

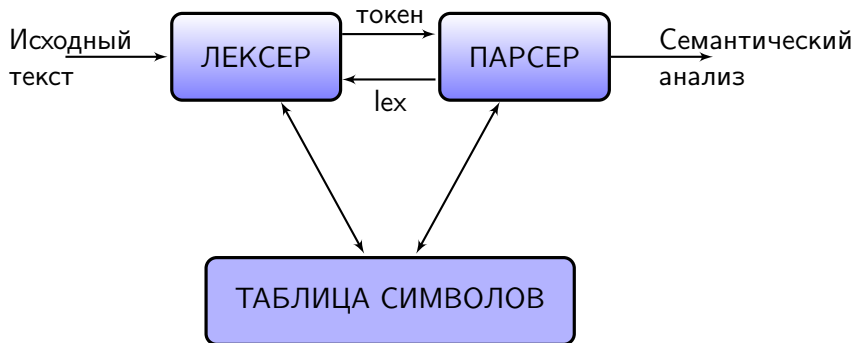
Теория и практика компиляции программ: лексический анализ

Юхин Кирилл

iLab, Intel Corp

МФТИ, 26 сентября 2018 года

Положение лексического анализа



Токен, лексема и шаблон

Токен – пара: имя и, опционально, прикрепленное значение.
Имя – это элемент некоторого конечного множества.

Лексема – последовательность знаков, соответствующая определённому токenu.

Шаблон – описание всех лексем, соответствующих токенам с одним и тем же именем.

Примеры

Имя токена	Шаблон	Лексемы
BOOLEAN	"#t" или "#f"	#t, #f
LPAREN	"(" или "["	(, [
SYMBOL	набор букв и цифр, содержащий букву	alala, 3k34dk, if, let
NUMBER	непустой набор цифр и точки	32, 0.0, .121, 3.14
STRING	всё, кроме " окружённое "	"мечты" "error"
SPACE	непустая последовательность пробелов	

Атрибуты токенов

```
(format #t "~a" 18.3)
```

- (LPAREN)
- (SYMBOL, индекс "format" в таблице символов)
- (SPACE)
- (BOOLEAN, true)
- (SPACE)
- (STRING, индекс "~a" в таблице строк)
- (SPACE)
- (NUMBER, 18.3)
- (RPAREN)

Лексические ошибки

- Некоторые ошибки лексер не способен обнаружить, например:
 - ▶ `(let ((a 3))) a)` → лишняя скобка
 - ▶ `(/ 1 "2")` → неправильный тип
- Однако некоторые, может:
 - ▶ `"abcd < EOF >` → незаконченная строка
 - ▶ `13a` → недопустимая конструкция в Си
- Такие ошибки как правило сигнализируются когда ни один шаблон не может отыскать лексему в тексте

Восстановление после ошибок

Паника Игнорирование всех последующих знаков. До конца файла, либо до некоторого знака (например пробела)

Догадка

- Удаление одного знака из входного потока,
- Вставка недостающего знака во входной поток
- Замена одного знака другим
- Перестановка двух соседних знаков

и повторная попытка найти лексему.

Общий механизм анализа

- Лексер считывает и запоминает по одному знаку.
- Когда информации становится достаточно, он передаёт парсеру очередной токен.
- Иногда лексеру требуется знать последующие символы, чтобы определить текущий токен. Например:
 - ▶ В Fortran: $DO\ 5\ I = 1.25 \leftrightarrow DO\ 5\ I = 1,25$
 - ▶ В Си: -, <, =

Для этого применяется двухбуферная схема.

Определения

Алфавит – конечное множество знаков. Обозначим Σ .

Язык – набор строк знаков из алфавита Σ .

Грамматика – набор правил, выделяющий из всех последовательностей алфавитных знаков, строки принадлежащие языку.

Функция значения $L(e) = M$ – язык M соответствующий грамматике e .

