

# Теория автоматов и формальных языков

## Регулярные языки

**Лектор:** Екатерина Вербицкая

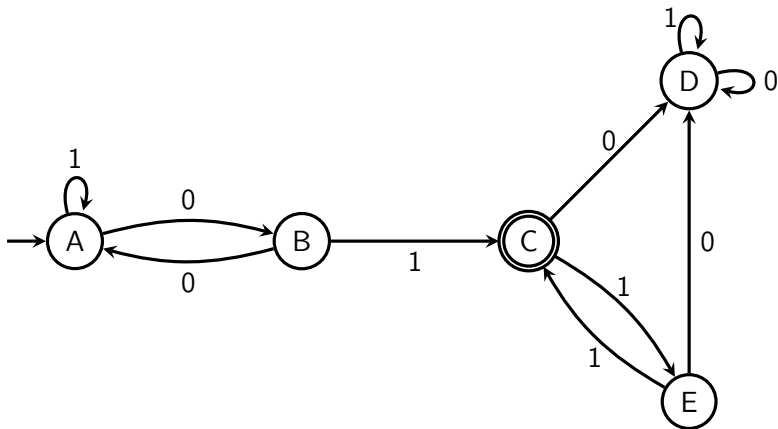
Санкт-Петербургский государственный электротехнический университет «ЛЭТИ»

5 октября 2021

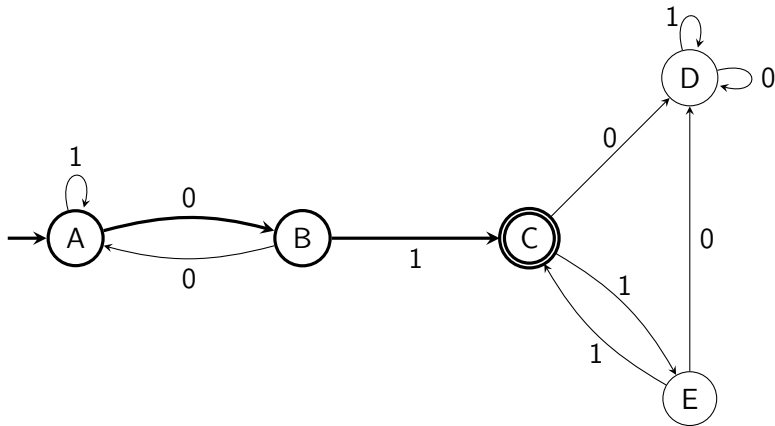
**Детерминированный конечный автомат** —  $\langle Q, \Sigma, \delta, q_0, F \rangle$

- $Q \neq \emptyset$  — конечное множество состояний
- $\Sigma$  — Конечный входной алфавит
- $\delta$  — функция переходов: отображение типа  $Q \times \Sigma \rightarrow Q$
- $q_0 \in Q$  — начальное состояние
- $F \subseteq Q$  — множество конечных состояний

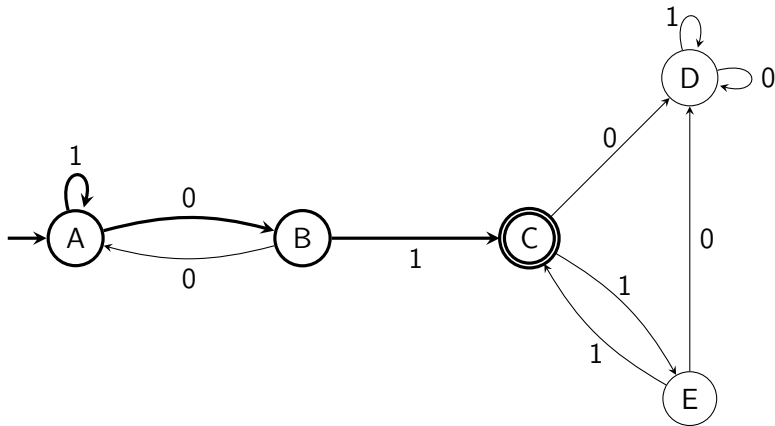
В предыдущей серии: ДКА



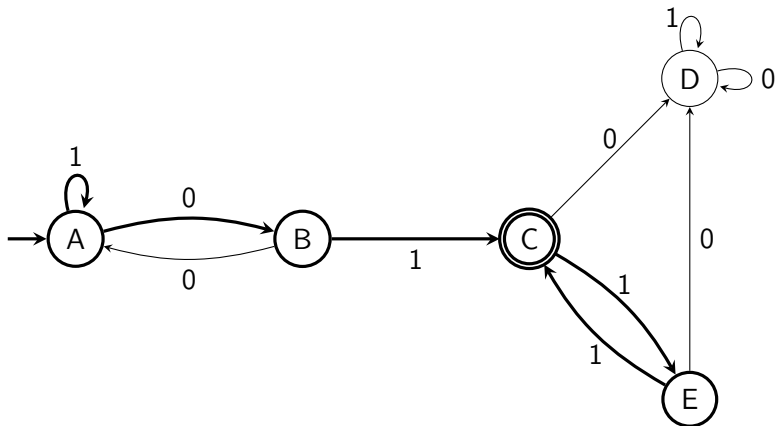
## В предыдущей серии: распознавание слова ДКА



