

# Звіт щодо виконання лабораторної роботи по реалізації кінцевого автомата для обробки регулярних виразів

## Зміст

<b>1</b>	<b>Вступ</b>	<b>2</b>
<b>2</b>	<b>Загальна архітектура</b>	<b>2</b>
<b>3</b>	<b>Базовий клас State</b>	<b>2</b>
3.1	Опис та призначення	2
3.2	Обґрунтування реалізації	3
<b>4</b>	<b>Типи станів</b>	<b>3</b>
4.1	StartState	3
4.2	TerminationState	3
4.3	DotState	3
4.4	AsciiState	4
4.5	StarState	4
4.6	PlusState	4
4.7	CharacterClass	5
<b>5</b>	<b>Клас RegexFSM</b>	<b>6</b>
5.1	Метод <code>_compile</code>	6
5.2	Метод <code>_handle_repetition</code>	7
5.3	Методи керування станами та переходами	7
5.4	Метод <code>_epsilon_closure</code>	8
5.5	Метод <code>check_string</code>	8
5.6	Метод <code>is_full_match</code>	9
<b>6</b>	<b>Алгоритми та методи</b>	<b>10</b>
6.1	Алгоритм побудови НСА (Недетермінованого Скінченного Автомата)	10
6.2	Алгоритм симуляції NFA	11
6.3	Використання епсилон-замикань	11
<b>7</b>	<b>Аналіз ефективності</b>	<b>11</b>
7.1	Просторова складність	11
7.2	Часова складність	11
<b>8</b>	<b>Висновки</b>	<b>11</b>

# 1 Вступ

Представлена реалізація впроваджує кінцевий автомат (Finite State Machine, FSM) для обробки регулярних виразів. Це класичний підхід до розпізнавання шаблонів у тексті, заснований на теорії формальних мов. Реалізація включає базовий абстрактний клас стану, спеціалізовані стани для різних конструкцій регулярних виразів та головний клас для компіляції та обробки шаблонів.

## 2 Загальна архітектура

Реалізація використовує об'єктно-орієнтовану архітектуру з наступними компонентами:

1. Абстрактний базовий клас **State** визначає інтерфейс для усіх типів станів
2. Спеціалізовані класи станів для обробки різних елементів регулярних виразів
3. Головний клас **RegexFSM**, який керує компіляцією шаблону в автомат та виконанням пошуку

Такий підхід обрано для забезпечення розширюваності та зручності додавання нових типів станів або операторів. Ця архітектура є класичною для реалізації скінченних автоматів і дозволяє чітко відокремити відповідальність між обробкою окремих символів та керуванням автоматом у цілому.

## 3 Базовий клас State

### 3.1 Опис та призначення

```
1 class State(ABC):
2     @abstractmethod
3     def __init__(self) -> None:
4         pass
5
6     @abstractmethod
7     def check_self(self, char: str) -> bool:
8         pass
9
10    def check_next(self, next_char: str) -> State | Exception:
11        for state in self.next_states:
12            if state.check_self(next_char):
13                return state
14        raise NotImplementedError("rejected string")
```

Цей абстрактний базовий клас визначає два основних методи:

1. **check\_self(char)** – абстрактний метод, який перевіряє, чи обробляє даний стан вказаний символ.
2. **check\_next(next\_char)** – метод, який передає символ наступним станам і повертає стан, що може його обробити.

## 3.2 Обґрунтування реалізації

Такий підхід використовує принцип єдиної відповідальності: кожен стан відповідає лише за власну логіку обробки, а методи `check_self` і `check_next` забезпечують послідовну перевірку переходів. Це дозволяє реалізувати недетермінований скінченний автомат (NFA), що є природним для регулярних виразів.

Зберігання посилань на наступні стани в атрибуті `next_states` спрощує роботу з переходами між станами без необхідності централізованої таблиці переходів.

## 4 Типи станів

### 4.1 StartState

```
1 class StartState(State):
2     next_states: List[State] = []
3     def __init__(self):
4         self.next_states = []
5     def check_self(self, char):
6         return False
```

Стартовий стан не обробляє жодних символів (`check_self` завжди повертає `False`), він слугує лише як точка входу в автомат. Ця реалізація підкреслює, що старт регулярного виразу не споживає жодних символів.

### 4.2 TerminationState

```
1 class TerminationState(State):
2     next_states: List[State] = []
3
4     def __init__(self):
5         self.next_states = []
6
7     def check_self(self, char):
8         return False
```

Термінальний стан також не обробляє символи, він лише позначає успішне завершення розбору. Цей підхід дозволяє чітко визначити момент, коли регулярний вираз було повністю розпізнано.