Грамматика задаёт только один язык. Но язык может быть определён несколькими грамматиками. $L(e)$ - сюръекция.

Регулярные грамматики

- Задают языки, включающие:
 - ▶ Пустую строку ε : $\{''\}$ или \emptyset
 - ▶ Одиночные символы из алфавита Σ : $\{''a''\}, \{''b''\}$
 - ▶ Объединение: $A + B = \{a|a \in A\} \cup \{b|b \in B\}$
 - ▶ Конкатенация: $AB = \{ab|a \in A \wedge b \in B\}$
 - ▶ Повторение: $A^* = \bigcup_{i \geq 0} A^i, \{A^0 = \varepsilon, A^i = A \dots A \text{ } i \text{ раз}\}$
- Контекстно-свободные
 - ▶ C++ шаблон: $A \langle B \langle C \rangle \rangle a$;
 - ▶ C++ сдвиг: $C \rangle \rangle a$;
- Задаются регулярными выражениями
- Задаются конечными автоматами

Регулярные выражения

- База:

- ▶ Пустая строка ε : $L(\varepsilon) = \{\emptyset\}$
- ▶ Один символ $a \in \Sigma$: $L(a) = \{a\}$

- Композиция:

- ▶ Объединение: $L(A + B) = L(A) \cup L(B)$
- ▶ Конкатенация: $L(AB) = \{ab | a \in L(A) \wedge b \in L(B)\}$
- ▶ Повторение: $L(A^*) = \bigcup_{i \geq 0} A^i$

- Расширения:

- ▶ $A^+ = AA^*$
- ▶ $A? = A + \varepsilon$
- ▶ $A^k = \underbrace{AA \dots A}_{k \text{ раз}}$

POSIX нотация

- $"" \equiv \varepsilon$
- $"abc" \equiv \{ "a" \} | \{ "b" \} | \{ "c" \}$
- $(A) \equiv A$
- $A|B \equiv A + B$
- AB
- $\cdot \equiv \Sigma \setminus \{ \text{"перевод строки"} \}$
- $[abc] \equiv "a" | "b" | "c"$
- $[j - p] \equiv \{ \text{все буквы коды которых лежат между } j \text{ и } p \text{ включительно} \}$
- $[^{\wedge}...] \equiv \Sigma \text{ кроме } [...]$
- $A\{n\} \equiv A^n$
- $A\{n, m\} \equiv A^n + A^{n+1} + \dots + A^m$
- $A\{, m\} \equiv A\{0, m\}; A\{n, \} \equiv A^n A^*$

Примеры

- Цифра: '0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9' $\equiv [0 - 9]$
- Целое неотрицательное число(Z_0^+):
- Целое число(Z):
- Десятичное дробное число:
- Идентификатор в Си:

Примеры

- Цифра: '0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9' $\equiv [0 - 9]$
- Целое неотрицательное число(Z_0^+): $[0 - 9]^+$
- Целое число(Z):
- Десятичное дробное число:
- Идентификатор в Си:

Примеры

- Цифра: $'0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9' \equiv [0 - 9]$
- Целое неотрицательное число(Z_0^+): $[0 - 9]^+$
- Целое число(Z): $'-'?[0 - 9]^+$
- Десятичное дробное число:
- Идентификатор в Си:

Примеры

- Цифра: $'0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9' \equiv [0 - 9]$
- Целое неотрицательное число(Z_0^+): $[0 - 9]^+$
- Целое число(Z): $'-'?[0 - 9]^+$
- Десятичное дробное число: $'-'?[0 - 9]^*.'[0 - 9]^+$
- Идентификатор в Си:

Примеры

- Цифра: $'0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9' \equiv [0 - 9]$
- Целое неотрицательное число(Z_0^+): $[0 - 9]^+$
- Целое число(Z): $'-'?[0 - 9]^+$
- Десятичное дробное число: $'-'?[0 - 9]^*.'[0 - 9]^+$
- Идентификатор в Си: $[a - zA - Z_][0 - 9a - zA - Z_]^*$

А имплементация?

