

Теория алгоритмических языков и трансляторов (09.03.04 – Программная инженерия)

Лектор:

Крючкова Елена Николаевна

Кафедра прикладной математики АлтГТУ

Что будем рассматривать?

Что такое ЯЗЫК?

Какие существуют ФОРМАЛЬНЫЕ СПОСОБЫ описания языков?

Что такое ЯЗЫК ПРОГРАММИРОВАНИЯ с точки зрения разработчика компилятора?

Как строится компилятор?

Чем программа компилятора отличается от программы интерпретатора?

Как реализовать интерпретатор или компилятор?

Структура курса (учебный план 09.03.04)

Семестр 16 недель

Лекции: 32 часа (1 пара в неделю)

Лабораторные работы: 32 часа (1 пара в неделю)

Экзамен

Литература

Крючкова Е.Н. Методы анализа в теории формальных языков. / Учебное пособие – 276 с.

Лабораторные работы

- 1 (2 часа) Синтез КС-грамматик
- 2 (2 часа) КС-грамматика языка программирования
- 3 (2 часа) Лексика (таблица лексем, конечный автомат лексики, конечный автомат игнорируемых символов)
- 4 (2 часа) Программа лексического анализатора (сканер)
- 5 (2 часа) Построение синтаксических диаграмм и их преобразование
- 6 (2 часа) Синтаксический анализатор (метод рекурсивного спуска)
- 7 (2 часа) Анализ контекстных условий языка программирования.
- 8 (2 часа) Реализация семантического дерева
- 9 (2 часа) Полный семантический анализ
- 10 (2 часа) *Данные интерпретатора*
- 11 (2 часа) *Запись и выборка данных при интерпретации*
- 12 (2 часа) *Интерпретация выражений*
- 13 (2 часа) *Работа с флагом интерпретации - проектирование*
- 14 (2 часа) *Работа с флагом интерпретации - реализация*
- 15 (2 часа) *Промежуточный код - проект конструкций*

Отчеты по работам

Отчет по каждой работе в электронном виде должен быть оформлен и сдан:

1. Отчет по каждой работе оформляется в виде одного файла *.txt или *.doc или *.docx.
2. Название любого такого файла стандартное:
<номер работы>_<номер группы>_<Фамилия студента>
3. Примеры: 06_92_Иванов.zip, 03_92_Ivanov.doc
4. Отчеты отсылаются на адрес
kruchkova_elen@mail.ru

Задания №1 и №2

Задание №1 (работа №1) – в учебнике.

Задание №2 (работы 2-15) - в файле в ЛК.

1. Лабораторная работа 1 – синтез КС-грамматик
2. Лабораторные работы №2-15 - методы трансляции языков программирования

Задание по трансляторам

В каждом задании описывается некоторый очень усеченный вариант известных языков программирования Java и C++. В задании указывается:

- структура программы,
- типы данных, которые могут использоваться,
- допустимые операции над этими данными,
- операторы,
- операции и операнды, из которых строятся выражения,
- все виды констант, которые могут использоваться в выражениях.

Пример задания по трансляторам

Программа: главная программа языка C++.
Допускается описание функций без параметров, функции возвращают значение.

Типы данных: int (в том числе short , long long, long) .

Операции: арифметические, сравнения.

Операторы: присваивания и do {} while().

Операнды: простые переменные и константы.

Константы: целые в 10 с/с и 16 с/с .

Обратите внимание на полноту реализации задания

1. во всех заданиях предполагается использование *составного и пустого оператора*,
2. всегда разрешается описание *глобальных данных*,
3. все *перечисленные элементы языка должны использоваться в программе* (например, если разрешается описание функций, то, безусловно, в перечень операторов Вам необходимо включить вызовы функций).

Тема 1

Формальные грамматики и языки

Базовые понятия

Алфавит - непустое конечное множество символов.

Цепочка над алфавитом Σ - конечная последовательность символов алфавита. Длина цепочки x - число ее символов, обозначается $|x|$. Цепочка нулевой длины называется пустой цепочкой и обозначается ε .

Множество всех цепочек (включая ε) над алфавитом Σ обозначается Σ^*

Язык L над алфавитом Σ - подмножество множества Σ^* , то есть $L \subseteq \Sigma^*$

Примеры алфавита

Алфавит $\Sigma = \{0,1\}$

$\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, \dots\}$

Алфавит $\Sigma = \{0,1,\dots,9,x,X,a,\dots,f,A,\dots,F\}$

$0Xff8a4 \in \Sigma^* \quad |0Xff8a4| = 7$

$aa4Xff8x \in \Sigma^*$

Множество часел в 16с/с – подмножество Σ^* ,
но не равно Σ^*

Алфавит языка программирования C++ - все
символы клавиатуры

Грамматика

Порождающая грамматика - упорядоченная четверка

$$G = (V_T, V_N, P, S), \text{ где}$$

V_T - алфавит терминальных символов;

V_N - алфавит нетерминальных символов;

P --- конечное множество правил вывода,

и $u \rightarrow v$, где $u, v \in (V_T \cup V_N)^*$;

S - начальный нетерминальный символ - аксиома грамматики.

