

Модели вычислений

В.А. Захаров, Р.И. Подловченко

Лекция 3.

1. Регулярные выражения. Алгебра регулярных выражений.
2. Теорема Клини о соответствии между регулярными выражениями и конечными автоматами.
3. Задача поиска по шаблону. Алгоритм Ахо-Корасик.
4. Двусторонние конечные автоматы.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

А МОЖНО ЛИ ОПИСАТЬ УСТРОЙСТВО
АВТОМАТНОГО ЯЗЫКА БЕЗ ПОМОЩИ
АВТОМАТОВ?

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

А МОЖНО ЛИ ОПИСАТЬ УСТРОЙСТВО АВТОМАТНОГО ЯЗЫКА БЕЗ ПОМОЩИ АВТОМАТОВ?

Существует и более удобный способ описания автоматных языков при помощи алгебраических выражений (формул) специального вида — регулярных выражений.

Именно регулярные выражения используются в большинстве компьютерных приложений — в текстовых редакторах, синтаксических анализаторах, интерпретаторах командных строк и др. — для описания простых формальных языков.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Регулярные выражения над алфавитом $\Sigma = \{a_1, a_2, \dots, a_n\}$ — это формулы, которые определяются над множеством констант

1. **0, 1**,
2. **a_1, a_2, \dots, a_n** ,

и множеством операций, состоящим из

1. двухместных операций **$\cdot, +$** ,
2. одноместной операции **$*$** .

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Регулярным выражением называется всякая формула, которая удовлетворяет следующим требованиям:

1. каждая константа является регулярным выражением;
2. если формулы R_1 и R_2 являются регулярными выражениями, то формулы

$$(R_1 \cdot R_2),$$

$$(R_1 + R_2),$$

$$(R_1^*)$$

также являются регулярными выражениями.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Примеры регулярных выражений.

$$(0 + 1),$$

$$(a_1 + a_2),$$

$$((a_1 \cdot a_2) + (a_2 \cdot a_1)),$$

$$(((a_1 \cdot ((a_1 + a_2)^*)) + a_2)^*).$$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Примеры регулярных выражений.

$$(0 + 1),$$

$$(a_1 + a_2),$$

$$((a_1 \cdot a_2) + (a_2 \cdot a_1)),$$

$$(((a_1 \cdot ((a_1 + a_2)^*)) + a_2)^*).$$

Для упрощения записи введем приоритет операций: высокий для $*$, средний для \cdot , низкий для $+$. Будем опускать некоторые скобки.

Например, запись

$$a_1 \cdot a_2^* + 1^* \cdot a_1$$

обозначает регулярное выражение

$$((a_1 \cdot (a_2^*)) + ((1^*) \cdot a_1)).$$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Значением каждого регулярного выражения R является формальный язык $L(R)$, $L(R) \subseteq \Sigma^*$, который определяется по следующим правилам:

1. $L(\mathbf{0}) = \emptyset$,
2. $L(\mathbf{1}) = \{\varepsilon\}$,
3. $L(\mathbf{a}_i) = \{a_i\}$, $1 \leq i \leq n$,
4. $L(R_1 \cdot R_2) = L(R_1)L(R_2)$ (конкатенация),
5. $L(R_1 + R_2) = L(R_1) \cup L(R_2)$ (объединение),
6. $L(R_1^*) = L^*(R_1)$ (итерация).

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Таким образом,

$$L(\mathbf{0} + \mathbf{1}) = \{\varepsilon\},$$

$$L(\mathbf{a}_1 + \mathbf{a}_2) = \{a_1, a_2\},$$

$$L(\mathbf{a}_1 \cdot \mathbf{a}_2 + \mathbf{a}_2 \cdot \mathbf{a}_1) = \{a_1a_2, a_2a_1\},$$

$$L((\mathbf{a}_1 \cdot ((\mathbf{a}_1 + \mathbf{a}_2)^*) + \mathbf{a}_2)^*) = \{a_1, a_2\}^* = L((\mathbf{a}_1 + \mathbf{a}_1)^*)$$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Язык называется **регулярным**, если он является значением некоторого регулярного выражения.

Например, язык L , состоящий из всех слов четной длины над алфавитом $\Sigma = \{a, b\}$, является регулярным, поскольку $L = L(((a + b) \cdot (a + b))^*)$.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ

Язык называется **регулярным**, если он является значением некоторого регулярного выражения.

Например, язык L , состоящий из всех слов четной длины над алфавитом $\Sigma = \{a, b\}$, является регулярным, поскольку $L = L(((a + b) \cdot (a + b))^*)$.

А как в этом убедиться?

Например, воспользоваться тождествами алгебры регулярных выражений.

Два регулярных выражения считаются равными, если их значения одинаковы.

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Утверждение 3.1. Для регулярных выражений справедливы следующие тождества

1. $F + G = G + F$;

2. $F + \mathbf{0} = F$;

3. $F + (G + H) = (F + G) + H$;

4. $F \cdot \mathbf{1} = F$;

5. $\mathbf{1} \cdot F = F$;

6. $F \cdot (G \cdot H) = (F \cdot G) \cdot H$;

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Утверждение 3.1. (Продолжение)

Для регулярных выражений справедливы следующие тождества

$$7 \quad F \cdot (G + H) = (F \cdot G) + (F \cdot H);$$

$$8 \quad (G + H) \cdot F = (G \cdot F) + (H \cdot F);$$

$$9 \quad F \cdot \mathbf{0} = \mathbf{0};$$

$$10 \quad \mathbf{0} \cdot F = \mathbf{0};$$

$$11 \quad F + F = F;$$

Доказательство. Самостоятельно

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Утверждение 3.1. (Продолжение)

Для регулярных выражений справедливы следующие тождества

$$7 \quad F \cdot (G + H) = (F \cdot G) + (F \cdot H);$$

$$8 \quad (G + H) \cdot F = (G \cdot F) + (H \cdot F);$$

$$9 \quad F \cdot \mathbf{0} = \mathbf{0};$$

$$10 \quad \mathbf{0} \cdot F = \mathbf{0};$$

$$11 \quad F + F = F;$$

Доказательство. Самостоятельно

Как видно из этих тождеств, регулярные выражения образуют ассоциативное полукольцо.