### 4.3 DotState

```
1 class DotState(State):
2     next_states: List[State] = []
3
4     def __init__(self):
5         self.next_states = []
6
7     def check_self(self, char: str):
8         return True
```

Стан `.` приймає будь-який символ, тому `check_self` завжди повертає `True`. Це безпосередньо відображає семантику оператора `.` в регулярних виразах.

## 4.4 AsciiState

```
1 class AsciiState(State):
2     next_states: List[State] = []
3     curr_sym = ""
4
5     def __init__(self, symbol: str) -> None:
6         self.next_states = []
7         self.curr_sym = symbol
8
9     def check_self(self, curr_char: str) -> bool:
10        return curr_char == self.curr_sym
```

Стан обробляє конкретний символ із ASCII-таблиці. Метод `check_self` порівнює вхідний символ з очікуваним символом `curr_sym`. Даний підхід дозволяє створити окремі стани для кожного літерального символу в регулярному виразі.

## 4.5 StarState

```
1 class StarState(State):
2     next_states: List[State] = []
3
4     def __init__(self, checking_state: State):
5         self.next_states = []
6         self.checking_state = checking_state
7
8     def check_self(self, char):
9         if self.checking_state.check_self(char):
10            return True
11        for state in self.next_states:
12            if state.check_self(char):
13                return True
14        return False
```

Стан `*` реалізує квантифікатор “нуль або більше” повторень. Він містить посилання на стан, який повинен повторюватися (`checking_state`). Метод `check_self` перевіряє, чи підходить символ для базового стану або для будь-якого з наступних станів.

Така реалізація відображає недетермінованість автомата для операції `*`, де можуть бути активні одночасно кілька станів.

## 4.6 PlusState

```
1 class PlusState(State):
2     next_states: List[State] = []
3
4     def __init__(self, checking_state: State):
```

```

5     self.next_states = []
6     self.checking_state = checking_state
7     self.matched_at_least_one = False
8
9     def check_self(self, char):
10        if self.checking_state.check_self(char):
11            self.matched_at_least_one = True
12            return True
13        if self.matched_at_least_one:
14            for state in self.next_states:
15                if state.check_self(char):
16                    return True
17        return False

```

Стан `+` реалізує квантифікатор “один або більше” повторень. Він відрізняється від `StarState` наявністю прапорця `matched_at_least_one`, що гарантує принаймні одне співпадіння перед переходом до наступних станів. Це безпосередньо відображає семантичну різницю між операторами `*` і `+`.

## 4.7 CharacterClass

```

1 class CharacterClass(State):
2     next_states: List[State] = []
3
4     def __init__(self, class_definition: str):
5         self.next_states = []
6         self.chars = set()
7         self._parse_class(class_definition)
8
9     def _parse_class(self, definition: str):
10        i = 0
11        while i < len(definition):
12            if i + 2 < len(definition) and definition[i+1] == '-':
13                start_char = definition[i]
14                end_char = definition[i+2]
15                for char_code in range(ord(start_char), ord(
16end_char) + 1):
17                    self.chars.add(chr(char_code))
18                    i += 3
19            else:
20                self.chars.add(definition[i])
21                i += 1
22
23        def check_self(self, char: str) -> bool:
24            return char in self.chars

```

Стан для класу символів `[]` обробляє діапазони символів, як-от `[a-z0-9]`. Метод `_parse_class` розбирає визначення класу, обробляючи як окремі символи, так і діапазони, та зберігає їх у множині `chars`. Метод `check_self` просто перевіряє наявність символу в множині.

Використання множини (`set`) забезпечує швидкий пошук символу  $O(1)$ , що важливо для ефективної роботи з класами символів, які можуть містити багато еле-

ментів.