Вывод цепочек языка из S по правилам грамматики

Типы $G = (V_T, V_N, P, S)$

Тип грамматики определяется сложностью правил $P = \{u \rightarrow v \mid u, v \in (V_T \cup V_N)^*\}$

Тип 0 – нет ограничений

Тип 1 - НС-грамматики

$\alpha A \beta \rightarrow \alpha \phi \beta, A \in V_N, \alpha, \phi, \beta \in (V_T \cup V_N)^*$

Тип 2 – КС-грамматики

$A \rightarrow \phi, A \in V_N, \phi \in (V_T \cup V_N)^*$

Тип 3 – Автоматные грамматики

- праволинейные

$A \rightarrow xB$ или $A \rightarrow x, A, B \in V_N, x \in V_T$

- леволинейные

Пример грамматики

$G = (V_T, V_N, P, S)$, где

$V_T = \{0, 1, \dots, 9, x, a, \dots, f\}$

$V_N = \{S, H, C\}$

$P = \{S \rightarrow 0xH, H \rightarrow HC, H \rightarrow C, C \rightarrow 0, C \rightarrow 1, \dots, C \rightarrow 9, C \rightarrow a, \dots, C \rightarrow f\}$

Стандартная запись грамматики

$G: S \rightarrow 0xH$

$H \rightarrow HC \mid C$

$C \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \mid a \mid b \mid c \mid d \mid e \mid f$

Пример вывода

$S \rightarrow 0xH \rightarrow 0xHC \rightarrow 0xHCC \rightarrow 0xCCCC \rightarrow 0xaCC \rightarrow 0xa8C$
 $\rightarrow 0xa8f$

Обозначение нетерминалов

Вариант 1 - большие латинские буквы

(при условии их отсутствия в V_T)

G: S \square 0xH

H \square HC | C

C \square 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | a | b | c | d | e | f

Вариант 2 – текст в угловых скобках

G: <число 16> \square 0x <цифры 16>

<цифры 16> \square <цифры 16> <одна цифра 16>

| <одна цифра 16>

<одна цифра 16> \square 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

| a | b | c | d | e | f

Правила КС-грамматики

$$G = (V_T, V_N, P, S)$$

конструкция правила

$\langle \text{нетерминал} \rangle \rightarrow \langle \text{цепочка из } V_T \text{ и } V_N \rangle$

Пример

$$S \rightarrow SaB \mid c$$

$$B \rightarrow xBy \mid \varepsilon$$

ε - пустая цепочка, $|\varepsilon| = 0$

Пример вывода

$$S \rightarrow SaB \rightarrow SaBaB \rightarrow caBaB \rightarrow caaB \rightarrow caaxBy \rightarrow$$

$$caaxxByu \rightarrow caaxxxByuu \rightarrow caaxxxuuu$$

Вывели $ca^2x^3y^3$

КС-грамматика и дерево

Вывода

Правила вывода в КС-грамматике $A \rightarrow \phi$

$A \in V_N, \phi \in (V_T \cup V_N)^*$

$G: S \rightarrow aSAb \mid c$
 $A \rightarrow cA \mid d$

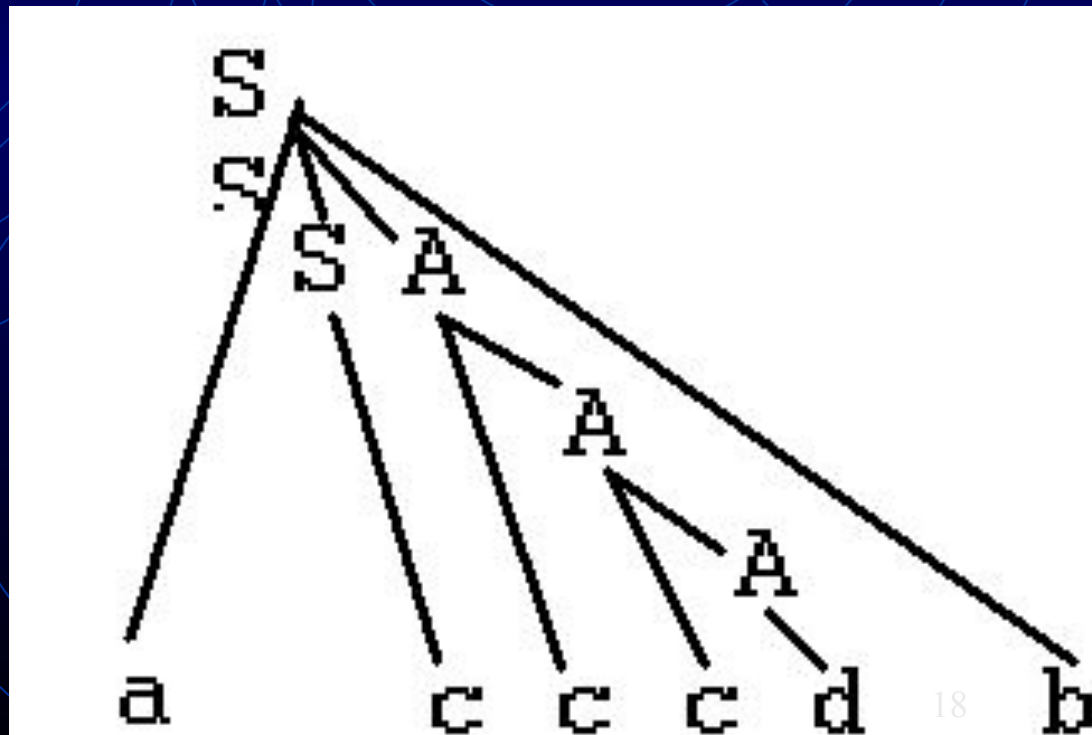
Вывод

$S \rightarrow aSAb \rightarrow acAb$

$\rightarrow accAb$

$\rightarrow acccAb \rightarrow acccdb$

Дерево вывода



Операции над языками

1) Обычные теоретико-множественные:

Объединение $L_1 \cup L_2$

Пересечение $L_1 \cap L_2$

Разность $L_1 \setminus L_2$

2) Специальные:

Произведение или конкатенация $L_1 L_2$

Возведение в степень L^n ,

Итерация L^* ,

Усеченная итерация L^+

Операции – инструмент синтеза КСГ

Теорема

Семейство КС-языков замкнуто относительно операций объединения, произведения, итерации и усеченной итерации.

