

Лекции по теории формальных языков

Лекция 4.

Контекстно-свободные грамматики и языки:
определения и примеры.

Лемма о накачке

Александр Сергеевич Герасимов

<http://gas-teach.narod.ru>

Кафедра математических и информационных технологий
Санкт-Петербургского академического университета
Российской академии наук.
Весенний семестр 2010/11 учебного года

4 марта 2011 г.

План

- 1 Контекстно-свободные грамматики и языки: определения и примеры
- 2 Лемма о накачке

План

1 Контекстно-свободные грамматики и языки: определения и примеры

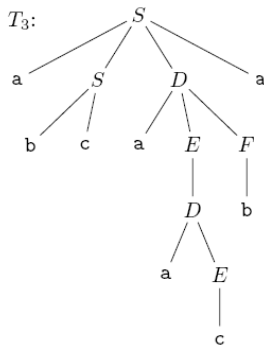
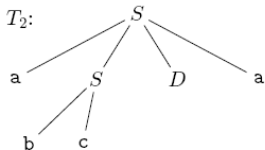
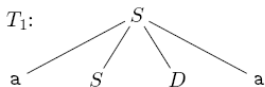
2 Лемма о накачке

Определения КС-грамматики и КС-языка. Примеры

- Грамматика называется *контекстно-свободной* (КС-грамматикой), если каждое её правило имеет вид $A \rightarrow \alpha$.
- Язык называется *контекстно-свободным* (КС-языком), если некоторая КС-грамматика его порождает.
- Если $G = (\Sigma, \Gamma, P, S)$ — грамматика, $S \Rightarrow_G^* \beta \Rightarrow_G^* w$, то цепочку β будем называть *формой*.
- $\{S \rightarrow aSb|ab\}$ — КС-грамматика, она порождает КС-язык $\{a^n b^n | n \geq 1\}$.
- Язык $\{a^n b^n a^m | m, n \geq 1\}$ контекстно-свободен, поскольку порождается КС-грамматикой $\{S \rightarrow CD, C \rightarrow aCb|ab, D \rightarrow aD|a\}$.
- Грамматика $\{S \rightarrow aSDa, S \rightarrow aba, aD \rightarrow Da, bD \rightarrow bb\}$ не являющаяся контекстно-свободной, порождает язык $\{a^n b^n a^n | n \geq 1\}$ (см. лекцию 3). Существует ли КС-грамматика, порождающая этот язык? Ответ — в текущей лекции.

Дерево вывода: пример

- Вывод в КС-грамматике удобно представлять в виде дерева.
- Для примера рассмотрим КС-грамматику G_1 с правилами $\{S \rightarrow aSDa|bc, D \rightarrow aE|aEF, E \rightarrow D|c, F \rightarrow b\}$ и вывод

$$\underline{S} \Rightarrow a\underline{S}Da \Rightarrow abc\underline{D}a \Rightarrow abca\underline{E}Fa \Rightarrow abca\underline{D}Fa \Rightarrow abcaac\underline{E}Fa \Rightarrow abcaac\underline{F}a \Rightarrow abcaacba = w .$$


Дерево вывода: определение

Деревом вывода (или *деревом разбора*) цепочки $w \in \Sigma^*$ в грамматике $G = (\Sigma, \Gamma, P, S)$ называется такое упорядоченное дерево, что

- его корень помечен S ;
- если его внутренний узел помечен $A \in \Gamma$ и $X_1, \dots, X_k \in \Sigma \cup \Gamma$ — перечисленные слева направо пометки всех сыновей этого внутреннего узла, то $A \rightarrow X_1 \dots X_k \in P$;
- если его внутренний узел помечен $A \in \Gamma$ и ε — пометка единственного сына этого внутреннего узла, то $A \rightarrow \varepsilon \in P$;
- $w = a_1 \dots a_m$, где $a_1, \dots, a_m \in \Sigma \cup \{\varepsilon\}$ — перечисленные слева направо пометки всех листьев этого дерева.

На слайде 5 дерево T_3 — дерево вывода цепочки w .

Стандартное поддереве дерева вывода.