## В предыдущей серии: распознавание слова ДКА

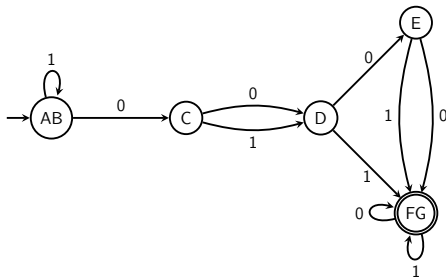
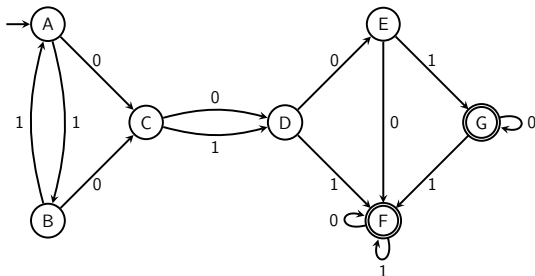


## В предыдущей серии: распознавание слова ДКА



Слово распознается за  $O(n)$

## В предыдущей серии: минимизация



**Недетерминированный конечный автомат** —  $\langle Q, \Sigma, \delta, q_0, F \rangle$

- $Q \neq \emptyset$  — конечное множество состояний
- $\Sigma$  — Конечный входной алфавит
- $\delta$  — отображение типа  $Q \times \Sigma \rightarrow 2^Q$ 
  - ▶  $\delta(q_i, x) = \{q_{j_0}, \dots, q_{j_k}\}$
- $q_0 \in Q$  — начальное состояние
- $F \subseteq Q$  — множество конечных состояний



# Недетерминированный КА: пример

$$\delta(q_0, a) = q_0$$

...

$$\delta(q_0, \kappa) = q_0$$

...

$$\delta(q_0, \text{я}) = q_0$$

$$\delta(q_0, \kappa) = q_1$$

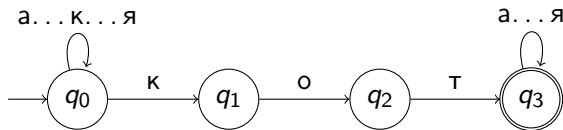
$$\delta(q_1, o) = q_2$$

$$\delta(q_2, \tau) = q_3$$

$$\delta(q_3, a) = q_3$$

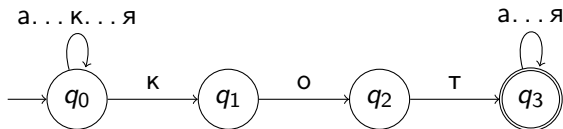
...

$$\delta(q_3, \text{я}) = q_3$$



- **Конфигурация (Мгновенное описание) КА** —  $\langle q, \omega \rangle$ , где  $q \in Q, \omega \in \Sigma^*$
- **Такт работы** — бинарное отношение  $\vdash$ : если  $q \in \delta(p, x)$  и  $\omega \in \Sigma^*$ , то  $\langle p, x\omega \rangle \vdash \langle q, \omega \rangle$
- Бинарное отношение  $\vdash^*$  — рефлексивное, транзитивное замыкание  $\vdash$
- **НКА допускает слово  $\alpha$** , если  $\exists t \in F : \langle s, \alpha \rangle \vdash^* \langle t, \varepsilon \rangle$
- **Язык НКА**  $L(A) = \{\omega \in \Sigma^* \mid \exists t \in F : \langle s, \omega \rangle \vdash^* \langle t, \varepsilon \rangle\}$
- **ДКА** — частный случай НКА

# Недетерминированный КА: пример



{кот, скот, котлета, мякоть, антрекот...}

## Алгоритм, определяющий допустимость слова

$$R(\alpha) = \{p \mid \langle q_0, \alpha \rangle \vdash^* \langle p, \varepsilon \rangle\}$$

$$R(\varepsilon) = \{q_0\}$$

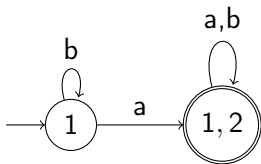
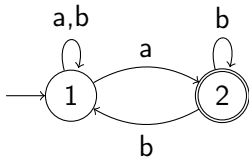
$$R(\alpha c) = \{q \mid q \in \delta(p, c), p \in R(\alpha)\}$$