Метод синтеза КСГ

- 1) Есть КСГ G_1 и G_2 простых языков L_1 и L_2
- 2) Строим нужную нам КСГ как вариант из $L_1 \cup L_2$, $L_1 L_2$, L_2^* , L_2^+

Доказательство для $L_1 L_2$

Конкатенация (умножение)

Дано:

$$G_1 = (V_{T1}, V_{N1}, P_1, S_1), \quad G_2 = (V_{T2}, V_{N2}, P_2, S_2)$$

Нужна грамматика языка $L(G) = L(G_1) L(G_2)$

Решение:

$$V_{N1} \cap V_{N2} = \emptyset$$

Новая грамматика с новым нетерминалом S

$$P = P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\}$$

Вывод в G

$$S \rightarrow S_1 S_2 \rightarrow \alpha S_2 \rightarrow \alpha \beta$$

α выводится в G_1 , β выводится в G_2

Доказательство для $L_1 \cup L_2$

Объединение

Дано:

$$G_1 = (V_{T1}, V_{N1}, P_1, S_1), \quad G_2 = (V_{T2}, V_{N2}, P_2, S_2)$$

Нужна грамматика языка $L(G) = L(G_1) \cup L(G_2)$

Решение:

$$V_{N1} \cap V_{N2} = \emptyset$$

Новая грамматика с новым нетерминалом S

$$P = P_1 \cup P_2 \cup \{S \rightarrow S_1 \mid S_2\}$$

Выводы в G (два варианта, так как два правила)

$$S \rightarrow S_1 \rightarrow \alpha \quad \text{или} \quad S \rightarrow S_2 \rightarrow \beta$$

α выводится в G_1 , β выводится в G_2

Доказательство для L^*

Итерация

Дано: $G_1 = (V_{T1}, V_{N1}, P_1, S_1)$

Нужна грамматика языка $L(G) = L(G_1)^*$

Решение:

Новая грамматика с новым нетерминалом S

$$P = P_1 \cup \{S \rightarrow SS_1, S \rightarrow \varepsilon\}$$

Выводы в G

$$S \rightarrow \varepsilon$$

$$S \rightarrow SS_1 \rightarrow S_1 \rightarrow \beta$$

$$S \rightarrow SS_1 \rightarrow SS_1S_1 \rightarrow S_1S_1 \rightarrow \dots$$

$$S \rightarrow SS_1 \rightarrow SS_1S_1 \rightarrow SS_1S_1S_1 \rightarrow S_1S_1S_1 \rightarrow \dots$$

Правая рекурсия для L^*

Итерация

Дано: $G_1 = (V_{T1}, V_{N1}, P_1, S_1)$

Нужна грамматика языка $L(G) = L(G_1)^*$

Решение:

Новая грамматика с новым нетерминалом S

$$P = P_1 \cup \{S \rightarrow S_1 S, S \rightarrow \varepsilon\}$$

Выводы в G

$$S \rightarrow \varepsilon$$

$$S \rightarrow S_1 S \rightarrow S_1 \rightarrow \beta$$

$$S \rightarrow S_1 S \rightarrow S_1 S_1 S \rightarrow S_1 S_1 \rightarrow \dots$$

$$S \rightarrow S_1 S \rightarrow S_1 S_1 S \rightarrow S_1 S_1 S_1 S \rightarrow S_1 S_1 S_1 \rightarrow \dots$$

Доказательство для L^+

Усеченная итерация

Дано: $G_1 = (V_{T1}, V_{N1}, P_1, S_1)$

Нужна грамматика языка $L(G) = L(G_1)^+$

Решение:

Новая грамматика с новым нетерминалом S

$$P = P_1 \cup \{S \rightarrow S S_1, S \rightarrow S_1\}$$

Выводы в G

$$S \rightarrow S_1 \rightarrow \beta$$

$$S \rightarrow S S_1 \rightarrow S_1 S_1 \rightarrow \dots$$

$$S \rightarrow S S_1 \rightarrow S S_1 S_1 \rightarrow S_1 S_1 S_1 \rightarrow \dots$$

...

Правая рекурсия для L^+

Усеченная итерация

Дано: $G_1 = (V_{T1}, V_{N1}, P_1, S_1)$

Нужна грамматика языка $L(G) = L(G_1)^+$

Решение:

Новая грамматика с новым нетерминалом S

$$P = P_1 \cup \{S \rightarrow S_1 S, S \rightarrow S_1\}$$

Выводы в G

$$S \rightarrow S_1 \rightarrow \beta$$

$$S \rightarrow S_1 S \rightarrow S_1 S_1 \rightarrow \dots$$

$$S \rightarrow S_1 S \rightarrow S_1 S_1 S \rightarrow S_1 S_1 S_1 \rightarrow \dots$$

...

Пример: язык $(ab)^*c \cup b^+$

G: $S \square A | B$ // объединение $(ab)^*c$ и b^+
A $\square Tc$ // произведение $(ab)^*$ и c
F $\square c$
T $\square Tab | \epsilon$ // итерация ab
B $\square Bb | b$ // усеченная итерация b

Можно проще:

G: $S \square A | B$
A $\square Tc$
T $\square Tab | \epsilon$
B $\square Bb | b$

Или еще проще*

G: $S \square Tc | B$
T $\square Tab | \epsilon$
B $\square Bb | b$

Пример: идентификатор С++

Идентификатор – это последовательность букв и цифр, начинающаяся с буквы. Знак подчеркивания рассматривается как буква.