Представление вывода деревом вывода

- Поддереве T' дерева вывода T называется *стандартным*, если корень T является корнем T' , и для любого узла x дерева T' либо x является листом T' , либо все сыновья узла x в дереве T также лежат в T' .
- Например, на слайде 5 деревья T_1 и T_2 — стандартные поддеревья дерева вывода T_3 .
- Будем говорить, что вывод $S \Rightarrow \alpha_1 \Rightarrow \dots \Rightarrow \alpha_n = w$ в грамматике G *представлен* деревом вывода T , если существуют такие стандартные поддеревья $T_1 \subset \dots \subset T_n$ дерева T , что для каждого $i = 1, \dots, n$ конкатенация пометок перечисленных слева направо листьев поддерева T_i есть форма α_i .
Также мы будем говорить, что вывод формы α_i *представлен* деревом T_i .
- Например, вывод на слайде 5 представлен деревом вывода T_3 .

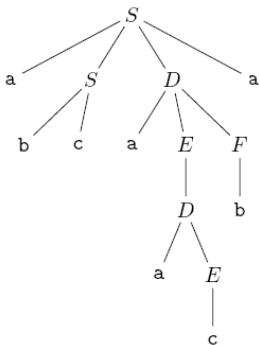
Разные выводы одной цепочки I

Продолжаем рассматривать КС-грамматику G_1 с правилами $\{S \rightarrow aSDa|bc, D \rightarrow aE|aEF, E \rightarrow D|c, F \rightarrow b\}$ со слайда 5.

Вывод

$$(1) \quad \underline{S} \Rightarrow a\underline{S}Da \Rightarrow abc\underline{D}a \Rightarrow abca\underline{E}Fa \Rightarrow abca\underline{D}Fa \Rightarrow \\ abca\underline{a}E\underline{F}a \Rightarrow abcaac\underline{F}a \Rightarrow abcaacba = w$$

представлен следующим деревом вывода:



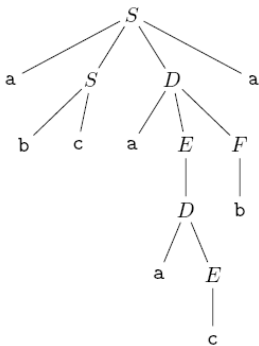
Разные выводы одной цепочки II

Продолжаем рассматривать КС-грамматику G_1 с правилами $\{S \rightarrow aSDa|bc, D \rightarrow aE|aEF, E \rightarrow D|c, F \rightarrow b\}$ со слайда 5.

Вывод

$$(2) \quad \underline{S} \Rightarrow a\underline{S}D\underline{a} \Rightarrow aS\underline{a}E\underline{F}a \Rightarrow aS\underline{a}E\underline{b}a \Rightarrow aS\underline{a}D\underline{b}a \Rightarrow \\ aS\underline{a}aE\underline{b}a \Rightarrow a\underline{S}aacsba \Rightarrow abcsaacba = w$$

представлен тем же деревом вывода, что и на слайде 8:



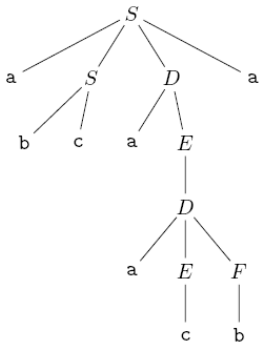
Разные выводы одной цепочки III

Рассматриваем КС-грамматику G_1 с правилами
 $\{S \rightarrow aSDa|bc, D \rightarrow aE|aEF, E \rightarrow D|c, F \rightarrow b\}$ со слайда 5.

Вывод

$$(3) \quad \underline{S} \Rightarrow a\underline{S}Da \Rightarrow abc\underline{D}a \Rightarrow abca\underline{E}a \Rightarrow abca\underline{D}a \Rightarrow \\ abcaa\underline{E}Fa \Rightarrow abcaac\underline{F}a \Rightarrow abcaacba = w$$

представлен следующим деревом вывода:

