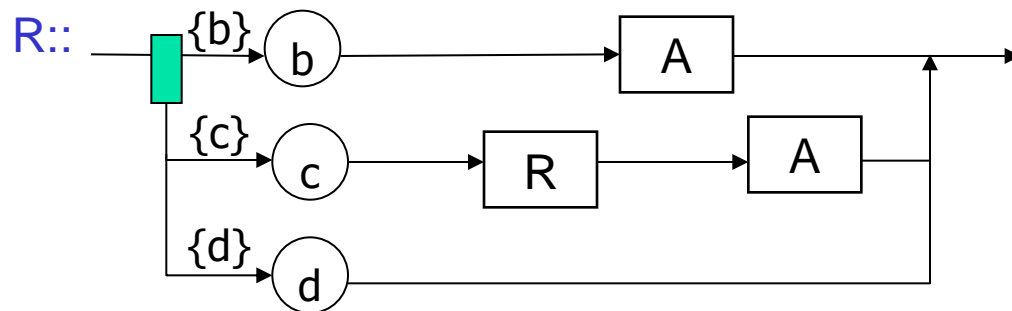
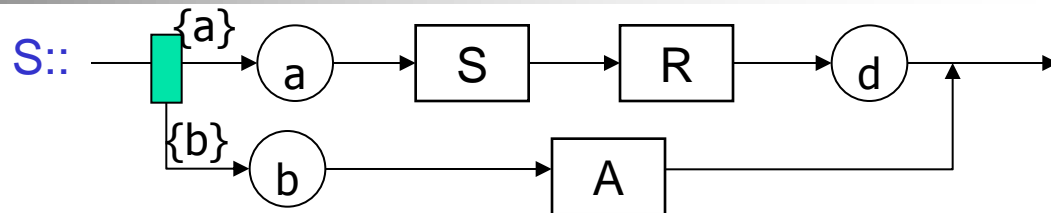


1. $S \rightarrow aSR$
2. $S \rightarrow bA$
3. $R \rightarrow bA$
4. $R \rightarrow cRA$
5. $R \rightarrow a$
6. $A \rightarrow aAb$
7. $A \rightarrow d$

Пример: s-грамматика и грамматика рекурсивного спуска

1. $S \rightarrow aSRd$
2. $S \rightarrow bA$
3. $R \rightarrow bA$
4. $R \rightarrow cRA$
5. $R \rightarrow d$
6. $A \rightarrow aAb$
7. $A \rightarrow d$

В фигурных скобках - правило выбора пути по очередному терминалу



- Для каждого нетерминала строится своя распознающая процедура
- Альтернативные пути в синтаксической диаграмме помечены различными терминалами, следовательно при анализе цепочки путь вычислений в процедурах будет выбран однозначно по результату анализа одного очередного терминального символа
- При синтаксическом анализе для таких грамматик стек не строится. Стек образуется неявно при рекурсивных вызовах процедур в процессе анализа входной цепочки