G: <идент> □ <буква> <окончание>
<буква> □ a | b | ... | z | A | B | ... | Z | _
<окончание> □ <окончание> <буква>
 | <окончание> <цифра>
<цифра> □ 0 | 1 | 2 | ... | 9

Симметрия в языке

$$\underline{L = \alpha^n \gamma \beta^n, \quad n \geq 0}$$

$$S \square \alpha S \beta \quad | \quad \gamma$$

Вывод

$$S \square \alpha S \beta \square \alpha^2 S \beta^2 \square \alpha^3 S \beta^3 \square \dots \square \alpha^n S \beta^n \square \alpha^n \gamma \beta^n$$

Не путать с $\underline{L = \alpha^* \gamma \beta^*}$

$$S \square A \gamma B$$

$$A \square A \alpha \quad | \quad \varepsilon$$

$$B \square B \beta \quad | \quad \varepsilon$$

Простой пример симметрии

Симметричные конструкции

$S \sqsubseteq ASB \mid X$

$$L(S) = L(A)^n L(X) L(B)^n \quad n \geq 0$$

Пример

$(ab)^n c (xyz)^n \quad n \geq 0$

$S \sqsubseteq abSxyz \mid c$

Пример симметрии

$$\underline{L = a^{2n} d (a^+ b^* c)^{n+1}}, \quad \underline{n \geq 0}$$

1. Выделим симметрию

$$\begin{aligned} L &= a^{2n} d (a^+ b^* c)^{n+1} \\ &= \underline{(aa)^n} \quad d \quad a^+ b^* c \quad \underline{(a^+ b^* c)^n} \end{aligned}$$

2. Строим КС-грамматику

G: S \square aaSB | dB

B \square Tnc

T \square Ta | a

H \square Hb | ε

Преобразования КСТ

1. Подстановки
2. Новые нетерминалы для фрагментов
3. Неукорачивание
4. Левая рекурсия правая
5. Правая рекурсия левая
6. Удаление одинаковых правых частей правил
7. Удаление правил «нетерминал нетерминал»
8. Удаление лишних нетерминалов
9. ...

Левая рекурсия правая

Были правила

$S \rightarrow SA \mid B$

Вывод:

$S \rightarrow SA \rightarrow SAA \rightarrow SA^* \rightarrow BA^*$

Новые правила

$S \rightarrow BH$

$H \rightarrow AH \mid \varepsilon$

Новый вывод

$S \rightarrow BH \rightarrow BAN \rightarrow BAAN \rightarrow \dots \rightarrow BA^*$

Левая рекурсия \square правая (Общий случай)

Были правила

$$S \square S\beta_1 \mid S\beta_2 \mid \dots \mid S\beta_k \mid \phi_1 \mid \phi_2 \mid \dots \mid \phi_m$$

Преобразуем :

$$S \rightarrow SX \mid Y \quad (\text{будет } S \square YN \quad N \square XN \mid \varepsilon)$$

$$X \square \beta_1 \mid \beta_2 \mid \dots \mid \beta_k$$

$$Y \square \phi_1 \mid \phi_2 \mid \dots \mid \phi_m$$

Новые правила строим так же, как и ранее

$$S \square \phi_1 N \mid \phi_2 N \mid \dots \mid \phi_m N$$

$$N \square \beta_1 N \mid \beta_2 N \mid \dots \mid \beta_k N \mid \varepsilon$$

Правая рекурсия \square левая (Общий случай)

Были правила

$$S \square \beta_1 S \mid \beta_2 S \mid \dots \mid \beta_k S \mid \phi_1 \mid \phi_2 \mid \dots \mid \phi_m$$

Преобразуем :

$$S \rightarrow XS \mid Y \quad (\text{будет } S \square NY \quad N \square NX \mid \varepsilon)$$

$$X \square \beta_1 \mid \beta_2 \mid \dots \mid \beta_k$$

$$Y \square \phi_1 \mid \phi_2 \mid \dots \mid \phi_m$$

Новые правила строим так же, как и ранее

$$S \square N \phi_1 N \mid N \phi_2 \mid \dots \mid N \phi_m$$

$$N \square N \beta_1 \mid N \beta_2 \mid \dots \mid N \beta_k \mid \varepsilon$$

Рекурсия

Грамматика конечна по определению.
Язык, как правило, бесконечен.

Рекурсия – единственный способ представить бесконечный язык конечными средствами.

Типы рекурсии:

- Левая $A \square Ab,$
- Правая $A \rightarrow bA,$
- Центральная $A \square aAb,$
- Неявная $A \square xBy \quad B \square uDv \quad D \square aAb .$

Рекурсия конструкций ЯП

<оператор> □ <простой оператор>

| <составной оператор>

| <пустой оператор>

<пустой оператор> □ ;

<составной оператор> □ {<опер и описания>}

<простой оператор> □ <if> | <while> | <for> | ...

<if> □ if (<выражение>) <оператор>

| if (<выражение>) <оператор> else <оператор>

<while> □ while (<выражение>) <оператор>

<опер и описания> □ <опер и описания> < оператор >

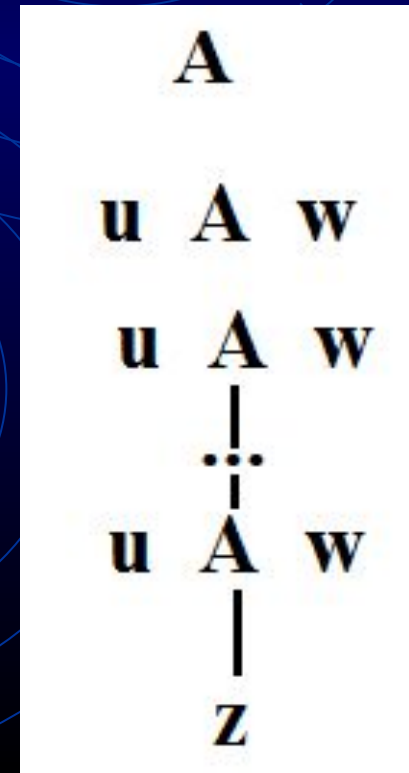
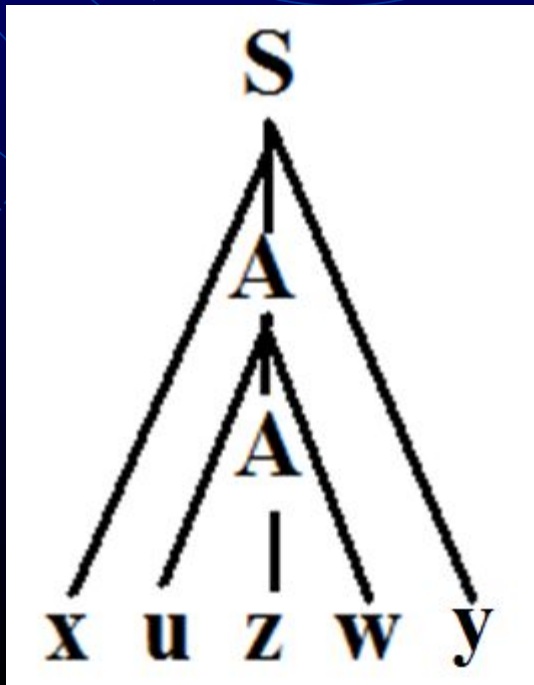
| <опер и описания> < описание >