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Упростим регулярное выражение

$$(a + b)(a + b) + aa + bb =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Упростим регулярное выражение

$$(a + b)(a + b) + aa + bb =$$

$$(a + b) \cdot a + (a + b) \cdot b + aa + bbb =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Упростим регулярное выражение

$$(a + b)(a + b) + aa + bb =$$

$$(a + b) \cdot a + (a + b) \cdot b + aa + bbb =$$

$$aa + ba + ab + bb + aa + bb =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Упростим регулярное выражение

$$(a + b)(a + b) + aa + bb =$$

$$(a + b) \cdot a + (a + b) \cdot b + aa + bbb =$$

$$aa + ba + ab + bb + aa + bb =$$

$$aa + aa + bb + bb + ba + ab =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Упростим регулярное выражение

$$(a + b)(a + b) + aa + bb =$$

$$(a + b) \cdot a + (a + b) \cdot b + aa + bbb =$$

$$aa + ba + ab + bb + aa + bb =$$

$$aa + aa + bb + bb + ba + ab =$$

$$aa + bb + ba + ab =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Упростим регулярное выражение

$$(a + b)(a + b) + aa + bb =$$

$$(a + b) \cdot a + (a + b) \cdot b + aa + bbb =$$

$$aa + ba + ab + bb + aa + bb =$$

$$aa + aa + bb + bb + ba + ab =$$

$$aa + bb + ba + ab =$$

$$(a + b)(a + b)$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Утверждение 3.2. Для регулярных выражений справедливы следующие тождества

$$1. F \cdot F^* = F^* \cdot F;$$

$$2. (F^*)^* = F^*;$$

$$3. F^* = \mathbf{1} + F \cdot F^*;$$

$$4. F^* = (\mathbf{1} + F + FF + \dots + F^{n-1}) \cdot (F^n)^*;$$

$$5. (F + G)^* = (F^* \cdot G)^* \cdot F^*;$$

$$6. (F \cdot G)^* = \mathbf{1} + F(G \cdot F)^* G.$$

Доказательство. Самостоятельно

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Используя указанные тождества, можно упрощать регулярные выражения и решать уравнения над регулярными выражениями.

Пример. Упростить регулярное выражение

$$(a^*b)^* + (b^*a)^*$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Используя указанные тождества, можно упрощать регулярные выражения и решать уравнения над регулярными выражениями.

Пример. Упростить регулярное выражение

$$(a^*b)^* + (b^*a)^*$$

$$(a^*b)^* + (b^*a)^* = \mathbf{1} + a^*b(a^*b)^* + \mathbf{1} + b^*a(b^*a)^* =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Используя указанные тождества, можно упрощать регулярные выражения и решать уравнения над регулярными выражениями.

Пример. Упростить регулярное выражение

$$(a^*b)^* + (b^*a)^*$$

$$(a^*b)^* + (b^*a)^* = \mathbf{1} + a^*b(a^*b)^* + \mathbf{1} + b^*a(b^*a)^* = \\ \mathbf{1} + \mathbf{1} + (a^*b)^*a^*b + (b^*a)^*b^*a =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Используя указанные тождества, можно упрощать регулярные выражения и решать уравнения над регулярными выражениями.

Пример. Упростить регулярное выражение

$$(a^*b)^* + (b^*a)^*$$

$$(a^*b)^* + (b^*a)^* = \mathbf{1} + a^*b(a^*b)^* + \mathbf{1} + b^*a(b^*a)^* =$$

$$\mathbf{1} + \mathbf{1} + (a^*b)^*a^*b + (b^*a)^*b^*a =$$

$$\mathbf{1} + (a + b)^*b + (b + a)^*a = \mathbf{1} + (b + a)^*b + (b + a)^*a =$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Используя указанные тождества, можно упрощать регулярные выражения и решать уравнения над регулярными выражениями.

Пример. Упростить регулярное выражение

$$(a^*b)^* + (b^*a)^*$$

$$\begin{aligned}(a^*b)^* + (b^*a)^* &= \mathbf{1} + a^*b(a^*b)^* + \mathbf{1} + b^*a(b^*a)^* = \\ \mathbf{1} + \mathbf{1} + (a^*b)^*a^*b + (b^*a)^*b^*a &= \\ \mathbf{1} + (a + b)^*b + (b + a)^*a &= \mathbf{1} + (b + a)^*b + (b + a)^*a = \\ \mathbf{1} + (b + a)^*(b + a) &= \mathbf{1} + (b + a)(b + a)^* = (b + a)^*\end{aligned}$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Пусть F и G — произвольные регулярные выражения.

Рассмотрим уравнение $X = X \cdot F + G$ над регулярными языками с неизвестной X .

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Пусть F и G — произвольные регулярные выражения.

Рассмотрим уравнение $X = X \cdot F + G$ над регулярными языками с неизвестной X .

Утверждение 3.3. Регулярное выражение $G \cdot F^*$ является решением уравнения $X = X \cdot F + G$.

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Пусть F и G — произвольные регулярные выражения.

Рассмотрим уравнение $X = X \cdot F + G$ над регулярными языками с неизвестной X .

Утверждение 3.3. Регулярное выражение $G \cdot F^*$ является решением уравнения $X = X \cdot F + G$.

Доказательство.

$$G \cdot F^* \cdot F + G = G \cdot (\mathbf{1} + F^* \cdot F) = G \cdot F^*$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

ЗАДАЧА 8

Доказать, что в случае $\varepsilon \notin L(F)$, это решение единственное.