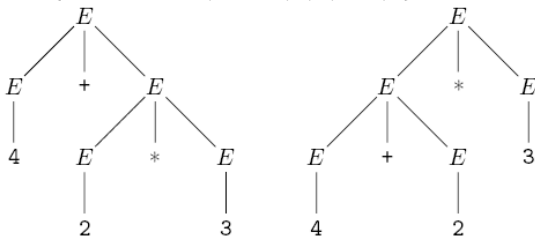


Левый и правый выводы

- Вывод в КС-грамматике называется *левым*, или *левосторонним*, (соответственно, *правым*, или *правосторонним*), если правило применяется к самому левому (соответственно, правому) вхождению нетерминала в каждую форму этого вывода.
- Произвольную форму из левого (соответственно, правого) вывода будем называть *l-формой* (соответственно, *r-формой*).
- Например, вывод (1) на слайде 8 и вывод (3) на слайде 10 являются левыми, а вывод (2) на слайде 9 — правым.
- Любое дерево вывода цепочки представляет единственный левый и единственный правый вывод этой цепочки.
- Цепочка имеет два или более различных дерева вывода тогда и только тогда, когда эта цепочка имеет два или более различных левых (соответственно, правых) вывода.

Однозначные и неоднозначные грамматики

- КС-грамматика G называется *однозначной*, если каждая цепочка из $L(G)$ имеет единственное дерево вывода. В противном случае КС-грамматика G называется *неоднозначной*.
- КС-грамматика G однозначна тогда и только тогда, когда каждая цепочка из $L(G)$ имеет единственный левый вывод.
- Рассматривавшаяся выше КС-грамматика G_1 неоднозначна.
- КС-грамматика $\{E \rightarrow E + E | E * E | 0 | 1 | \dots | 9\}$ неоднозначна:



- Желательно, чтобы грамматика, определяющая язык программирования, была однозначной. Иначе смысл некоторых программ можно понимать по-разному.

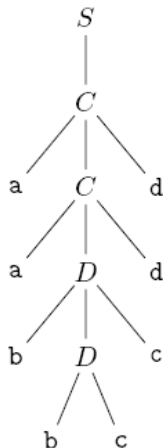
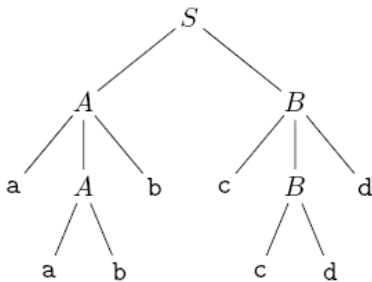
Неразрешимость проблемы однозначности.

Однозначные и неоднозначные языки

- Доказано, что не существует алгоритма, определяющего по произвольной КС-грамматике, однозначна она или нет.
- КС-язык называется *однозначным*, если некоторая однозначная КС-грамматика порождает его. В противном случае КС-язык называется *неоднозначным*.
- Установлено, что КС-язык $L = \{a^m b^m c^n d^n \mid m, n \geq 1\} \cup \{a^m b^n c^n d^m \mid m, n \geq 1\}$ неоднозначен.
- КС-язык L порождается, например, КС-грамматикой с правилами $S \rightarrow AB|C$, $A \rightarrow aAb|ab$, $B \rightarrow cBd|cd$, $C \rightarrow aCd|aDd$, $D \rightarrow bDc|bc$.

Два дерева вывода одной цепочки

Цепочка $a^2b^2c^2d^2 \in L$ (см. предыдущий слайд) имеет два дерева вывода в грамматике с правилами $S \rightarrow AB|C$, $A \rightarrow aAb|ab$, $B \rightarrow cBd|cd$, $C \rightarrow aCd|aDd$, $D \rightarrow bDc|bc$:



КС-грамматики, порождающие скобочный язык LB

- Алфавит языка LB : $\{[,]\}$.
Определение языка LB :
 - ▶ $\varepsilon \in LB$;
 - ▶ если $u \in LB$ и $v \in LB$, то $[u] \in LB$ и $uv \in LB$.
- Грамматика $GB_1 = \{S \rightarrow \varepsilon \mid [S] \mid SS\}$ неоднозначна, поскольку цепочка ε имеет два различных левых вывода:
 - ▶ $S \Rightarrow \varepsilon$ и
 - ▶ $S \Rightarrow SS \Rightarrow \varepsilon S \Rightarrow \varepsilon\varepsilon = \varepsilon$.
- Грамматика $GB_2 = \{S \rightarrow \varepsilon \mid S[S]\}$ однозначна, но содержит правило вида $A \rightarrow A\alpha$ (так называемое леворекурсивное правило).
- Грамматика $GB_3 = \{S \rightarrow \varepsilon \mid [S]S\}$ однозначна.