НКА допускает слово  $\alpha \Leftrightarrow \exists t \in F : t \in R(\alpha)$

# Построение ДКА по НКА

- Помещаем в *Queue* множество  $\{q_0\}$
- Пока очередь не пуста, выполняем:
  - ▶  $q = \text{Queue.pop}()$
  - ▶ Строим множество  $q' = \{t = \delta(s, c) \mid s \in q, c \in \Sigma\}$ . Если  $q' \notin \text{Queue}$ , добавить его в очередь. Каждое такое множество — новая вершина ДКА; добавляем переходы по соответствующим символам
  - ▶ Если во множестве есть хотя бы одна вершина, являющаяся терминальной в данном НКА, то соответствующая вершина ДКА будет конечной
- Результат:  $\langle \Sigma, Q_d, q_{d_0} \in Q_d, F_d \subset Q_d, \delta_d : Q_d \times \Sigma \rightarrow Q_d \rangle$ 
  - ▶  $Q_d = \{q_d \mid q_d \subset 2^Q\}$
  - ▶  $q_{d_0} = \{q_0\}$
  - ▶  $F_d = \{q \in Q_d \mid \exists p \in F : p \in q\}$
  - ▶  $\delta_d(q, c) = \{\delta(a, c) \mid a \in q\}$

## Детерминизация НКА: пример



# Эквивалентность языков, распознаваемых ДКА и НКА

## Теорема

*ДКА и НКА распознают один и тот же класс языков*

## Доказательство.

$\Rightarrow$ : очевидно

$\Leftarrow$ : Рассмотрим произвольный НКА и покажем, что алгоритм строит по нему эквивалентный ДКА.

$\forall q \in q_d, \forall c \in \Sigma, \forall p \in \delta(q, c) : p \in \delta_d(q_d, c)$

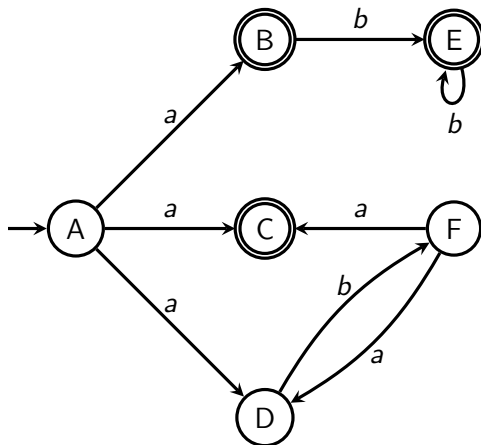
Рассмотрим  $\langle q_0, w_1 w_2 \dots w_m \rangle \vdash \langle u_1, w_2 \dots w_m \rangle \vdash^* \langle u_m, \varepsilon \rangle, u_m \in F$

$\forall i : u_i \in u_{d_i}$ , где  $(q_{d_0}, w_1 w_2 \dots w_m) \vdash (u_{d_1}, w_2 \dots w_m) \vdash^* (u_{d_m}, \varepsilon)$

$\Rightarrow u_m \in u_{d_m}$



# Распознавание слова НКА



Слово распознается за  $O(|\omega| \sum_{t \in Q} \sum_{c \in \Sigma} |\delta(t, c)|)$



# Произведение автоматов

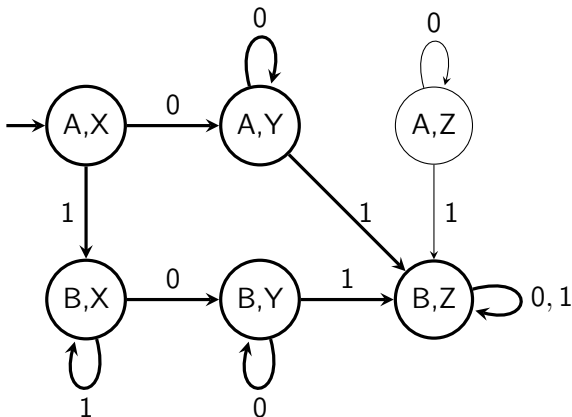
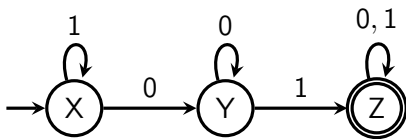
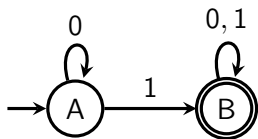
$A_1 = \langle \Sigma_1, Q_1, q_{10}, \delta_1, F_1 \rangle$  и  $A_2 = \langle \Sigma_2, Q_2, q_{20}, \delta_2, F_2 \rangle$  — КА

Произведением автоматов назовем  $A = \langle \Sigma, Q, q_0, \delta, F \rangle$ , где