Какие еще решения имеет уравнение

$\mathcal{X} = \mathcal{X} \cdot F + G$ в случае $\varepsilon \in L(F)$?

ЗАДАЧА 9

Какие решения имеет уравнение $\mathcal{X} = F \cdot \mathcal{X} + G$?

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Перед нами открывается интересная возможность:
языки можно задавать неявно, как решения
уравнений в алгебре регулярных выражений!

Например, язык L определяется как наименьшее
решение уравнения

$$(\mathcal{X} \cdot a + b \cdot \mathcal{X})^* = \mathbf{1} + b \cdot \mathcal{X}^2 \cdot a.$$

АЛГЕБРА РЕГУЛЯРНЫХ ВЫРАЖЕНИЙ

Перед нами открывается интересная возможность: языки можно задавать неявно, как решения уравнений в алгебре регулярных выражений!

Например, язык L определяется как наименьшее решение уравнения

$$(X \cdot a + b \cdot X)^* = 1 + b \cdot X^2 \cdot a.$$

ЗАДАЧА 10 [Трудная] Докажите или опровергните гипотезу:

Каково бы ни было уравнение над регулярными выражениями с одной переменной, если это уравнение имеет решение, то хотя бы одно из его минимальных решений — это регулярный язык.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

КАК СООТНОСЯТСЯ ДРУГ С ДРУГОМ
КЛАССЫ АВТОМАТНЫХ И РЕГУЛЯРНЫХ
ЯЗЫКОВ?

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Теорема 3.4. Каждый регулярный язык является автоматным языком.

Доказательство. 1) Языки, которые задаются константами **0, 1** и **a**, $a \in \Sigma$, очевидно являются автоматными.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Теорема 3.4. Каждый регулярный язык является автоматным языком.

Доказательство. 1) Языки, которые задаются константами **0**, **1** и **a**, $a \in \Sigma$, очевидно являются автоматными.

2) Класс автоматных языков замкнут относительно операции объединения **+**.

Покажем, что он замкнут относительно операций конкатенации и итерации Клини.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Пусть $L_i = L(\mathcal{A}_i)$, $i = 1, 2, \dots$, где $\mathcal{A}_i = (\Sigma, S_i, I_i, F_i, T_i)$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Пусть $L_i = L(\mathcal{A}_i)$, $i = 1, 2$, где $\mathcal{A}_i = (\Sigma, S_i, I_i, F_i, T_i)$

Рассмотрим конечный автомат

$\mathcal{A}_0 = (\Sigma, S_1 \cup S_2, I_1, F_2, T_1 \cup T_2 \cup T_0)$, где

$T_0 = \{(s', \varepsilon, s'') : s' \in F_1, s'' \in I_2\}$.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Пусть $L_i = L(\mathcal{A}_i)$, $i = 1, 2$, где $\mathcal{A}_i = (\Sigma, S_i, I_i, F_i, T_i)$

Рассмотрим конечный автомат

$\mathcal{A}_0 = (\Sigma, S_1 \cup S_2, I_1, F_2, T_1 \cup T_2 \cup T_0)$, где

$T_0 = \{(s', \varepsilon, s'') : s' \in F_1, s'' \in I_2\}$.

Тогда $w \in L_1 \cdot L_2 \Leftrightarrow w = uv$, $u \in L_1$, $v \in L_2 \Leftrightarrow$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Пусть $L_i = L(\mathcal{A}_i)$, $i = 1, 2$, где $\mathcal{A}_i = (\Sigma, S_i, I_i, F_i, T_i)$

Рассмотрим конечный автомат

$\mathcal{A}_0 = (\Sigma, S_1 \cup S_2, I_1, F_2, T_1 \cup T_2 \cup T_0)$, где

$T_0 = \{(s', \varepsilon, s'') : s' \in F_1, s'' \in I_2\}$.

Тогда $w \in L_1 \cdot L_2 \Leftrightarrow w = uv$, $u \in L_1$, $v \in L_2 \Leftrightarrow$

существуют успешные вычисления $s'_0 \xrightarrow{u}^* s'$ и

$s''_0 \xrightarrow{v}^* s''$ автоматов \mathcal{A}_1 и $\mathcal{A}_2 \Leftrightarrow$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Пусть $L_i = L(\mathcal{A}_i)$, $i = 1, 2$, где $\mathcal{A}_i = (\Sigma, S_i, I_i, F_i, T_i)$

Рассмотрим конечный автомат

$\mathcal{A}_0 = (\Sigma, S_1 \cup S_2, I_1, F_2, T_1 \cup T_2 \cup T_0)$, где

$T_0 = \{(s', \varepsilon, s'') : s' \in F_1, s'' \in I_2\}$.

Тогда $w \in L_1 \cdot L_2 \Leftrightarrow w = uv$, $u \in L_1$, $v \in L_2 \Leftrightarrow$
существуют успешные вычисления $s'_0 \xrightarrow{u} s'$ и
 $s''_0 \xrightarrow{v} s''$ автоматов \mathcal{A}_1 и $\mathcal{A}_2 \Leftrightarrow$ существует
успешное вычисление $s'_0 \xrightarrow{u} s' \xrightarrow{\varepsilon} s''_0 \xrightarrow{v} s''$
автомата $\mathcal{A}_0 \Leftrightarrow$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Пусть $L_i = L(\mathcal{A}_i)$, $i = 1, 2$, где $\mathcal{A}_i = (\Sigma, S_i, I_i, F_i, T_i)$

Рассмотрим конечный автомат

$\mathcal{A}_0 = (\Sigma, S_1 \cup S_2, I_1, F_2, T_1 \cup T_2 \cup T_0)$, где

$T_0 = \{(s', \varepsilon, s'') : s' \in F_1, s'' \in I_2\}$.