| ε

Синтаксический анализ

Задача любого транслятора – построить дерево грамматического разбора.

Если язык бесконечен, а КС-грамматика конечна, то в дереве СА достаточно длинной цепочки обязательно на некотором пути из корня в лист повторится некоторый нетерминал $A \square u^n Aw^n \square u^n zw^n$



Лемма о разрастании

Для любой КС - грамматики, порождающей бесконечный язык, существуют такие натуральные числа p и q , что каждая цепочка $\gamma \in L(G)$, $|\gamma| > p$ может быть представлена в виде $\gamma = xuzw^ny$, где $|uw| > q$ и для любого $n > 0$ цепочка $xu^nzw^ny \in L(G)$.

Мы получили на предыдущем слайде нетерминал $A \Rightarrow u^nAw^n \Rightarrow u^nzw^n$

A diagram illustrating the derivation of the string u^nzw^n from the non-terminal A . The derivation is shown as a vertical sequence of strings connected by arrows. The top string is A . Below it is uAw . Below that is uAw again. A vertical ellipsis indicates that this process repeats $n-2$ times. Below the ellipsis is uAw . A final arrow points down to the string uz .

Следствие – теорема о языке $a^n b^n c^n$

Теорема

Язык $a^n b^n c^n$ не является КС-языком.

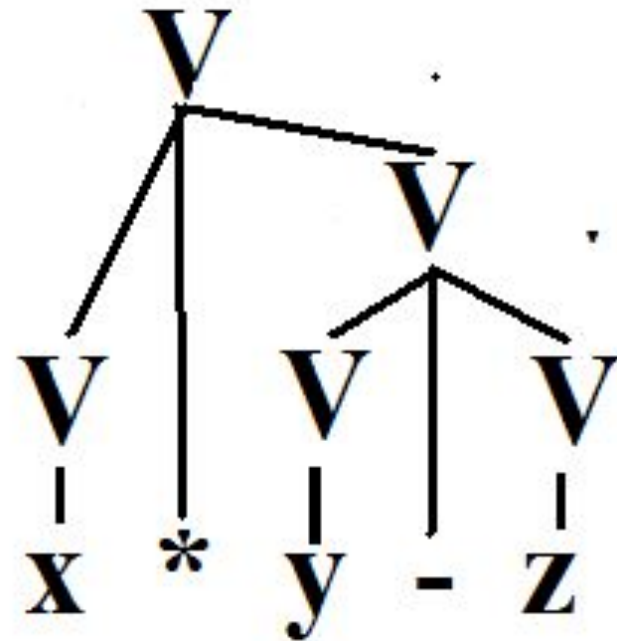
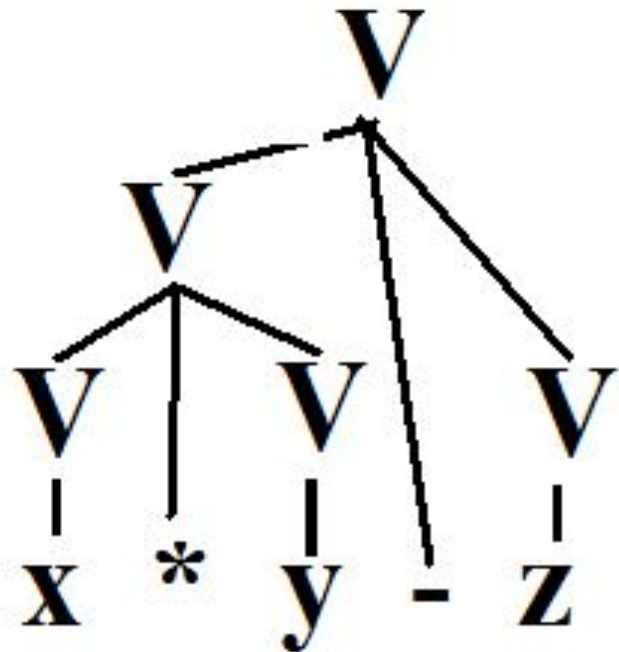
Вывод: КС-грамматика не может синхронизировать более двух фрагментов.

Пример

```
int aaaa;  
    // идентификаторы состоят только из букв а  
int main() {  
    cin >> aaaa;  
    cout << aaaa * 2;  
}
```


Недетерминированность КСГ недопустима !!!