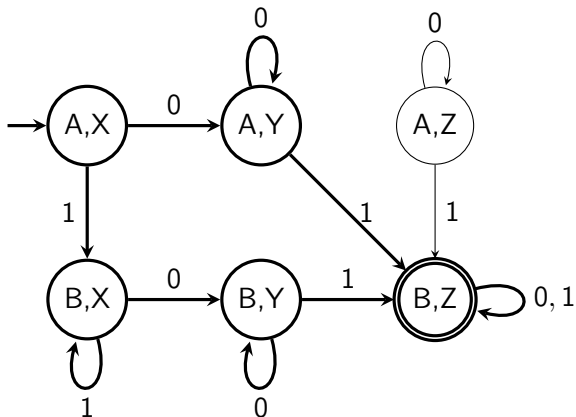
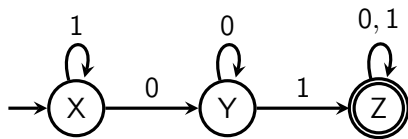
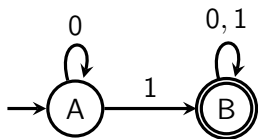
- $\Sigma = \Sigma_1 \cup \Sigma_2$
- $Q = Q_1 \times Q_2$
- $q_0 = (q_{10}, q_{20})$
- $F \subseteq Q$ 
  - ▶  $F = F_1 \times F_2$  — распознает **пересечение** языков
  - ▶  $F = (F_1 \times Q_2) \cup (Q_1 \times F_2)$  — распознает **объединение** языков
  - ▶  $F = F_1 \times (Q_2 \setminus F_2)$  — распознает **разность** языков
- $\delta((q_1, q_2), c) = (\delta_1(q_1, c), \delta_2(q_2, c))$

Интуиция: ищем пути в двух автоматах одновременно

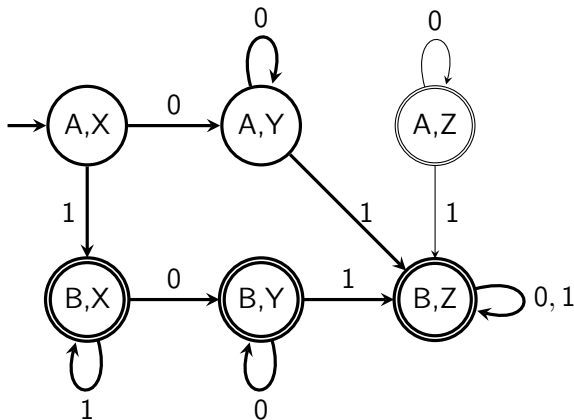
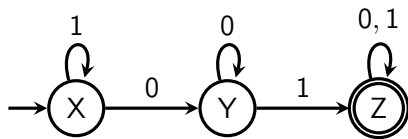
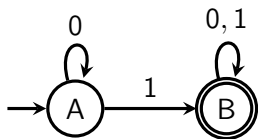
## Произведение автоматов: пример



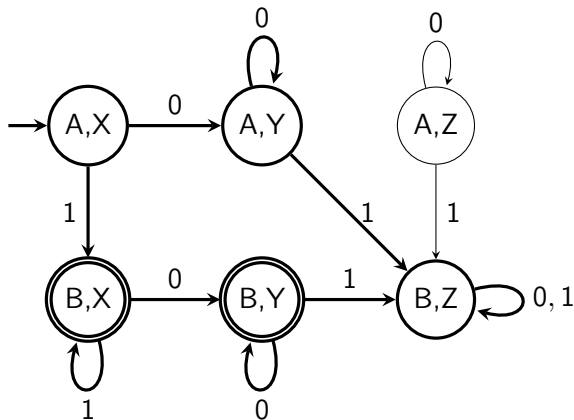
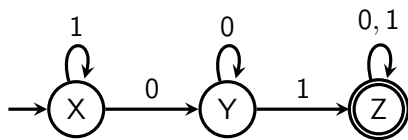
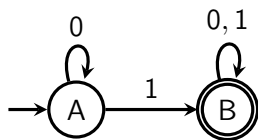
# Пересечение языков



# Объединение языков



# Разность языков



# Замкнутость автоматных языков относительно операций

Автоматные языки замкнуты относительно операций:

- Объединения
- Пересечения
- Разности
- Дополнения
  - ▶  $\bar{X} = \Sigma^* \setminus X$

# Регулярное множество (регулярный язык)

**Регулярное множество** в алфавите  $\Sigma$  определяется итеративно:

- $\emptyset$  — регулярное множество в алфавите  $\Sigma$
- $\{a\}$  — регулярное множество в алфавите  $\Sigma$  для каждого  $a \in \Sigma$
- $\{\varepsilon\}$  — регулярное множество в алфавите  $\Sigma$
- Если  $P$  и  $Q$  — регулярные множества в алфавите  $\Sigma$ , то регулярны
  - ▶  $P \cup Q$  (объединение)
  - ▶  $PQ = \{pq \mid p \in P, q \in Q\}$  (конкатенация)
  - ▶  $P^* = \{\varepsilon\} \cup P \cup PP \cup PPP \cup \dots$  (итерация)
- Ничто другое не является регулярным множеством в алфавите  $\Sigma$
- Множество всех регулярных языков обозначим  $\mathbb{R}$