Тогда $w \in L_1 \cdot L_2 \Leftrightarrow w = uv$, $u \in L_1$, $v \in L_2 \Leftrightarrow$
существуют успешные вычисления $s'_0 \xrightarrow{u}_* s'$ и
 $s''_0 \xrightarrow{v}_* s''$ автоматов \mathcal{A}_1 и $\mathcal{A}_2 \Leftrightarrow$ существует
успешное вычисление $s'_0 \xrightarrow{u}_* s' \xrightarrow{\varepsilon} s''_0 \xrightarrow{v}_* s''$
автомата $\mathcal{A}_0 \Leftrightarrow w \in L(\mathcal{A}_0)$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Покажите самостоятельно, что в случае $L = L(\mathcal{A})$ существует такой конечный автомат \mathcal{A}' , для которого $L^* = L(\mathcal{A}')$, и завершите тем самым доказательство теоремы.

QED

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Теорема 3.5. Каждый автоматный язык является регулярным языком.

Доказательство. Пусть $\mathcal{A} = (\Sigma, S, I, F, T)$ — конечный автомат без ε -переходов, и при этом $I \cap F = \emptyset$. Построим регулярное выражение, описывающее автоматный язык $L(\mathcal{A})$.

Для этого составим и решим систему уравнений над регулярными выражениями.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Теорема 3.5. Каждый автоматный язык является регулярным языком.

Доказательство. Пусть $A = (\Sigma, S, I, F, T)$ — конечный автомат без ε -переходов, и при этом $I \cap F = \emptyset$. Построим регулярное выражение, описывающее автоматный язык $L(A)$.

Для этого составим и решим систему уравнений над регулярными выражениями.

1). Для каждого состояния $s, s \in S$, введем переменную \mathcal{X}_s , а также множество пар $In(s) = \{(s', a) : s' \xrightarrow{a} s \in T\}$.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

2). Для каждого состояния s , $s \in S$, составим уравнение

$$\mathcal{X}_s = \sum_{(s',a) \in \text{In}(s)} \mathcal{X}_{s'} \cdot \mathbf{a} + \delta_s,$$

где $\delta_s = \mathbf{1}$, если $s \in I$, и $\delta_s = \mathbf{0}$ иначе.

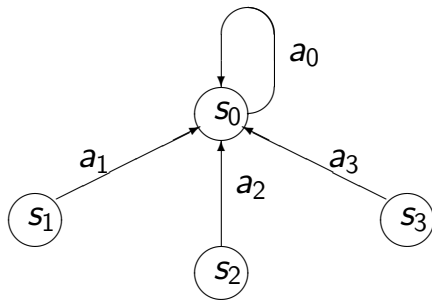
РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

А откуда взялись эти уравнения?

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

А откуда взялись эти уравнения?

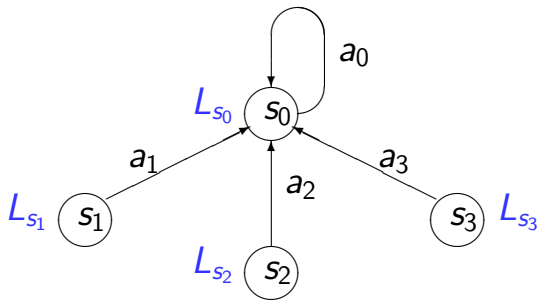
Рассмотрим состояние s_0 автомата.



РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

А откуда взялись эти уравнения?

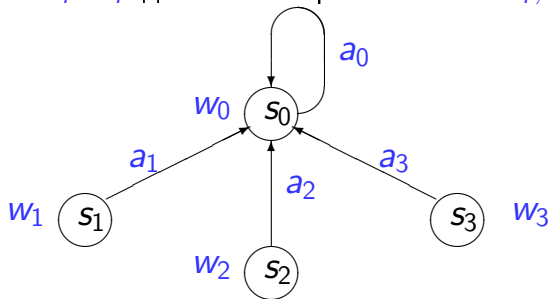
Рассмотрим состояние s_0 автомата.



И пусть L_{s_i} — это язык, состоящий из слов, которые прочитывает автомат, достигнув состояния s_i , $0 \leq i \leq 3$.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

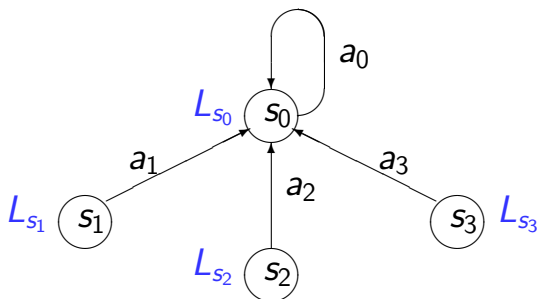
Тогда каждое слово w языка L_{S_0} представимо в виде $w = w_j \cdot a_j$ для некоторого слова w_j , $w_j \in L_j$.



РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Тогда справедливо равенство

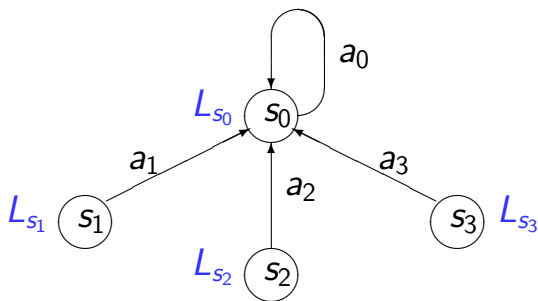
$$L_{s_0} = L_{s_0} \cdot a_0 + L_{s_1} \cdot a_1 + L_{s_2} \cdot a_2 + L_{s_3} \cdot a_3$$



РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Тогда справедливо равенство

$$L_{s_0} = L_{s_0} \cdot a_0 + L_{s_1} \cdot a_1 + L_{s_2} \cdot a_2 + L_{s_3} \cdot a_3$$



Отсюда и возникает уравнение

$$\mathcal{X}_{s_0} = \mathcal{X}_{s_0} \cdot a_0 + \mathcal{X}_{s_1} \cdot a_1 + \mathcal{X}_{s_2} \cdot a_2 + \mathcal{X}_{s_3} \cdot a_3$$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

3). Найдем решение этой системы уравнений

$$\{\mathcal{X}_s = R_s : s \in S\},$$

методом Гаусса (исключением переменных),
используя утверждение 3.3.