Идея алгоритма рекурсивного спуска

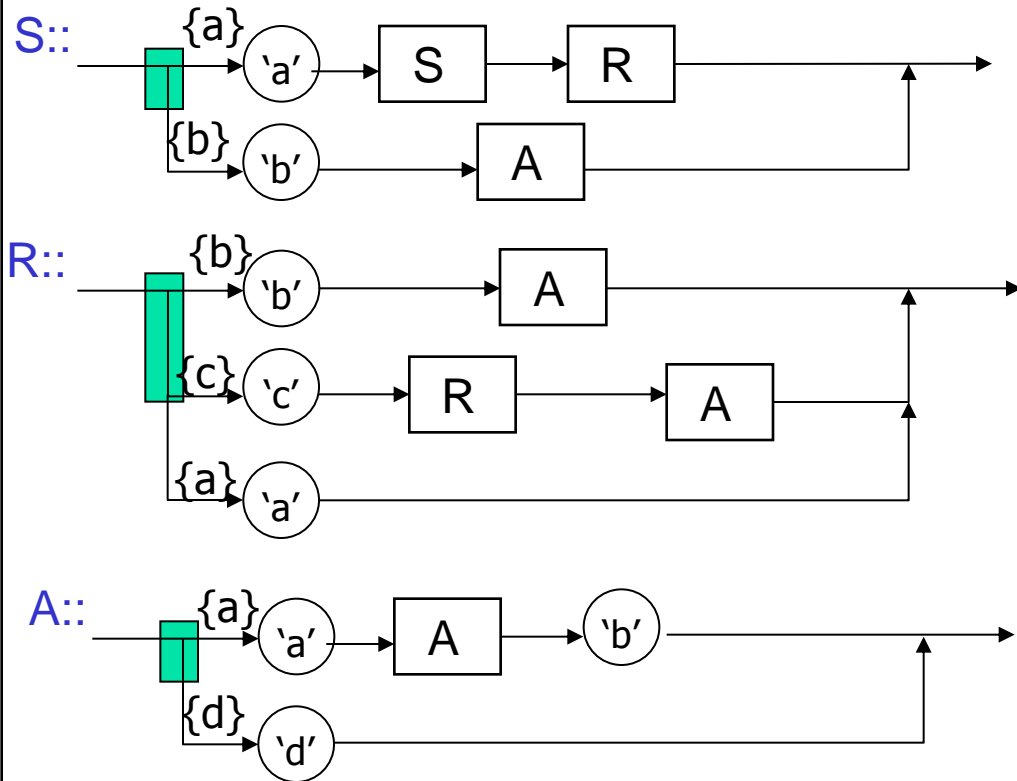
```
int Yk:=0; // указ-ль на очередной символ
```

```
S () {  
  if( s[Yk] == 'a') then { Yk++; S(); R(); }  
  else  
    { if( s[Yk] == 'b') then { Yk++; A(); }  
      else { Error1( ); }  
    }  
}
```

```
A () {  
  if( s[Yk] == 'a') then {  
    Yk++; A();  
    if( s[Yk] == 'b') then { Yk++; break; }  
    else { Error2( ); }  
  }  
  else { if( s[Yk] == 'd') then { Yk++; break; }  
        else { Error3( ); }  
  }  
}
```

При входе в процедуру N Yk установлен на первый символ подцепочки, выводящейся из нетерминала N

Ю.Г.Карпов

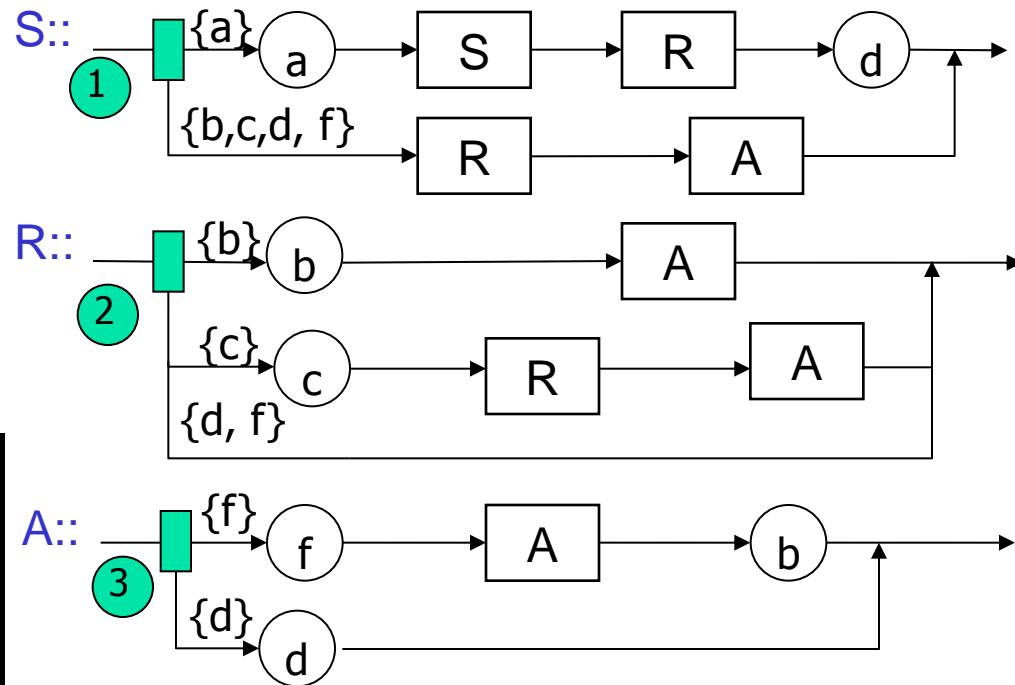


- Каждому нетерминалу сопоставляется процедура. Начальному нетерминалу – main
- Каждая процедура распознает максимальную подстроку своего языка