# Примеры регулярных языков

- Все конечные языки
  - ▶  $\{-2147483648, -2147483647, \dots, 2147483647\}$  — все 32-разрядные целые числа
- $L_a = \{a^k \mid k - odd\}$
- $L_b = \{b^l \mid l - even\}$
- $L_{ab} = \{a^k b^l \mid k - odd, l - even\} = L_a L_b$
- $L = \{a^*\} = L_a^*$



# Регулярное выражение

**Регулярное выражение** — способ записи регулярного множества

- $\emptyset$  — обозначает  $\emptyset$
- $a$  — обозначает  $\{a\}$
- $\varepsilon$  — обозначает  $\{\varepsilon\}$
- Если  $p$  и  $q$  обозначают  $P$  и  $Q$ , то:
  - ▶  $p \mid q$  обозначает  $P \cup Q$
  - ▶  $pq$  обозначает  $PQ$
  - ▶  $p^*$  обозначает  $P^*$

# Примеры регулярных выражений

- $-2147483648 \mid -2147483647 \mid \dots \mid 2147483647$  — все 32-разрядные целые числа
- $a(aa)^* : L_a = \{a^k \mid k - odd\}$
- $(bb)^* : L_b = \{b^l \mid l - even\}$
- $a(aa)^*(bb)^* : L_{ab} = \{a^k b^l \mid k - odd, l - even\} = L_a L_b$
- $a^* : L = \{a^*\} = L_a^*$

# Замкнутость регулярных языков относительно операций

Регулярные языки замкнуты ( $A \in \mathbb{R}, B \in \mathbb{R} \Rightarrow A \diamond B \in \mathbb{R}$ ) относительно операций:

- Конкатенации ( $L_1 L_2$ ), объединения ( $L_1 \cup L_2$ ), итерации ( $L^*$ )
- Пересечения ( $L_1 \cap L_2$ ), дополнения ( $\neg L$ ), разности ( $L_1 \setminus L_2$ )
- Обращения ( $L_{rev} = \{\omega^R = a_m a_{m-1} \dots a_1 \mid a_1 a_2 \dots a_m = \omega \in L\}$ )
- Гомоморфизма цепочек  $\phi$ 
  - ▶  $\phi(\varepsilon) = \varepsilon$
  - ▶  $\phi(\alpha\beta) = \phi(\alpha)\phi(\beta)$
- Обратного гомоморфизма цепочек

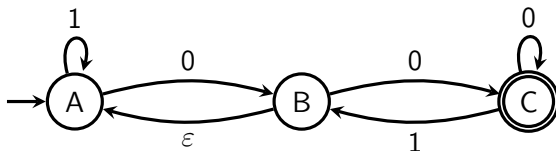
# Теорема Клини

## Теорема

*Классы автоматных и регулярных языков эквивалентны*

# НКА с $\varepsilon$ -переходами: почему бы и нет?

$$\delta : Q \times (\Sigma \cup \varepsilon) \rightarrow 2^Q$$

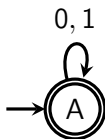
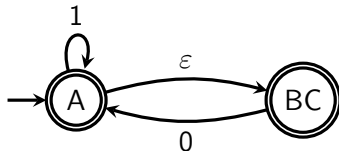
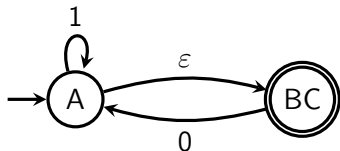
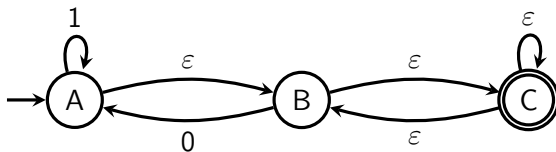


Ничего не поломалось?

# Эквивалентность НКА с $\varepsilon$ -переходами и НКА без $\varepsilon$ -переходов

- НКА без  $\varepsilon$ -переходов — частный случай НКА с  $\varepsilon$ -переходами
- В обратную сторону — можно построить  $\varepsilon$ -замыкание
  - ▶ Транзитивное замыкание: для каждого подграфа, состоящего только из  $\varepsilon$ -переходов, делаем  $\varepsilon$ -замыкание
  - ▶ Добавление терминальных состояний: для  $\varepsilon$ -перехода из состояния  $u$  в  $v$ , где  $v$  — терминальное, добавляем  $u$  в терминальные
  - ▶ Добавление ребер:  $\forall u, v, c, w : \delta(u, \varepsilon) = v, \delta(v, c) = w$ , добавим переход  $\delta(u, c) = w$
  - ▶ Устранение  $\varepsilon$ -переходов

## $\epsilon$ -замыкание



# Теорема Клини: доказательство $\Leftarrow$

## Теорема

*Классы автоматных и регулярных языков эквивалентны*

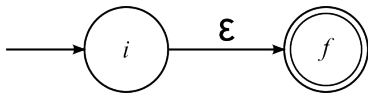
## Доказательство.