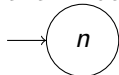
Для строки s и выражения e , $s \in L(e)$?

Конечный автомат

- Входной алфавит Σ
- Конечный набор состояний S



- Начальное состояние $n \in S$



- Набор конечных состояний $F \subseteq S$



- Набор переходов из состояния $s_1 \in S$ в состояние $s_2 \in S$ по входу $a \in \Sigma$: $s_1 \xrightarrow{a} s_2$

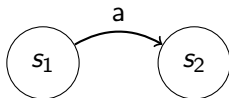
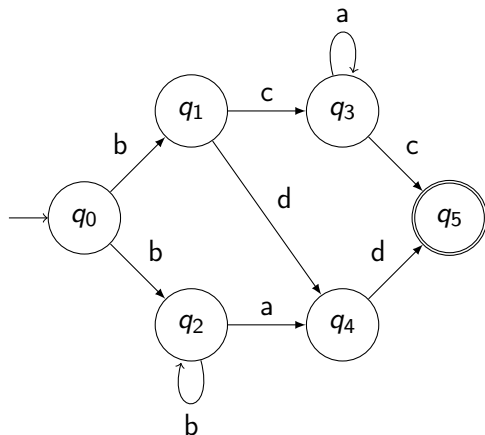
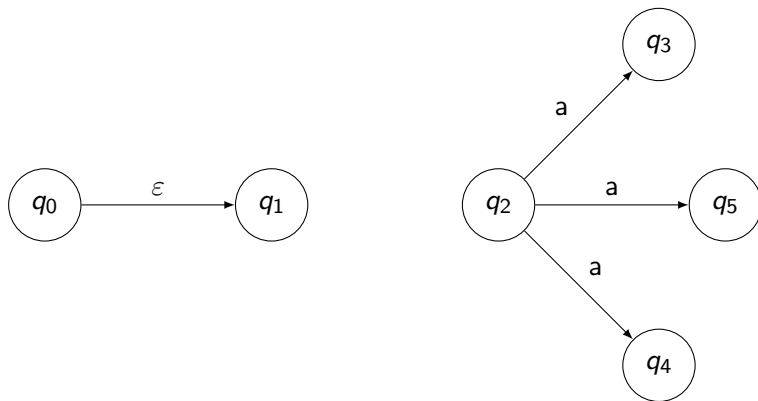


Таблица и диаграмма переходов



	a	b	c	d
q_0	\emptyset	q_1, q_2	\emptyset	\emptyset
q_1	\emptyset	\emptyset	q_3	q_4
q_2	q_4	q_2	\emptyset	\emptyset
q_3	q_3	\emptyset	q_5	\emptyset
q_4	\emptyset	\emptyset	\emptyset	q_5
q_5	\emptyset	\emptyset	\emptyset	\emptyset

множественные и ε - переходы



ДКА и НКА

Не детерминированный конечный автомат (НКА)

Конечный автомат, который может одновременно быть в нескольких состояниях.

Может содержать:

- ε -переходы
- множественные переходы

НКА как правило лаконичнее.

Детерминированный конечный автомат (ДКА)

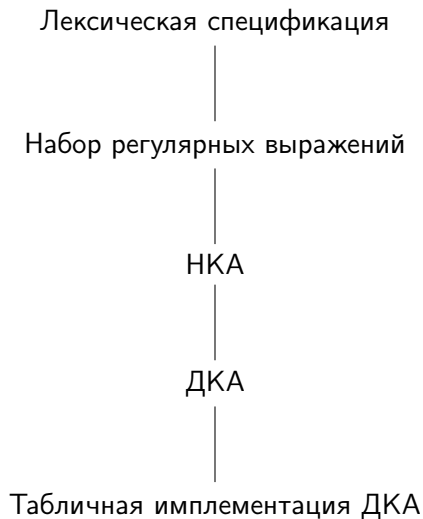
Конечный автомат, находящийся в каждый момент времени только в одном состоянии.