Пример грамматики рекурсивного спуска (не s-грамматики)

1. $S \rightarrow aSRd$
2. $S \rightarrow RA$
3. $R \rightarrow bA$
4. $R \rightarrow cRA$
5. $R \rightarrow \varepsilon$
6. $A \rightarrow fAb$
7. $A \rightarrow d$

В методе рекурсивного спуска поскольку одна из альтернатив R – пустой символ, нужно проверить, какой символ идет после нетерминала R (Follow(R)) и проверить, пересекаются ли множества



- Выбор 3: определяется очередным терминалом f или d . Выбор однозначен
- Выбор 2: третий путь определяется терминалами Follow(R), т.е. теми терминалами, которые могут встретиться **после** R. Это множество $\{f, d\}$ тех, с которых начинаются цепочки, выводимые из A, поскольку после R стоит только A. Выбор пути однозначен
- Выбор 1: второй путь определяется терминалами $\{b, c, f, d\}$, т.е. теми, с которых могут начаться цепочки, выводимые из R, и теми, которые могут встретиться **после** R, поскольку есть правило $R \rightarrow \varepsilon$. Выбор однозначен

Распознаватель языка перечисления номеров страниц

```
int Ук:=0; // указ-ль на очередной символ
```

```
main( ) {
```

```
  m: Группа( );
```

```
    if( s[Ук] == 'eot') then { Ук++; break };
    else if( s[Ук] == ',') then { Ук++; goto m };
    else Error1( );
```

```
Группа( ) {
```

```
  Целое( );
```

```
  if( s[Ук] == '-') then { Ук++; Целое( ); }
  else skip
```

```
};
```

```
Целое( ) {
```

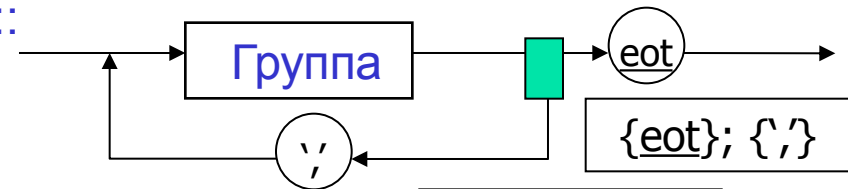
```
  if( s[Ук] == 'ц') then Ук++;
```

```
  else Error2( );
```

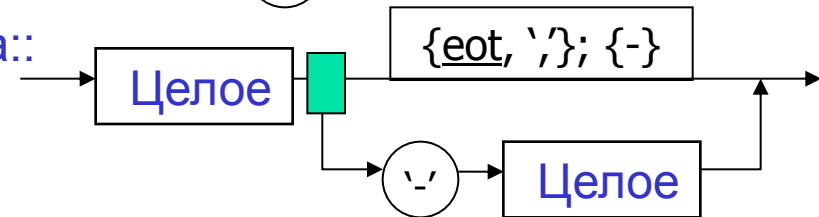
```
  m: if( s[Ук] == 'ц') then { Ук++; goto m };
  else skip
```

```
};
```

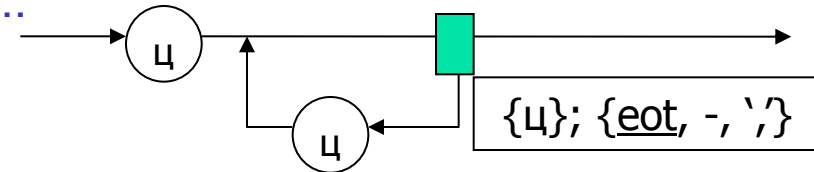
Прог::