КС-грамматики, порождающие язык списков LL

- Алфавит языка LL : $\{a, ;, [,]\}$.
Определение языка LL :
 - ▶ $a \in LL$;
 - ▶ если $u \in LL$ и $v \in LL$, то $[u] \in LL$ и $u;v \in LL$.
- Грамматика $GL_1 = \{S \rightarrow a \mid [S] \mid S; S\}$ неоднозначна, поскольку цепочка $a; a; a$ имеет два различных левых вывода:
 - ▶ $S \Rightarrow S; S \Rightarrow a; S \Rightarrow a; S; S \Rightarrow a; a; S \Rightarrow a; a; a$ и
 - ▶ $S \Rightarrow S; S \Rightarrow S; S; S \Rightarrow a; S; S \Rightarrow a; a; S \Rightarrow a; a; a$.
- Грамматика $GL_2 = \{S \rightarrow S; L \mid L, L \rightarrow a \mid [S]\}$ однозначна.
- Грамматика $GL_3 = \{S \rightarrow L; S \mid L, L \rightarrow a \mid [S]\}$ однозначна.

КС-грамматика, порождающая язык описаний типов LD

- Алфавит языка LD : $\{i, :, ;, \text{int}, \text{real}\}$.
Каждая цепочка языка LD имеет вид:
 - ▶ непустой список букв i , разделённых точкой с запятой,
 - ▶ затем двоеточие,
 - ▶ наконец, один из символов int или real .
- Грамматика $GD = \{D \rightarrow L : T, L \rightarrow L; i \mid i, T \rightarrow \text{int} \mid \text{real}\}$.

КС-грамматики, порождающие языки двоичных чисел LN и LN'

- Алфавит языка LN : $\{0, 1, .\}$.
Каждая цепочка языка LN имеет вид
 - ▶ u (где $u \in \{0, 1\}^+$),
 - ▶ $.v$ (где $v \in \{0, 1\}^+$) или
 - ▶ $u.v$ (где $u, v \in \{0, 1\}^+$).
- $GN = \{S \rightarrow L \mid .L \mid L.L, L \rightarrow B \mid LB, B \rightarrow 0 \mid 1\}$.
- Язык LN' получаем, если дополнительно потребуем, чтобы целая часть (u) числа начиналась с 1, дробная часть (v) заканчивалась на 1.
- $GN' = \{S \rightarrow L \mid .R \mid L.R, L \rightarrow 1 \mid L0 \mid L1, R \rightarrow 0R \mid 1R \mid 1\}$.

КС-грамматики, порождающие язык арифметических выражений LA

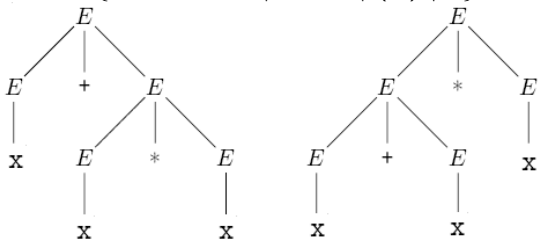
- Алфавит языка LA : $\{x, +, *, (,)\}$.

Определение языка LA :

▶ $x \in LA$;

▶ если $u \in LA$ и $v \in LA$, то $(u) \in LA$, $u + v \in LA$ и $u * v \in LA$.

- Грамматика $GA_1 = \{E \rightarrow E + E \mid E * E \mid (E) \mid x\}$ неоднозначна:



- Грамматика

$GA_2 = \{E \rightarrow T \mid E + T, T \rightarrow F \mid T * F, F \rightarrow (E) \mid x\}$

однозначна, она учитывает естественный приоритет операций и их левую ассоциативность.