Не содержит:

- ε -переходы
- множественные переходы

НКА может быть преобразован в ДКА принимающий тот же язык

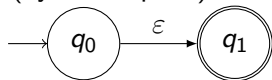
Общая схема генерации сканера



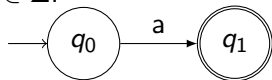
Построение НКА для регулярного выражения

Автоматы для базовых элементов:

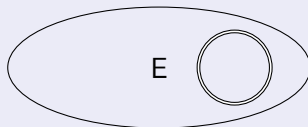
- ε (пустая строка):



- $a \in \Sigma$:



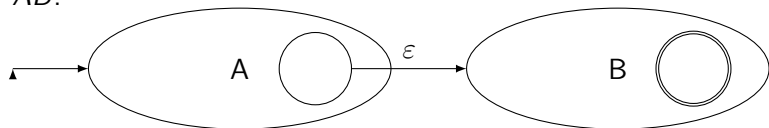
Общее обозначение НКА для выражения E



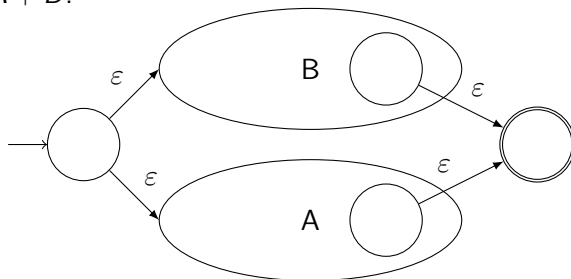
Построение НКА для регулярного выражения

Комбинации:

- AB :

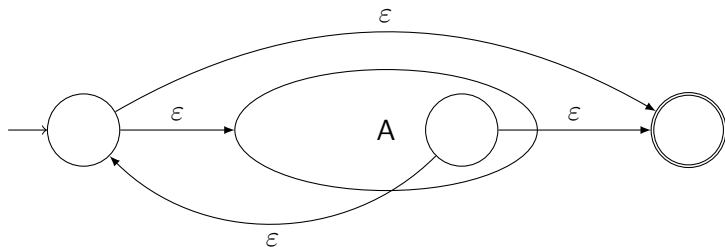


- $A + B$:



Построение НКА для регулярного выражения

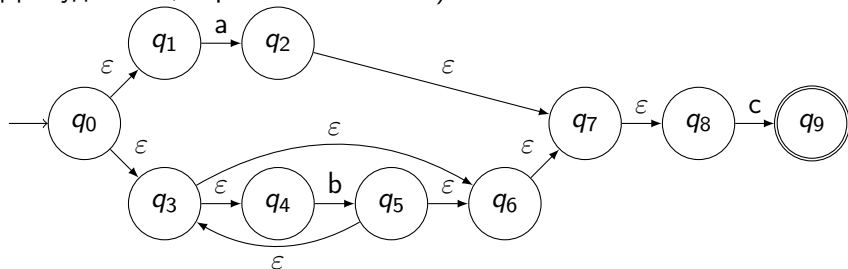
- A^* :



Пример

Выражение: $(a + b^*)c$

(Для удобства, вершины помечены)