$\Leftarrow$ : Построим по регулярному выражению КА (НКА с  $\varepsilon$ -переходами)

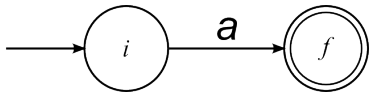




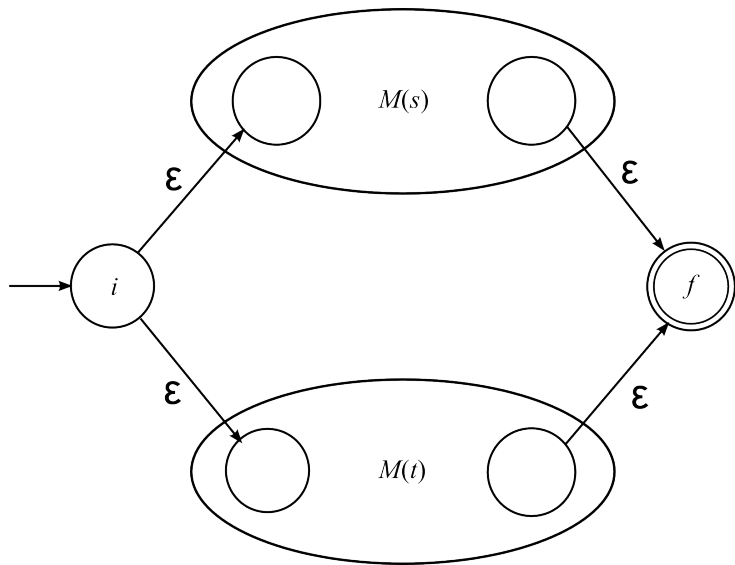
## Построение КА по РВ: $\varepsilon$



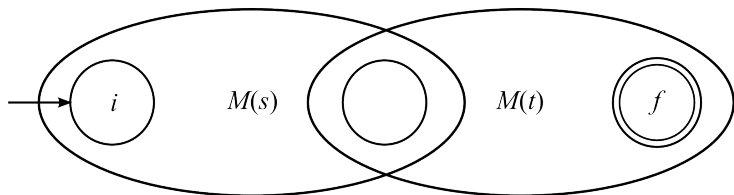
## Построение КА по РВ: символ



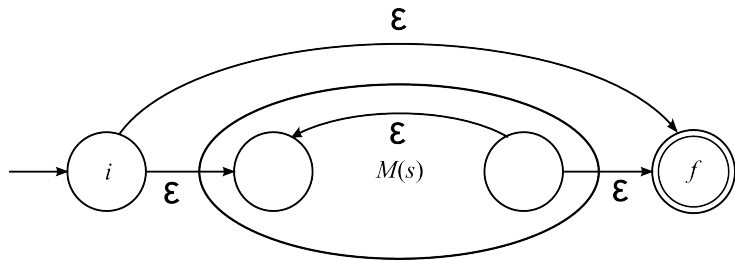
## Построение КА по РВ: объединение $p \mid q$



## Построение КА по РВ: конкатенация $pq$



## Построение КА по РВ: итерация $p^*$



# Теорема Клини: доказательство $\Rightarrow$

## Теорема

*Классы автоматных и регулярных языков эквивалентны*

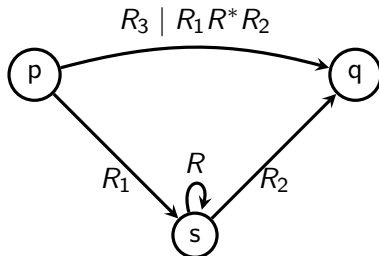
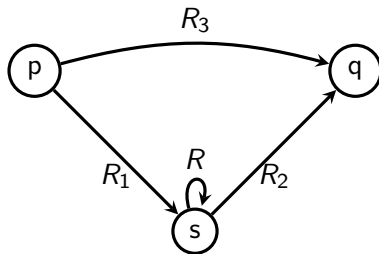
## Доказательство.