Пусть $V == \langle \text{выражение} \rangle$ $a == \langle \text{идент} \rangle$
 $G: V \square V+V \mid V-V \mid V*V \mid V/V \mid (V) \mid a$



Выражения

Нет понятия логического, арифметического и т.п. выражения:

Работает семантика контекстных условий для контроля корректности выражений:

```
int x, z[3][100], T[100];
```

```
bool y;
```

```
x[i] = z * y(6.7) - 5; // семантические ошибки
```

```
T[i] = x * sin(6.7) - 5; // верная конструкция
```

```
double a = x - 7.78; // приведение типов
```

Синтаксис выражений

- 1 Упорядочить группы операций по приоритетам, начиная с низкоприоритетных .
- 2 Каждой группе поставить в соответствие нетерминал $A_1, A_2, A_3, \dots, A_n$
- 3 Для каждого нетерминала A_i записать правила для каждой операции группы (см. следующий слайд)
- 4 Добавить правила, соответствующие неприменению операции
 $A_i \square A_{i+1}$
- 5 Определить правила для элементарного выражения A_{n+1}

Правила для группы многократных операций

Операции многократные, т.е. разрешается запись $a*a*a$

Операция бинарная

выполнение слева направо $A_i \square A_i * A_{i+1}$

справа налево $A_i \square A_{i+1} * A_i$

Операция префиксная унарная $A_i \square * A_i$

постфиксная унарная $A_i \square A_i *$

Правила для группы однократных операций

Операции однократные, т.е. НЕ разрешается запись $a*a*a*a$

Операция бинарная

$$A_i \square A_{i+1} * A_{i+1}$$

Операция префиксная унарная $A_i \square * A_{i+1}$

постфиксная унарная $A_i \square A_{i+1} *$

Пример выражения с бинарными операциями

A1 > < ==

A2 + -

A3 * /

A1 □ A1 > A2 | A1 < A2 | A1 == A2 | A2

A2 □ A2 + A3 | A2 - A3 | A3

A3 □ A3 * A4 | A3 / A4 | A4

A4 □ <идентификатор> | <константа> | (A1)

Лабораторная работа №1

1. Построить КС-грамматику по заданию в учебнике (№ задания == номер в списке группы)
2. Построить деревья разбора для двух разных цепочек

Учебное пособие (содержит теоретический материал и пример выполнения задания):

- глава 1 - материал по грамматикам

Лабораторная работа № 2

1. Построить таблицу приоритетов операций языка программирования индивидуального задания
2. Построить КС-грамматику языка программирования индивидуального задания

Учебное пособие (содержит теоретический материал и пример выполнения задания):

- глава 3

Языки, грамматика, автоматы

1. Грамматика порождает язык
2. Автомат распознает язык
3. Программист использует (1) при создании программы
4. Транслятор использует (2) в процессе работы
5. Вывод: (1) и (2) должны быть эквивалентны

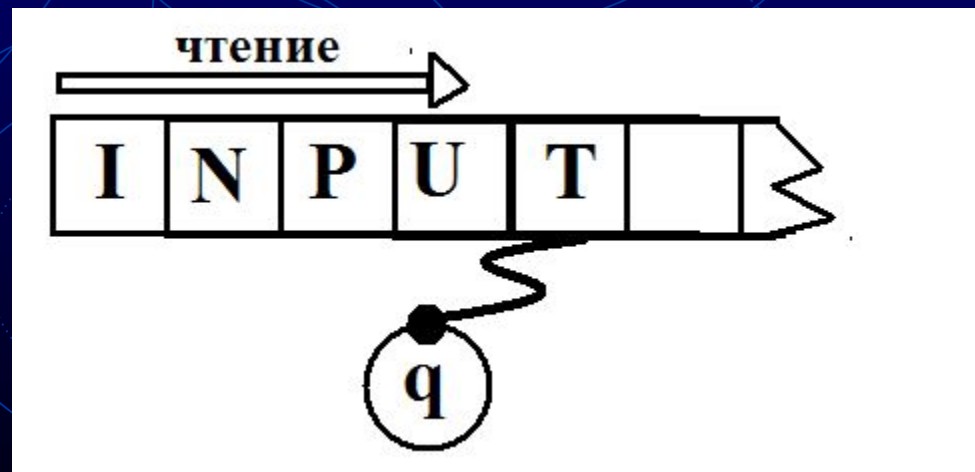
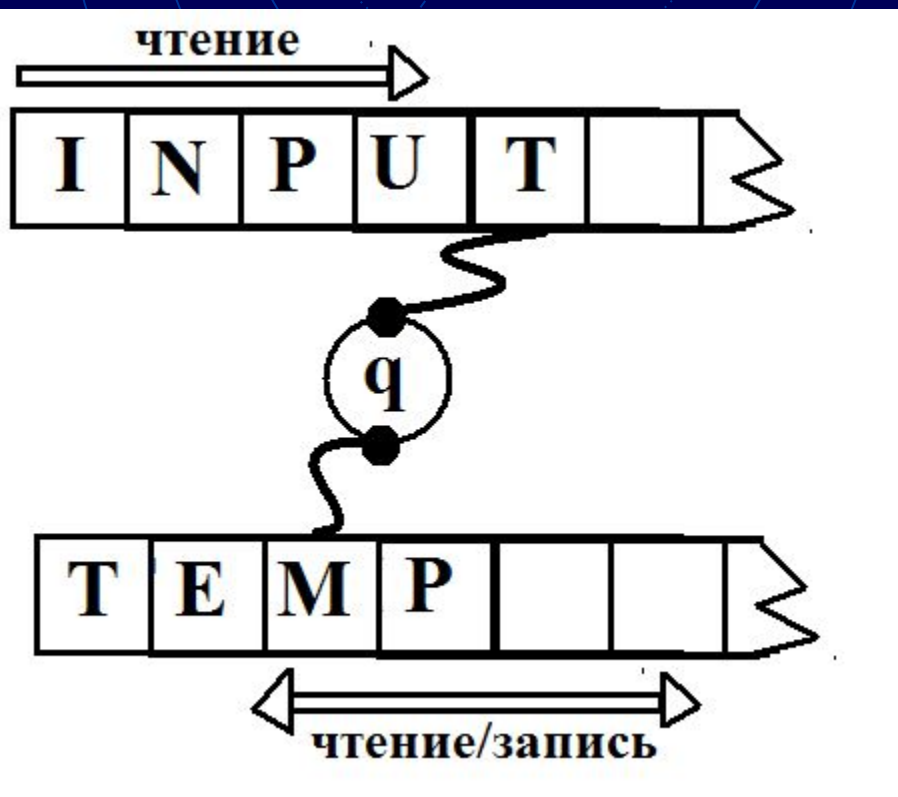
Разработчик транслятора создает грамматику, выбирает метод распознавания, реализует метод распознавания в соответствии с грамматикой

Типы автоматов при трансляции

Машина Тьюринга == любой алгоритм.