План

1 Контекстно-свободные грамматики и языки: определения и примеры

2 Лемма о накачке

Лемма о накачке

Лемма (о накачке)

Для любого КС-языка L существуют числа n и m такие, что каждая цепочка $w \in L$ при $|w| > n$ представима в виде $w = xuzv^k y$, где

- $uv \neq \varepsilon$,
- $|uzv| \leq m$,
- $xu^kzv^ky \in L$ для любого $k \geq 0$.

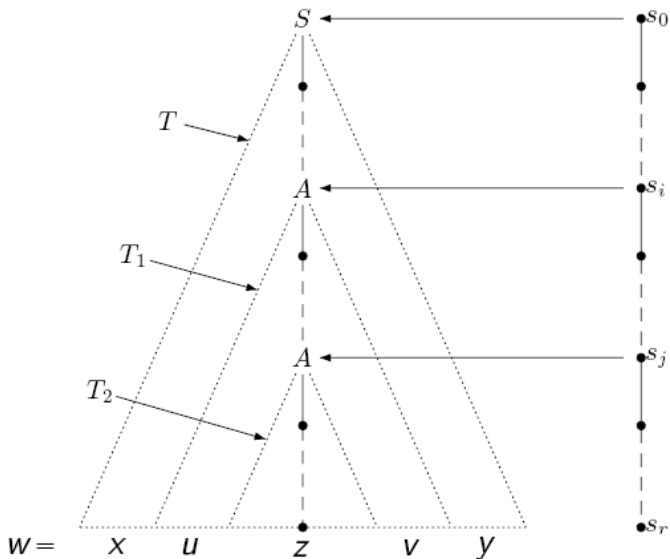
Доказательство.

- Для $L = \emptyset$ лемма верна. Далее считаем, что $L \neq \emptyset$.
- Пусть КС-грамматика $G = (\Sigma, \Gamma, P, S)$ порождает язык L .
- Для каждой цепочки $w \in L$ можно выбрать её дерево вывода в G с минимальным числом узлов; такое дерево вывода мы будем называть минимальным.

Лемма о накачке: продолжение доказательства

- Пусть n — максимальная длина цепочки с минимальным деревом вывода высоты не более $|\Gamma|$. (Хотя бы одна такая цепочка существует, так как $L \neq \emptyset$.)
- Рассмотрим цепочку $w \in L$, $|w| > n$, а также минимальное дерево вывода T цепочки w .
- Высота дерева T больше $|\Gamma|$, поэтому в длиннейшем пути от корня до листа найдутся 2 различных узла, помеченные одним и тем же нетерминалом.
- Пусть этот путь имеет вид $s_0, \dots, s_i, \dots, s_j, \dots, s_r$, а узлы s_i и s_j ($i < j$) помечены нетерминалом A . Можно считать, что $r - i \leq |\Gamma| + 1$.
- Введём обозначения следующих поддеревьев дерева T :
 - ▶ T_1 , состоящее из узла s_i и всех его потомков,
 - ▶ T_2 , состоящее из узла s_j и всех его потомков,
 - ▶ T' , получаемое из T удалением всех потомков узла s_i ,
 - ▶ T'_1 , получаемое из T_1 удалением всех потомков узла s_j .

Лемма о накачке: продолжение доказательства



Лемма о накачке: продолжение доказательства

- Дерево T представляет вывод $S \Rightarrow^* w$.
- Стандартное поддереву T' дерева T представляет вывод $S \Rightarrow^* xAy$.
- Поддереву T_1 с корнем s_i дерева T представляет вывод $A \Rightarrow^* uzv$.
- Стандартное поддереву T'_1 дерева T_1 представляет вывод $A \Rightarrow^* uAv$.
- Поддереву T_2 с корнем s_j дерева T представляет вывод $A \Rightarrow^* z$.
- Докажем, что представление $w = xuzvu$ — искомое.
- $S \Rightarrow^* xAy \Rightarrow^* xuAvy \Rightarrow^* xu^2Av^2y \Rightarrow^* \dots \Rightarrow^* xu^kAv^ky \Rightarrow^* xu^kzv^ky$, так что $xu^kzv^ky \in L$ для любого $k \geq 0$.
- Если бы $uv = \varepsilon$, то дерево T'_1 представляло бы вывод $A \Rightarrow^* A$. Тогда, заменив в дереве вывода T поддереву T_1 на T_2 , мы получили бы дерево вывода цепочки w с меньшим, чем в T , числом узлов; это противоречило бы минимальности T . Следовательно, $uv \neq \varepsilon$.