$\Rightarrow$ : Построим регулярное выражение по конечному автомату методом исключения состояний

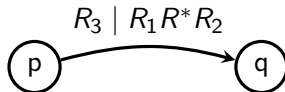
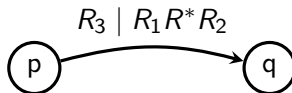
Идея: на ребрах пишем регулярные выражения, соответствующие путям между вершинами, последовательно исключаем состояния



## Исключение состояния $s$

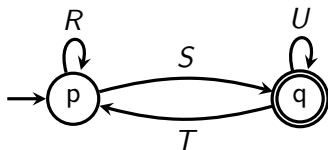


## Исключение состояния $s$ : удаление ребер и вершины



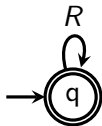


## Исключение состояний: последний шаг



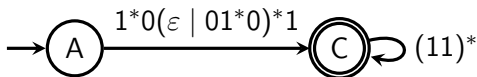
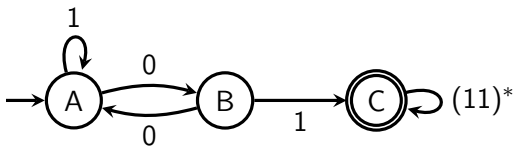
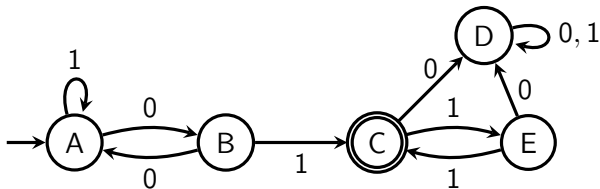
$$(R^* \mid SU^*T)^*SU^*$$

## Исключение состояний: последний шаг



$R^*$

## Исключение состояний: пример



$$1^*0(\epsilon \mid 01^*0)^*1(11)^*$$

# Свойства регулярных выражений

- $a \mid a = a$
- $a \mid \emptyset = a$
- $a \mid b = b \mid a$
- $a \mid (b \mid c) = (a \mid b) \mid c$
- $a(bc) = (ab)c$
- $\{\varepsilon\}a = a\{\varepsilon\} = a$
- $\emptyset a = a\emptyset = \emptyset$
- $a(b \mid c) = ab \mid ac$
- $(a \mid b)c = ac \mid bc$
- $\{\varepsilon\} \mid aa^* \subseteq a^*$
- $\{\varepsilon\} \mid a^*a \subseteq a^*$
- $ab \subseteq b \Rightarrow a^*b \subseteq b$
- $ab \subseteq a \Rightarrow ab^* \subseteq a$

# Регулярная грамматика

**Праволинейная грамматика** — грамматика, все правила которой имеют следующий вид:

- $A \rightarrow aB, A \rightarrow a$  или  $A \rightarrow \varepsilon$ , где  $A, B \in V_N, a \in V_T$

**Левوليнейная грамматика** — грамматика, все правила которой имеют следующий вид:

- $A \rightarrow Ba, A \rightarrow a$  или  $A \rightarrow \varepsilon$ , где  $A, B \in V_N, a \in V_T$

# Регулярная грамматика

**Праволинейная грамматика** — грамматика, все правила которой имеют следующий вид:

- $A \rightarrow aB$ ,  $A \rightarrow a$  или  $A \rightarrow \varepsilon$ , где  $A, B \in V_N$ ,  $a \in V_T$

**Левوليнейная грамматика** — грамматика, все правила которой имеют следующий вид:

- $A \rightarrow Ba$ ,  $A \rightarrow a$  или  $A \rightarrow \varepsilon$ , где  $A, B \in V_N$ ,  $a \in V_T$

## Теорема

Пусть  $L$  — формальный язык.

$\exists G_r$  — праволинейная грамматика, т.ч.  $L = L(G_r) \Leftrightarrow$

$\exists G_l$  — левوليнейная грамматика, т.ч.  $L = L(G_l)$

# Регулярная грамматика

**Праволинейная грамматика** — грамматика, все правила которой имеют следующий вид:

- $A \rightarrow aB$ ,  $A \rightarrow a$  или  $A \rightarrow \varepsilon$ , где  $A, B \in V_N$ ,  $a \in V_T$

**Левوليнейная грамматика** — грамматика, все правила которой имеют следующий вид:

- $A \rightarrow Ba$ ,  $A \rightarrow a$  или  $A \rightarrow \varepsilon$ , где  $A, B \in V_N$ ,  $a \in V_T$

## Теорема

Пусть  $L$  — формальный язык.