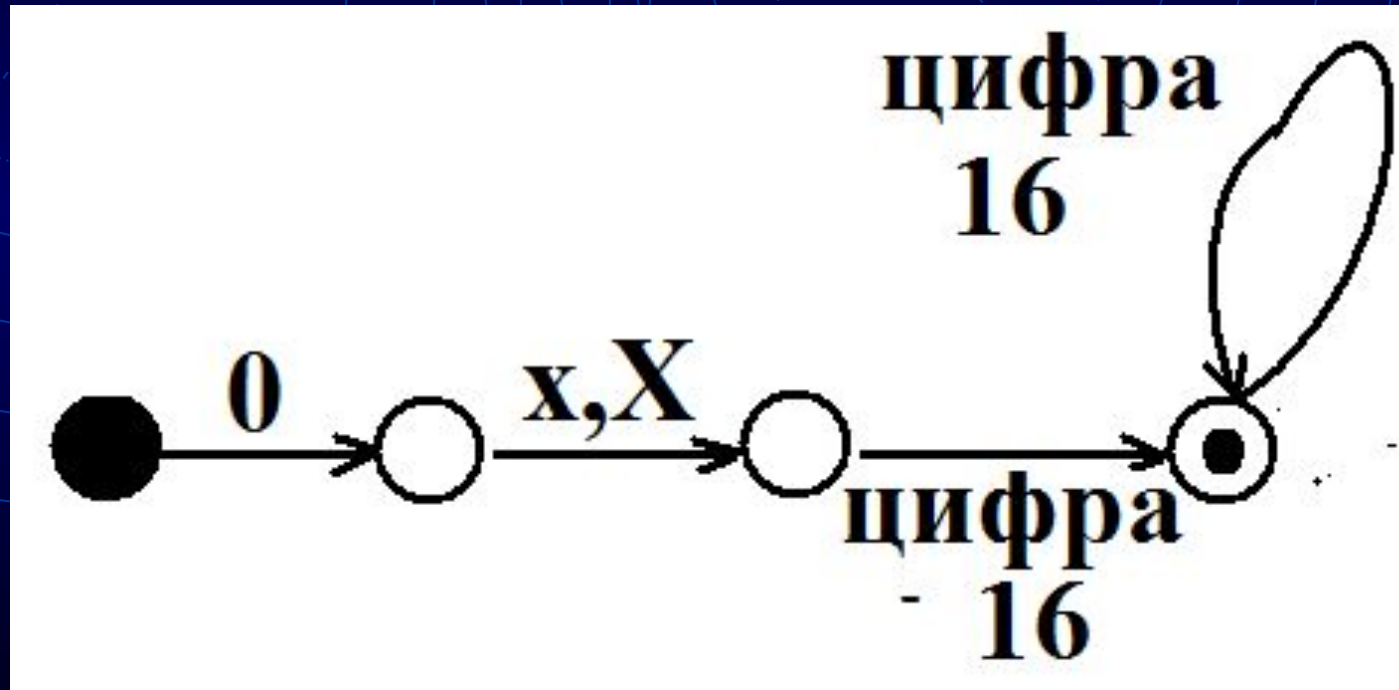
Требование трансляции - эффективность

$$T(v) = O(n)$$



Автоматы при трансляции

Конечный автомат лексика языка
(переходит из состояния в состояние, читая символы)