Лемма о накачке: окончание доказательства

- Положим $m = M^{|\Gamma|+1}$, где M — наибольшая длина правой части правила из P , и покажем, что $|uzv| \leq m$.
- Высота дерева T_1 равна $(r - i)$. Поэтому в T_1 не более M^{r-i} листьев.
- Таким образом, $|uzv| \leq M^{r-i} \leq M^{|\Gamma|+1} = m$.



Язык $L_1 = \{a^n b^n a^n \mid n \geq 1\}$ не является КС-языком

- Предположим, что L_1 является КС-языком.
- Тогда по лемме о накачке для некоторого m верно $a^m b^m a^m = xuzv^m$, $uv \neq \varepsilon$ и $xu^2zv^2y \in L_1$.
- Мы будем говорить, что цепочка $a^n b^n a^n$ состоит из трёх блоков a^n , b^n и a^n .
- Если u или v пересекается с двумя блоками, то в цепочке xu^2zv^2y более двух перемен букв. Поэтому как u , так и v входит лишь в один блок.
- Но тогда в xu^2zv^2y один из блоков короче другого. Противоречие.



Язык $L_2 = \{wsw \mid w \in \{a, b\}^*\}$ не является КС-языком: начало доказательства

- Предположим, что L_2 является КС-языком.
- В силу леммы о накачке цепочка $w_1 = a^m b^m c a^m b^m \in L_2$ при некотором m представима в виде $w_1 = xuzv$, где $uv \neq \varepsilon$, $|uzv| \leq m$ и $w_2 = xu^2zv^2y \in L_2$.
- Буква c не входит ни в u , ни в v , так как иначе эта буква будет входить в w_2 дважды.
- Цепочки u и v располагаются в w_1 по разные стороны от буквы c , поскольку в противном случае с одной стороны от c в цепочке w_2 будет больше букв, чем с другой.
- Теперь очевидно, что буква c входит в z .
- Поскольку $|uzv| \leq m$, то
 - ▶ u входит в подцепочку b^m цепочки w_1 слева от буквы c ,
 - ▶ v входит в подцепочку a^m цепочки w_1 справа от буквы c .

Язык $L_2 = \{wsw \mid w \in \{a, b\}^*\}$ не является
КС-языком: окончание доказательства

- Тогда

- (1) при $u \neq \varepsilon$ в w_2 число вхождений b слева от s больше числа вхождений b справа от s ,
- (2) при $v \neq \varepsilon$ в w_2 число вхождений a справа от s больше числа вхождений a слева от s .

Так как мы имеем $uv \neq \varepsilon$ и в каждом из случаев (1) и (2) заключаем, что $w_2 \notin L_2$, то мы получили противоречие.



Теорема о языках в однобуквенных алфавитах и периодических множествах

Множество $M \subseteq \{0, 1, 2, \dots\}$ называется *периодическим*, если существуют целые положительные числа n_0 (*индекс*) и d (*период*) такие, что для любого $n \geq n_0$ условие $n \in M$ влечёт $n + d \in M$.

Теорема

Для произвольного языка L над алфавитом $\{a\}$ следующие утверждения эквивалентны:

- (1) L — КС-язык;
- (2) L — регулярный язык;
- (3) $M = \{n \geq 0 \mid a^n \in L\}$ — периодическое множество.

Теорема: начало доказательства

(1) влечёт (3).

- По лемме о накачке для КС-языка L существуют числа n и m такие, что каждая цепочка $w \in L$ при $|w| > n$ представима в виде $w = xizvuy$, где
 - ▶ $uv \neq \varepsilon$,
 - ▶ $|uzv| \leq m$,
 - ▶ $xu^kzv^ky \in L$ для любого $k \geq 0$.
- Пусть $j = |uv|$; очевидно, $0 < j \leq m$. Для любой цепочки $w \in L$ при $|w| > n$ имеем $w = xzya^j$ и $xu^{k'+1}zv^{k'+1}y = wa^{jk'} \in L$ для любого $k' \geq 0$.
- Положим $n_0 = n + 1$ и $d = m!$. Тогда для любой цепочки $w \in L$ при $|w| \geq n_0$ выполняется $wa^d \in L$.
- Следовательно, n_0 и d можно взять в качестве индекса и периода множества M соответственно.