Преобразование НКА в ДКА

ϵ -замыкание

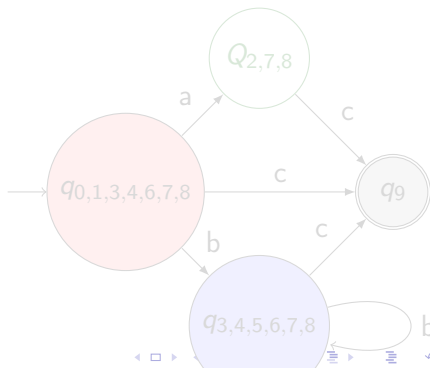
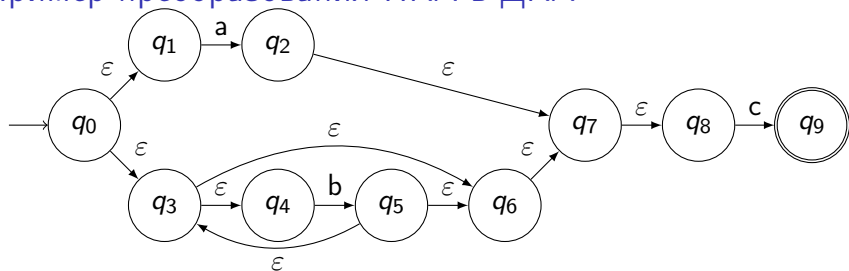
ϵ -замыкание — Множество состояний НКА, в которые можно перейти из данного по цепочке ϵ -переходов.

$$\epsilon closure(q_i) = \{q_i\} \cup \{\epsilon closure(q_j) | \exists q_j \xrightarrow{\epsilon} q_i\}$$

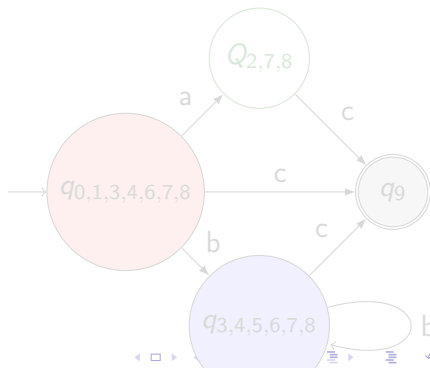
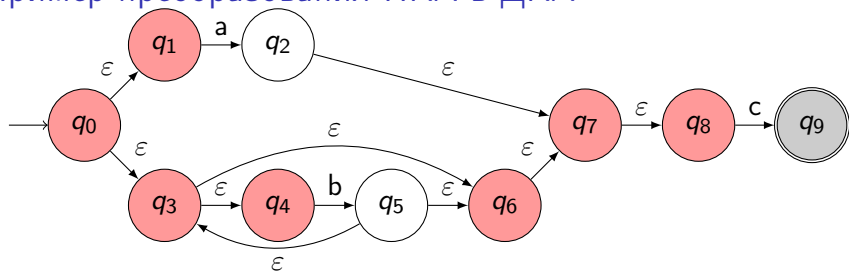
Пусть S — множество состояний НКА, s — его начальное состояние, F — все конечные.

- Каждое состояние ДКА (X) — это подмножество (S) состояний НКА. $X \subseteq S$
- Начальное состояние: $\epsilon closure(s)$
- Конечные состояния: $\{X | X \cap F \neq \emptyset\}$
- Переходы: $\exists X \xrightarrow{a} Y \Leftrightarrow Y = \bigcup_{q_i} \epsilon closure(q_i), q_i \in X$

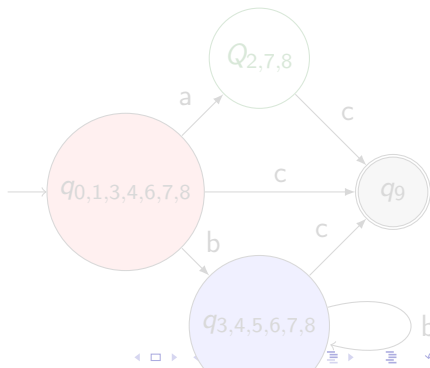
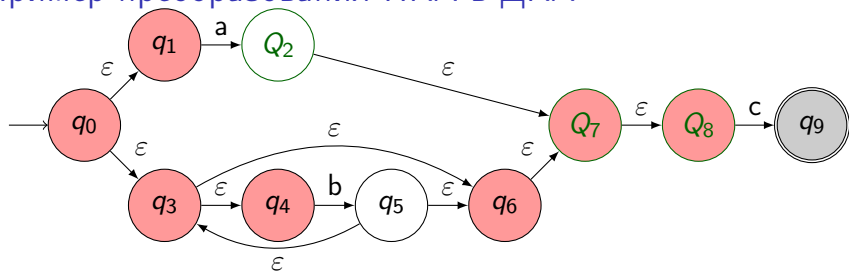
Пример преобразования НКА в ДКА