Согласно этому утверждению (а также
последующей задаче) такое решение будет
единственным.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

4) Индукцией по длине слова w покажем, что $w \in L(R_s)$ тогда и только тогда, когда слово w допускается конечным автоматом

$\mathcal{A}^s = (\Sigma, S, I, \{s\}, T)$: это легко видно из устройства уравнений

$$\mathcal{X}_s = \sum_{(s', a) \in In(s)} \mathcal{X}_{s'} \cdot \mathbf{a} + \delta_s,$$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

4) Индукцией по длине слова w покажем, что $w \in L(R_s)$ тогда и только тогда, когда слово w допускается конечным автоматом

$\mathcal{A}^s = (\Sigma, S, I, \{s\}, T)$: это легко видно из устройства уравнений

$$x_s = \sum_{(s', a) \in \text{In}(s)} x_{s'} \cdot a + \delta_s,$$

5). Тогда $L(\mathcal{A}) = L(\sum_{s \in F} R_s)$.

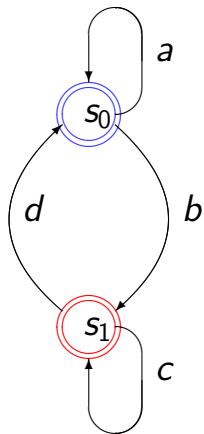
QED

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

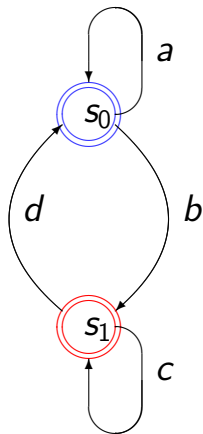
Задача 10. А зачем мне понадобилось в начале доказательства оговорить условие $I \cap F = \emptyset$?

И как нужно модифицировать доказательство теоремы 3.5, если это условие не выполнено?

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ



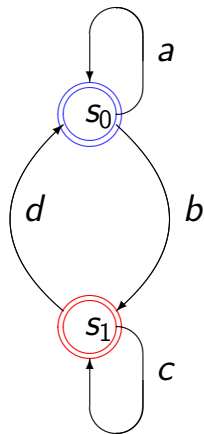
РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ



$$\mathcal{X}_0 = \mathcal{X}_0 \cdot a + \mathcal{X}_1 \cdot d + 1$$

$$\mathcal{X}_1 = \mathcal{X}_0 \cdot b + \mathcal{X}_1 \cdot c + 0$$

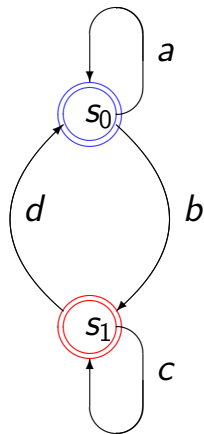
РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ



$$\mathcal{X}_0 = \mathcal{X}_0 \cdot a + (\mathcal{X}_1 \cdot d + 1)$$

$$\mathcal{X}_1 = \mathcal{X}_0 \cdot b + \mathcal{X}_1 \cdot c$$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

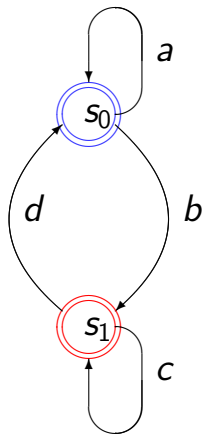


$$\mathcal{X}_0 = (\mathcal{X}_1 \cdot d + 1) \cdot a^*$$

$$\mathcal{X}_1 = \mathcal{X}_0 \cdot b + \mathcal{X}_1 \cdot c$$

Утверждение 3.3 о решении уравнений

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

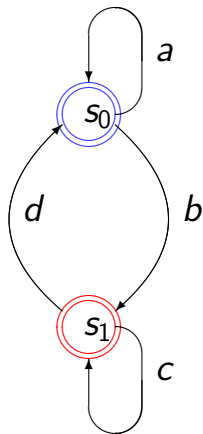


$$\mathcal{X}_0 = (\mathcal{X}_1 \cdot d + 1) \cdot a^*$$

$$\mathcal{X}_1 = (\mathcal{X}_1 \cdot d + 1) \cdot a^* \cdot b + \mathcal{X}_1 \cdot c$$

Исключение переменной \mathcal{X}_0

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

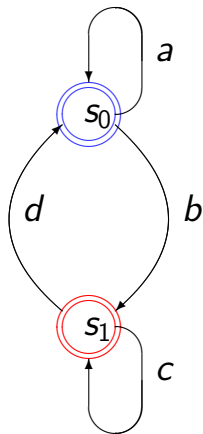


$$\mathcal{X}_0 = (\mathcal{X}_1 \cdot d + 1) \cdot a^*$$

$$\mathcal{X}_1 = \mathcal{X}_1 \cdot (da^*b + c) + a^*b$$

Приведение подобных слагаемых

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

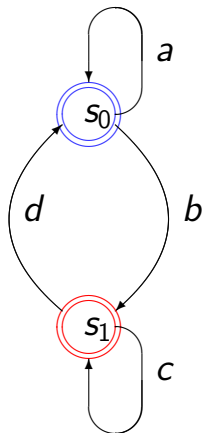


$$\mathcal{X}_0 = (\mathcal{X}_1 \cdot d + 1) \cdot a^*$$

$$\mathcal{X}_1 = a^* b (d a^* b + c)^*$$

Утверждение 3.3 о решении уравнений

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ



Искомое регулярное выражение

$$R = a^*b(da^*b + c)^*$$

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Теорема 3.6. [С. Клини] Формальный язык является автоматным тогда и только тогда, когда он является регулярным.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Теорема 3.6. [С. Клини] Формальный язык является автоматным тогда и только тогда, когда он является регулярным.