Группа::



Целое::



Для каждого нетерминала – своя
распознающая процедура

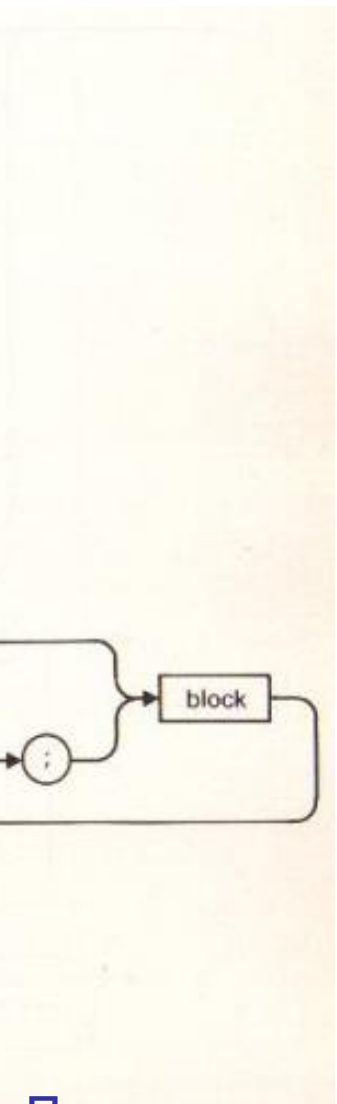
Эта грамматика рекурсивного спуска:
в каждом разветвлении синтаксической
диаграммы по очередному терминалу
путь определяется однозначно

Добавление семантических операций в синтаксическую диаграмму позволяет
построить транслятор



Lookahead символы

- Стандартный метод, примененный в трансляции рекурсивного спуска, заключается в использовании одиночного, упреждающего символа, объявленного глобально (Lookahead символа)
- В рамках процедуры инициализации обязательно выполняется функция, считывающая первый символ из входного потока, например, `GetChar()`. Каждый вызов функции `GetChar()` или аналогичной дает следующий символ из входного потока



- HO



Семантические вычисления при рекурсивном спуске

- При распознавании цепочки могут выполняться семантические процедуры
- Синтезируемые атрибуты нетерминалов являются определяемыми параметрами при вызове соответствующих процедур.
- Наследуемые атрибуты символов являются параметрами, устанавливаемыми при вызове соответствующей процедуры
- Пример: генерация кода. Пусть нетерминал имеет два атрибута:
N.Снач – начальный адрес фрагмента программы, который вычисляет N,
N.Скон – адрес, следующий за фрагментом, вычисляющим N.

Тогда С можно считать не глобальным, а локальным параметром, в main устанавливается Снач=0, а при вызовах текущее С передается как параметр Снач, а после вызова процедуры очередного нетерминала локальное С устанавливается по значению Скон.

- Пример – язык линейных алгебраических уравнений. Там как параметры передаются наследуемые атрибуты (например, номер обрабатываемого уравнения)



Заключение

- Задачей синтаксического анализа КС-языка является восстановление вывода (дерева вывода) цепочки языка из начального символа грамматики. Теория синтаксического анализа является наиболее разработанной областью информатики
- Существуют универсальные неэффективные алгоритмы синтаксического анализа КС-языков, для подклассов КС-грамматик существуют эффективные (линейные) алгоритмы распознавания
- Прежде чем строить синтаксический анализатор, в грамматике следует выбросить неприменяемые правила
- Существует два класса алгоритмов синтаксического анализа КС-языков: **НИСХОДЯЩИЕ** и **ВОСХОДЯЩИЕ** (по синтаксическому дереву) алгоритмы. Нисходящие алгоритмы восстанавливают **ЛЕВЫЙ** вывод цепочки языка, восходящие – **ПРАВЫЙ** вывод
- Эффективные нисходящие алгоритмы синтаксического анализа для:
 - s-грамматик,
 - LL(k)-грамматик,
 - грамматик рекурсивного спуска
- Алгоритм рекурсивного спуска часто используется при разработке трансляторов специализированных языков



Спасибо за внимание