$\exists G_r$  — праволинейная грамматика, т.ч.  $L = L(G_r) \Leftrightarrow$

$\exists G_l$  — левوليнейная грамматика, т.ч.  $L = L(G_l)$

**Регулярная грамматика** — праволинейная или левوليнейная грамматика

# Эквивалентность регулярной грамматики и НКА

Алгоритм построения НКА  $\langle Q, \Sigma, q_0, \delta, F \rangle$  по праволинейной грамматике  $\langle V_T, V_N, P, S \rangle$

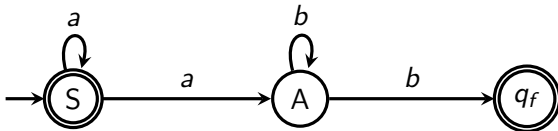
- $Q = V_N \cup \{q_f\}$
- $\forall (A \rightarrow aB) \in P : \delta(A, a) = B$
- $\forall (A \rightarrow a) \in P : \delta(A, a) = q_f$
- $q_0 = S$
- $\forall (B \rightarrow \varepsilon) \in P : B \in F$



## Пример построения НКА по регулярной грамматике

$$S \rightarrow aA \mid aS \mid \varepsilon$$

$$A \rightarrow bA \mid b$$

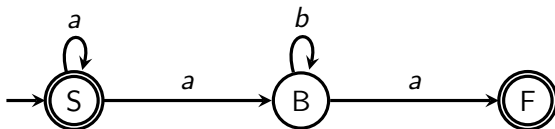


# Эквивалентность регулярной грамматики и НКА

Алгоритм построения праволинейной грамматики  $\langle V_T, V_N, P, S \rangle$  по НКА  $\langle Q, \Sigma, q_0, \delta, F \rangle$

- $V_N = Q$
- $V_T = \Sigma$
- $\forall \delta(A, a) = B : (A \rightarrow aB) \in P$
- $\forall B \in F : (B \rightarrow \varepsilon) \in P$
- $S = q_0$
- Опционально: удалить  $\varepsilon$ -правила и бесполезные символы

## Пример построения регулярной грамматики по НКА

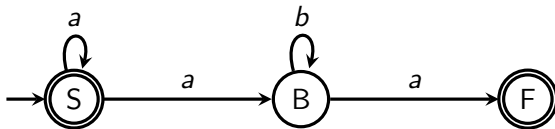


$$S \rightarrow aB \mid aS \mid \varepsilon$$

$$B \rightarrow bB \mid aF$$

$$F \rightarrow \varepsilon$$

## Пример построения регулярной грамматики по НКА



$$S \rightarrow aB \mid aS \mid \varepsilon$$

$$B \rightarrow bB \mid a$$

# Лемма о разрастании (о накачке)

## Теорема

$L$  — регулярный язык над  $\Sigma \Rightarrow \exists n : \forall \omega \in L, |\omega| > n$

$\exists x, y, z \in \Sigma^* : xyz = \omega, y \neq \varepsilon, |xy| \leq n,$

$\forall k \geq 0 : xy^kz \in L$

## Доказательство.

Строим автомат, распознающий  $L$ .

Обозначаем за  $n$  число состояний автомата.

Слово длины большей, чем  $n$ , обязано при разборе пройти через одно состояние дважды — получили цикл.

Метка цикла — искомое  $y$ , по циклу можно пройти сколько угодно раз.



## Использование леммы о накачке

$$L = \{({}^k)^k \mid k \geq 0\}$$

- Предполагаем, что  $L$  — регулярный язык, значит выполняется лемма о накачке
- Берем  $n$  из леммы, рассматриваем слово  $({}^n)^n$
- Его можно разбить на  $xuz$ ,  $u \neq \varepsilon$ ,  $|xu| \leq n$
- $|xu| \leq n \Rightarrow u = ({}^b, b > 0$
- Берем  $k = 2$  :  $xu^kz = ({}^{n+b})^n$ , что не принадлежит  $L$
- Получили противоречие  $\Rightarrow L$  не регулярен

# Использование леммы о накачке

$$L = \{a^{k^2} \mid k \geq 0\}$$

- Предполагаем, что  $L$  — регулярный язык, значит выполняется лемма о накачке
- Берем  $n$  из леммы, рассматриваем слово  $a^{(n+1)^2}$ 
  - ▶ Слово  $a^{n^2}$  не подойдет, потому что  $|a^{n^2}| = n$ , где  $n = 1$
- Его можно разбить на  $x y z$ ,  $y \neq \varepsilon$ ,  $|xy| \leq n$
- Берем  $k = 0$  :  $x y^0 z = xz$
- $n^2 < n^2 + n + 1 = (n+1)^2 - n \leq |xz| \leq (n+1)^2 - 1 < (n+1)^2$
- Длина слова  $xz$  находится между двумя соседними полными квадратами, поэтому это слово не в  $L$
- Получили противоречие  $\Rightarrow L$  не регулярен

- ДКА, НКА, НКА с  $\varepsilon$ -переходами, регулярные выражения, регулярные грамматики — все эти формализмы задают один класс (регулярных) языков и эквивалентны друг другу
- Проверка принадлежности слова регулярному языку осуществляется за  $O(n)$  и не требует дополнительной памяти
- Класс регулярных языков обладает хорошими свойствами, прост и нагляден
- С помощью леммы о накачке можно доказать нерегулярность языка