Эта теорема очень важна для построения алгоритмов поиска строк (текстов, потоков данных и пр.), подпадающих под заданный шаблон.

Области применения: интернет-поиск, построение таблиц маршрутизации, обнаружение вредоносного кода, выявление плагиата, и пр.

РЕГУЛЯРНЫЕ ВЫРАЖЕНИЯ И КОНЕЧНЫЕ АВТОМАТЫ

Принцип применения:

1. Интересующее множество строк описываем в виде шаблона P , которому соответствует некоторое регулярное выражение R .
2. Для регулярного выражения R строим (недетерминированный) конечный автомат A .
3. Проводим детерминизацию автомата A и получаем детерминированный автомат A_0 .
4. На вход автомата A_0 подаем строку w . Автомат A_0 достигает финального состояния тогда и только тогда, когда w подпадает под шаблон P .

ЗАДАЧА ПОИСКА ПО ШАБЛОНУ

Основные конструкции языка описания шаблонов:

- ▶ $?$ — любая буква;
- ▶ $*$ — любое слово;
- ▶ $\{a_1, a_2, a_3 : +\}$ — любая из букв a_1, a_2, a_3 ;
- ▶ $\{b_1, b_2 : -\}$ — любая из букв кроме b_1, b_2 ;
- ▶ $R[5 : 8]$ — шаблон R , повторяющийся подряд не менее 5 и не более 8 раз;
- ▶ $R_1 \& R_2$ — пересечение шаблонов R_1 и R_2 ;
- ▶ $R_1 + R_2$ — альтернативный выбор (объединение) шаблонов R_1 и R_2 .

ЗАДАЧА ПОИСКА ПО ШАБЛОНУ

Пример. Пусть $\Sigma = \{a, b, c\}$. Тогда шаблону

$$\{b : -\}[1 : 2] * b?$$

соответствует регулярное выражение

$$((a + c) + (a + c)(a + c))(a + b + c)^* b(a + b + c)$$

ЗАДАЧА ПОИСКА ПО ШАБЛОНУ

Наибольшую вычислительную трудность имеет этап преобразования недетерминированного конечного автомата в детерминированный, поскольку здесь возможен эффект «комбинаторного взрыва» числа состояний.

Однако для некоторых широко используемых шаблонов соответствующие детерминированные автоматы имеют простое устройство и строятся довольно эффективно.

АЛГОРИТМ АХО—КОРАСИК

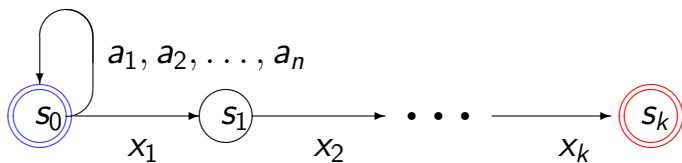
Рассмотрим задачу поиска вхождений заданного слова $w = x_1x_2 \dots x_k$ в произвольной строке.

АЛГОРИТМ АХО—КОРАСИК

Рассмотрим задачу поиска вхождений заданного слова $w = x_1x_2 \dots x_k$ в произвольной строке.

Формально эту задачу можно поставить так: построить алгоритм, распознающий все слова, имеющие суффикс w .

Слова этого языка описываются шаблоном $*x_1x_2 \dots x_k$ и распознаются недетерминированным конечным автоматом простого вида



АЛГОРИТМ АХО—КОРАСИК

Но как построить детерминированный автомат?

АЛГОРИТМ АХО—КОРАСИК

Но как построить детерминированный автомат?

Удобное решение придумали в 1965 году Альфред Ахо и Маргерет Корасик.

Минимальный детерминированный автомат, распознающий все тексты, оканчивающиеся строкой $w = x_1x_2 \dots x_k$, имеет вид

$A_w = (\Sigma, S, \{s_0\}, \{s_k\}, T)$, где

▶ $S = \{s_0, s_1, \dots, s_k\}$;

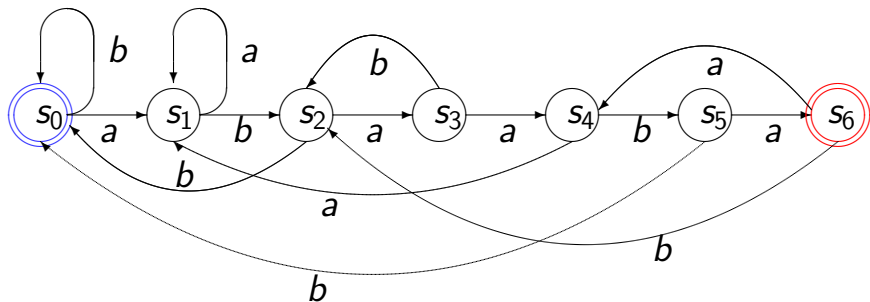
▶ для любых $s_i, s_j \in S$, и $y, y \in \Sigma$, в множестве T есть переход $s_i \xrightarrow{y} s_j$, где

$j = \max(\ell : x_1x_2 \dots x_\ell \in \text{Suff}(x_1x_2 \dots x_iy))$:

«ищи максимальный префикс w , являющийся суффиксом прочитанной части текста»

АЛГОРИТМ АХО—КОРАСИК

Если $\Sigma = \{a, b\}$ и $w = abaaba$, то минимальный детерминированный автомат A_w , распознающий язык $(a + b)^* abaaba$ имеет следующий вид.



АЛГОРИТМ АХО—КОРАСИК

Предложенный алгоритм строит автомат \mathcal{A}_w за время $O(|w|^2)$ и обнаруживает все вхождения строки w в тексте длины n за время $O(n)$.

ЗАДАЧА 11.