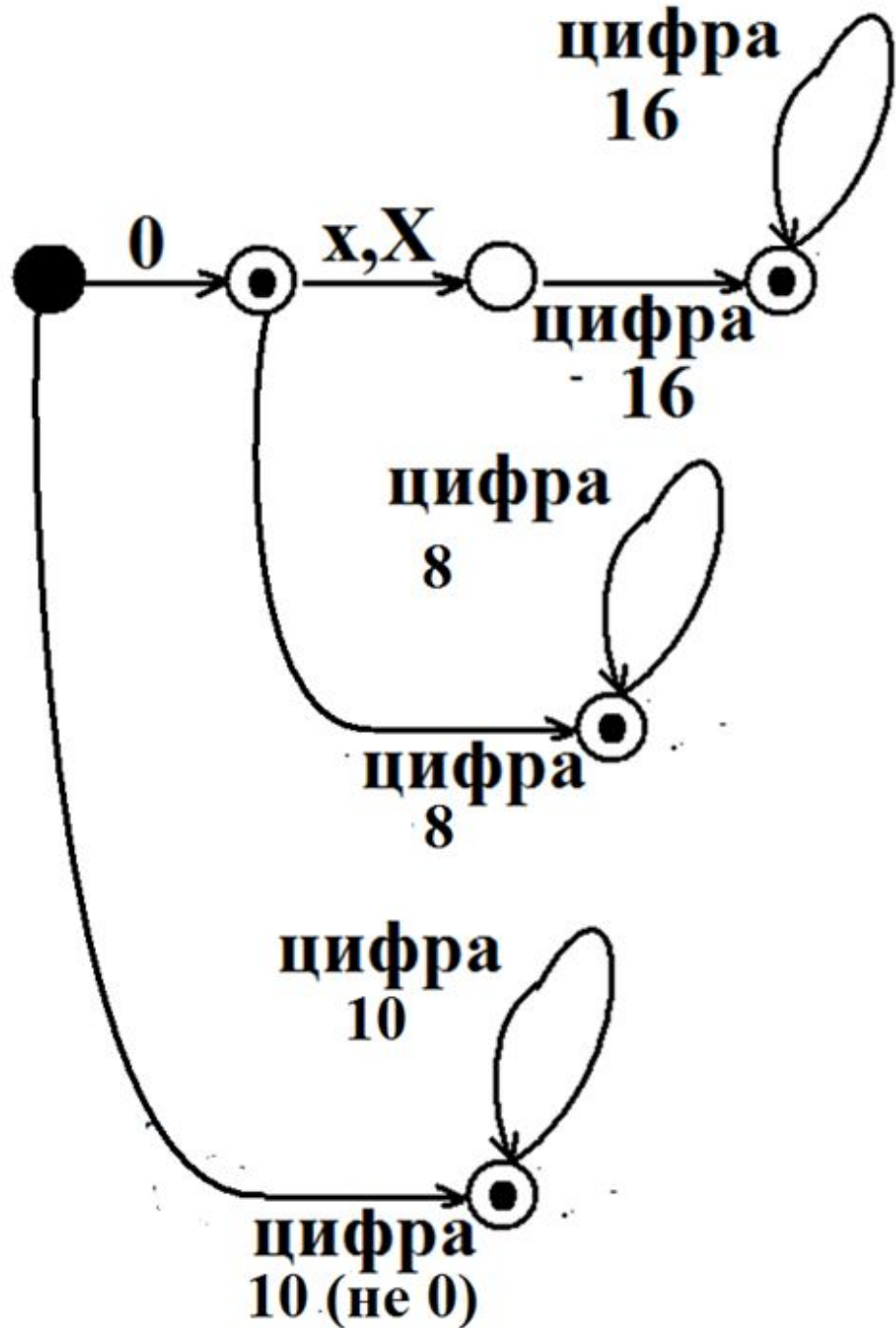


МП-автомат (автомат с магазинной памятью
синтаксис языка)

Синтез конечных автоматов

Учебное пособие (содержит теоретический материал и примеры):

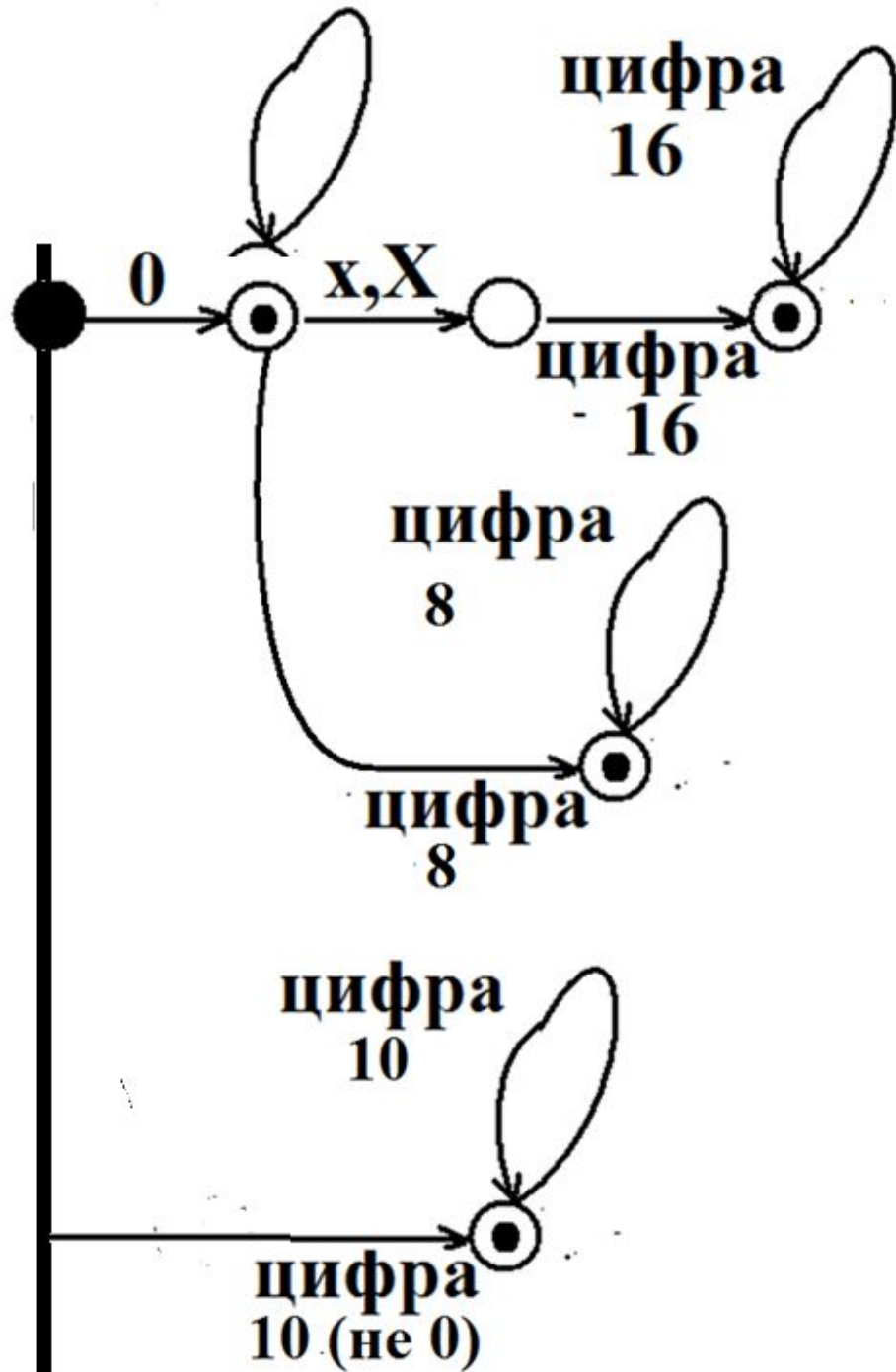
- глава 2



Целые константы:

- 1) 8с/с – начинается с 0,
- 2) 16с/с – начинаются 0x или 0X,
- 3) 10 с.с – начинаются с цифры, не равной нулю

Целые константы





The end

Лабораторные работы далее – лексический анализ