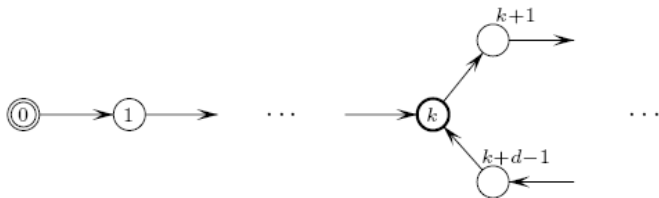
Теорема: продолжение доказательства

(3) влечёт (2).

- Если множество M конечно, то язык L конечен и потому регулярен. Далее считаем, что множество M бесконечно.
- Пусть n_0 и d являются соответственно индексом и периодом множества M .
- Для каждого i такого, что $0 \leq i < d$, определим (если оно существует) минимальное число k_i , удовлетворяющее условиям $k_i \in M$, $k_i \geq n_0$ и $k_i \equiv i \pmod{d}$. Хотя бы одно из k_i определено, поскольку M бесконечно.
- Обозначим через k наибольшее из всех определённых k_i .

Теорема: продолжение доказательства

- Построим ДКА \mathcal{A} , распознающий язык L . На диаграмме переходов автомата \mathcal{A} каждое ребро помечено символом a .



- Из начального состояния 0 до состояния n при $n < k$ читается только цепочка a^n , а при $n \geq k$ — множество всех цепочек вида a^{n+md} для каждого $m \geq 0$.
- Каждое состояние n такое, что $a^n \in L$, объявим заключительным.
- Таким образом, $L = L(\mathcal{A})$, значит, язык L является регулярным.

Теорема: окончание доказательства

(2) влечёт (1). Индукция по числу операций (объединения, произведения и итерации) в регулярном выражении, обозначающем язык L .

- База индукции. Языки \emptyset , $\{\varepsilon\}$ и $\{a\}$, обозначаемые регулярными выражениями без операций, контекстно-свободны.
- Индукционный переход. Пусть языки L_1 и L_2 порождаются КС-грамматиками $G_1 = (\Sigma_1, \Gamma_1, P_1, S_1)$ и $G_2 = (\Sigma_2, \Gamma_2, P_2, S_2)$. Можно считать, что $\Gamma_1 \cap \Gamma_2 = \emptyset$. Достаточно построить КС-грамматики, порождающие языки $L_1 \cup L_2$, $L_1 L_2$ и L_1^* . Выберем новый нетерминал $S \notin \Gamma_1 \cup \Gamma_2$.
 - ▶ $L_1 \cup L_2$ порождается КС-грамматикой с правилами $\{S \rightarrow S_1 | S_2\} \cup P_1 \cup P_2$.
 - ▶ $L_1 L_2$ порождается КС-грамматикой с правилами $\{S \rightarrow S_1 S_2\} \cup P_1 \cup P_2$.
 - ▶ L_1^* порождается КС-грамматикой с правилами $\{S \rightarrow S_1 S | \varepsilon\} \cup P_1$.



Литература

Основная литература

- Замятин А. П., Шур А. М. Языки, грамматики, распознаватели: Учебное пособие. Екатеринбург : Изд-во Урал. ун-та, 2007 (электронный вариант книги — на <http://elar.usu.ru>, поиск).

Дополнительная литература

- Ахо А., Лам М., Сети Р., Ульман Дж. Компиляторы: принципы, технологии и инструментарий. М.: ООО "И.Д. Вильямс", 2008.
- Ахо А., Ульман Дж. Теория синтаксического анализа, перевода и компиляции. М.: Мир, 1978.
- Мартыненко Б. К. Языки и трансляции: Учеб. пособие. СПб.: Издательство С.-Петербургского университета, 2004 (электронный вариант книги — на <http://www.math.spbu.ru/user/mbk>).