Разработайте обобщение предложенного метода, которое позволяет для произвольного множества строк $D = \{w_1, w_2, \dots, w_m\}$ построить минимальный детерминированный конечный автомат \mathcal{A}_D , распознающий все слова, оканчивающиеся хотя бы одной из строк множества D .

ТЕХНОЛОГИЯ ПРИМЕНЕНИЯ КОНЕЧНЫХ АВТОМАТОВ

**Сформулируйте задачу поиска или распознавания
на языке словарных шаблонов**

ТЕХНОЛОГИЯ ПРИМЕНЕНИЯ КОНЕЧНЫХ АВТОМАТОВ

Сформулируйте задачу поиска или распознавания
на языке словарных шаблонов



Транслируйте шаблон
в регулярное выражение

ТЕХНОЛОГИЯ ПРИМЕНЕНИЯ КОНЕЧНЫХ АВТОМАТОВ

Сформулируйте задачу поиска или распознавания на языке словарных шаблонов



Транслируйте шаблон в регулярное выражение



Постройте для регулярного выражения недетерминированный автомат

ТЕХНОЛОГИЯ ПРИМЕНЕНИЯ КОНЕЧНЫХ АВТОМАТОВ

Сформулируйте задачу поиска или распознавания на языке словарных шаблонов

Транслируйте шаблон в регулярное выражение

Постройте для регулярного выражения недетерминированный автомат

Детерминизируйте этот автомат

ТЕХНОЛОГИЯ ПРИМЕНЕНИЯ КОНЕЧНЫХ АВТОМАТОВ

Сформулируйте задачу поиска или распознавания на языке словарных шаблонов

Транслируйте шаблон в регулярное выражение

Постройте для регулярного выражения недетерминированный автомат

Детерминизируйте этот автомат

Минимизируйте полученный автомат

ТЕХНОЛОГИЯ ПРИМЕНЕНИЯ КОНЕЧНЫХ АВТОМАТОВ

Сформулируйте задачу поиска или распознавания на языке словарных шаблонов

Транслируйте шаблон в регулярное выражение

Постройте для регулярного выражения недетерминированный автомат

Детерминизируйте этот автомат

Минимизируйте полученный автомат

Вы построили оптимальную программу поиска!

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Конечный автомат можно рассматривать как механическое устройство, которое движется в одну сторону по ленте, на которой записано заданное слово в алфавите Σ . Лента оканчивается **граничным маркером** \vdash , $\vdash \notin \Sigma$. Слово допускается автоматом, если при достижении граничного маркера автомат оказывается в финальном состоянии.



\uparrow
 q_0



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Конечный автомат можно рассматривать как механическое устройство, которое движется в одну сторону по ленте, на которой записано заданное слово в алфавите Σ . Лента оканчивается **граничным маркером** \vdash , $\vdash \notin \Sigma$. Слово допускается автоматом, если при достижении граничного маркера автомат оказывается в финальном состоянии.



q_{i_1}



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Конечный автомат можно рассматривать как механическое устройство, которое движется в одну сторону по ленте, на которой записано заданное слово в алфавите Σ . Лента оканчивается **граничным маркером** \vdash , $\vdash \notin \Sigma$. Слово допускается автоматом, если при достижении граничного маркера автомат оказывается в финальном состоянии.

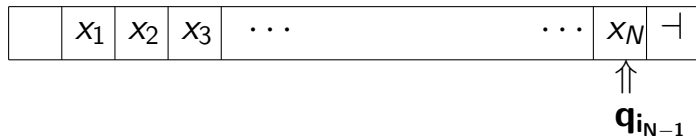


q_{i_2}



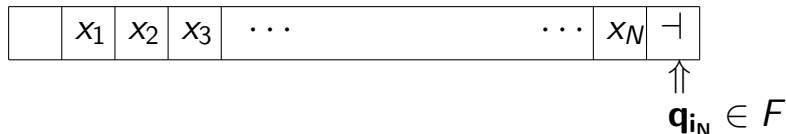
ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Конечный автомат можно рассматривать как механическое устройство, которое движется в одну сторону по ленте, на которой записано заданное слово в алфавите Σ . Лента оканчивается **граничным маркером** \vdash , $\vdash \notin \Sigma$. Слово допускается автоматом, если при достижении граничного маркера автомат оказывается в финальном состоянии.



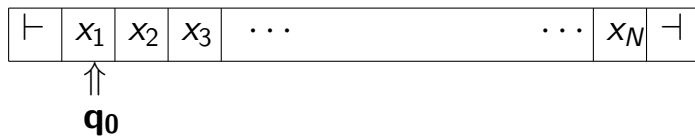
ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Конечный автомат можно рассматривать как механическое устройство, которое движется в одну сторону по ленте, на которой записано заданное слово в алфавите Σ . Лента оканчивается **граничным маркером** \vdash , $\vdash \notin \Sigma$. Слово допускается автоматом, если при достижении граничного маркера автомат оказывается в финальном состоянии.



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

А что произойдет, если разрешить такому автомату перемещать считывающую головку по ленте в обе стороны?



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

А что произойдет, если разрешить такому автомату перемещать считывающую головку по ленте в обе стороны?



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

А что произойдет, если разрешить такому автомату перемещать считывающую головку по ленте в обе стороны?



↑↑
 q_{i_2}



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

А что произойдет, если разрешить такому автомату перемещать считывающую головку по ленте в обе стороны?



↑↑
 q_{i_3}



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

А что произойдет, если разрешить такому автомату перемещать считывающую головку по ленте в обе стороны?



↑↑
 q_i



ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

А что произойдет, если разрешить такому автомату перемещать считывающую головку по ленте в обе стороны?



Насколько расширится класс языков, распознаваемых такими двусторонними автоматами?

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Двусторонним конечным автоматом называется система $\mathcal{B} = (\Sigma, S, I, F, T)$, в которой

- ▶ Σ — ленточный алфавит, S — конечное множество состояний, I, F — подмножества начальных и финальных состояний;
- ▶ $T \subseteq S \times \Sigma \cup \{\vdash, \dashv\} \times S \times \{-, +\}$ — отношение переходов.