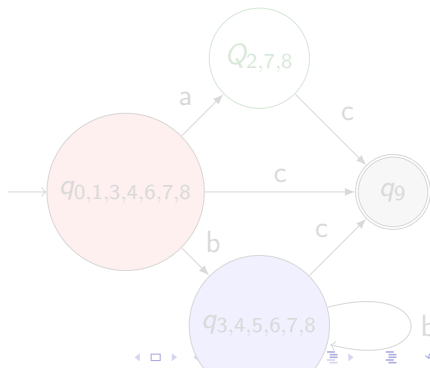
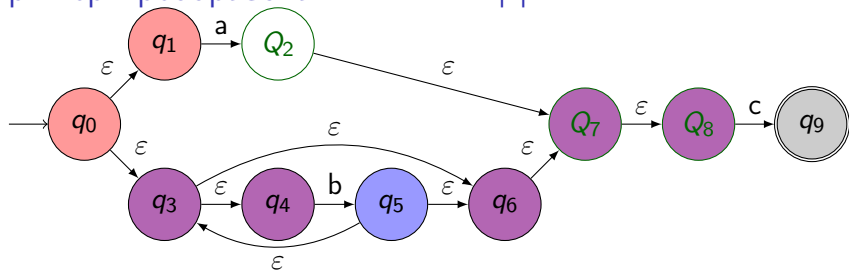
Пример преобразования НКА в ДКА



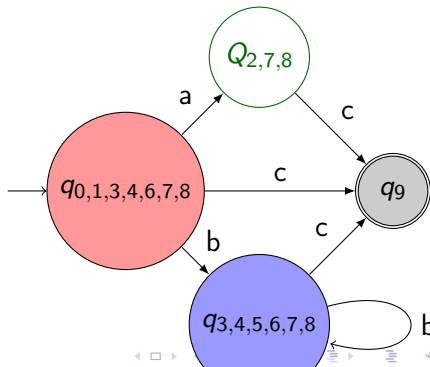
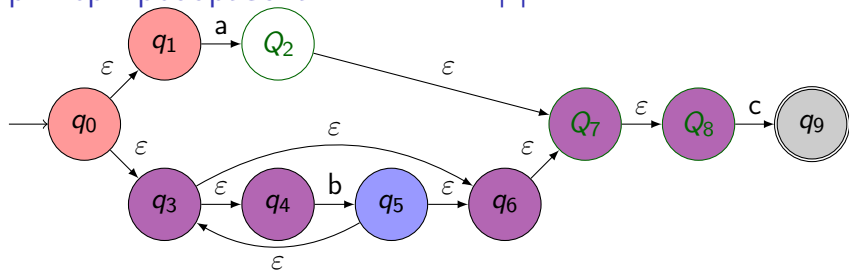
Пример преобразования НКА в ДКА



Пример преобразования НКА в ДКА



Пример преобразования НКА в ДКА



Табличная имплементация

A	a	b	c
q_1	q_3	q_2	q_4
q_2	\emptyset	q_2	q_4
q_3	\emptyset	\emptyset	q_4
q_4	\emptyset	\emptyset	\emptyset

```
(define (go state input)
  (let ((next (gethash '(state input))))
    (if next
        (go next (move (forward input)))
        #f)))

(define (accept input)
  (let ((rez (go n input)))
    (and rez (accepting? rez))))
```