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Двусторонним конечным автоматом называется система $B = (\Sigma, S, I, F, T)$, в которой

- ▶ Σ — ленточный алфавит, S — конечное множество состояний, I, F — подмножества начальных и финальных состояний;
- ▶ $T \subseteq S \times \Sigma \cup \{\vdash, \dashv\} \times S \times \{-, +\}$ — отношение переходов.

Четверка (переход) (s', y, s'', δ) означает, что автомат, пребывающий в состоянии s' и обзорающий букву или граничный маркер y , переходит в состояние s'' и перемещает считывающую головку к соседнему символу по направлению δ .

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Более формально вычисление двустороннего автомата определяется так.

Конфигурацией конечного двустороннего автомата $\mathcal{B} = (\Sigma, S, I, F, T)$ называется всякое слово одного из следующих видов: $s \vdash w \dashv$, $\vdash usv \dashv$, где w, u, v — слова, и s — состояние.

Конфигурация $\vdash sw \dashv$ называется начальной, если $s \in I$.

Конфигурация называется финальной, если $s \in F$.

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Отношение переходов $\alpha \longrightarrow \beta$ на множестве конфигураций автомата \mathcal{B} выполняется в том и только том случае, если

1. $\alpha = s \vdash w \dashv$, $\beta = \vdash s'w \dashv$, и $(s, \vdash, s', +1) \in T$;
2. $\alpha = \vdash usav \dashv$, $\beta = \vdash uas'v \dashv$, и $(s, a, s', +1) \in T$;
3. $\alpha = \vdash ubsav \dashv$, $\beta = \vdash us'bav \dashv$, и $(s, a, s', -1) \in T$
4. $\alpha = \vdash wbs \dashv$, $\beta = \vdash ws'b \dashv$, и $(s, \dashv, s', -1) \in T$.

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Отношение переходов $\alpha \longrightarrow \beta$ на множестве конфигураций автомата \mathcal{B} выполняется в том и только том случае, если

1. $\alpha = s \vdash w \dashv$, $\beta = \vdash s' w \dashv$, и $(s, \vdash, s', +1) \in T$;
2. $\alpha = \vdash u s a v \dashv$, $\beta = \vdash u a s' v \dashv$, и $(s, a, s', +1) \in T$;
3. $\alpha = \vdash u b s a v \dashv$, $\beta = \vdash u s' b a v \dashv$, и $(s, a, s', -1) \in T$;
4. $\alpha = \vdash w b s \dashv$, $\beta = \vdash w s' b \dashv$, и $(s, \dashv, s', -1) \in T$.

Вычисление двустороннего автомата \mathcal{B} — это последовательность переходов

$$\alpha_0 \longrightarrow \alpha_1 \longrightarrow \alpha_2 \longrightarrow \cdots \longrightarrow \alpha_n.$$

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Двусторонний конечный автомат \mathcal{B} допускает слово w , если существует такое вычисление этого автомата

$$\alpha_0 \longrightarrow \alpha_1 \longrightarrow \alpha_2 \longrightarrow \cdots \longrightarrow \alpha_n,$$

в котором $\alpha_0 = \vdash q_0 w \dashv$ — начальная конфигурация, а α_n — финальная конфигурация.

Язык $L(\mathcal{B})$ двустороннего конечного автомата \mathcal{B} — это множество всех слов, которые допускает этот автомат.

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

ЗАДАЧА 10 [Трудная]

Доказать, что для любого двустороннего конечного автомата \mathcal{B} язык $L(\mathcal{B})$ — это регулярный (автоматный) язык.

Таким образом, увеличение «подвижности» конечного автомата не приводит к расширению его вычислительных возможностей.

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

ЗАДАЧА 10 [Трудная]

Доказать, что для любого двустороннего конечного автомата \mathcal{B} язык $L(\mathcal{B})$ — это регулярный (автоматный) язык.

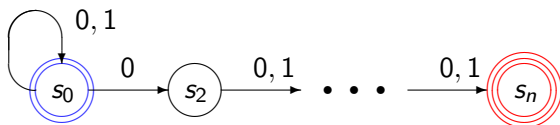
Таким образом, увеличение «подвижности» конечного автомата не приводит к расширению его вычислительных возможностей.

А упрощается ли при этом распознавание языков?

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Рассмотрим недетерминированный конечный автомат

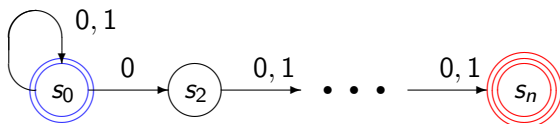
\mathcal{A}_n , $n \geq 1$,



который распознает язык, состоящий из всех двоичных слов, содержащих 0 на n -ой позиции справа.

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

Рассмотрим недетерминированный конечный автомат \mathcal{A}_n , $n \geq 1$,



который распознает язык, состоящий из всех двоичных слов, содержащих 0 на n -ой позиции справа.

Тогда минимальный детерминированный односторонний конечный автомат, эквивалентный автомату \mathcal{A}_n , имеет 2^n состояний, а минимальный детерминированный двусторонний конечный автомат, распознающий тот же язык, имеет $n + 2$ состояние.

Таким образом, двусторонние конечные автоматы могут быть гораздо компактнее односторонних. Однако до конца этот вопрос все еще не решен.

ДВУСТОРОННИЕ КОНЕЧНЫЕ АВТОМАТЫ

ОТКРЫТАЯ ПРОБЛЕМА

Существует ли такой регулярный язык, для распознавания которого можно построить недетерминированный двусторонний автомат, имеющий существенно меньший размер, нежели минимальный недетерминированный односторонний автомат, распознающий тот же язык?

КОНЕЦ ЛЕКЦИИ 3