## 5 Клас RegexFSM

Головний клас для компіляції регулярних виразів і виконання пошуку:

```
1 class RegexFSM:
2     def __init__(self, regex_expr: str) -> None:
3         self.pattern = regex_expr
4         self.curr_state = StartState()
5         self.start_state = self.curr_state
6         self.final_state = TerminationState()
7         self.states_map = {}
8         self.states_map[self.start_state] = {"transitions": {}, "
epsilon": set()}
9         self.states_map[self.final_state] = {"transitions": {}, "
epsilon": set()}
10        self._compile(regex_expr)
```

Конструктор ініціалізує:

1. Початковий стан (StartState)
2. Кінцевий стан (TerminationState)
3. Словник-мапу станів і переходів (states\_map)
4. Компілює регулярний вираз у граф станів через виклик \_compile

### 5.1 Метод \_compile

```
1 def _compile(self, pattern: str):
2     if not pattern:
3         self._add_epsilon_transition(self.start_state, self.
final_state)
4         return
5
6     current_state = self.start_state
7     i = 0
8     while i < len(pattern):
9         char = pattern[i]
10        # (
11        self._add_epsilon_transition(current_state, self.final_state)
```

Цей метод є ключовим для компіляції регулярного виразу в граф станів. Він обробляє патерн посимвольно, створюючи відповідні стани та переходи. Метод використовує алгоритм прямого аналізу (лексичний аналіз) регулярного виразу, що забезпечує простоту реалізації та підтримки.

Основні принципи реалізації:

1. Обробка спеціальних конструкцій (класи символів [], квантифікатори \*, +)
2. Створення станів для кожного символу та додавання переходів між ними
3. З'єднання останнього стану з кінцевим через епсилон-перехід

## 5.2 Метод \_handle\_repetition

```
1 def _handle_repetition(self, current_state, base_state, repeat_type
2 ):
3     self._add_state(base_state)
4     if repeat_type == '*':
5         loop_state = StarState(base_state)
6         self._add_state(loop_state)
7         self._add_epsilon_transition(current_state, loop_state)
8         self._add_transition(current_state, base_state, '')
9         self._add_epsilon_transition(base_state, loop_state)
10        self._add_epsilon_transition(loop_state, base_state)
11        return loop_state
12    else:
13        loop_state = PlusState(base_state)
14        self._add_state(loop_state)
15        self._add_transition(current_state, base_state, '')
16        self._add_epsilon_transition(base_state, loop_state)
17        self._add_epsilon_transition(loop_state, base_state)
18        return loop_state
```

Цей метод обробляє квантифікатори \* і +. Структура переходів відрізняється для операторів:

- Для \* додається епсилон-перехід в обхід базового стану (відображає можливість нуля повторень)
- Для + прямий перехід до базового стану (потрібне щонайменше одне повторення)

В обох випадках додаються епсилон-переходи для циклу між базовим станом і квантифікатором.

## 5.3 Методи керування станами та переходами

```
1 def _add_state(self, state):
2     if state not in self.states_map:
3         self.states_map[state] = {"transitions": {}, "epsilon": set
4         ()}
5
6 def _add_transition(self, from_state, to_state, char):
7     if char not in self.states_map[from_state]["transitions"]:
8         self.states_map[from_state]["transitions"][char] = set()
9         self.states_map[from_state]["transitions"][char].add(to_state)
10        from_state.next_states.append(to_state)
11
12 def _add_epsilon_transition(self, from_state, to_state):
13     self.states_map[from_state]["epsilon"].add(to_state)
14     from_state.next_states.append(to_state)
```

Ці допоміжні методи спрощують керування станами та переходами:

1. \_add\_state додає новий стан до загальної мапи станів

2. `_add_transition` додає символний перехід між станами
3. `_add_epsilon_transition` додає епсилон-перехід (перехід без споживання символу)

Використання множин (`set`) для переходів важливе, оскільки воно забезпечує відсутність дублювання станів і покращує пошук  $O(1)$ .

## 5.4 Метод `_epsilon_closure`

```
1 def _epsilon_closure(self, states):
2     """Find all states reachable through epsilon transitions"""
3     closure = set(states)
4     stack = list(states)
5
6     while stack:
7         state = stack.pop()
8         for next_state in self.states_map[state]["epsilon"]:
9             if next_state not in closure:
10                 closure.add(next_state)
11                 stack.append(next_state)
12     return closure
```

Цей метод реалізує пошук епсилон-замикання для множини станів. Він використовує алгоритм пошуку в глибину (DFS) для знаходження всіх станів, досяжних через епсилон-переходи. Цей метод є критичним для реалізації недетермінованого скінченного автомата (NFA), оскільки він дозволяє одночасно перебувати в декількох станах.

## 5.5 Метод `check_string`

```
1 def check_string(self, text: str) -> bool:
2     """
3     Check if the input string contains the regex pattern.
4     """
5     if self.is_full_match(text):
6         return True
7     for start_pos in range(len(text)):
8         current_states = self._epsilon_closure({self.start_state})
9         for i in range(start_pos, len(text)):
10             char = text[i]
11             next_states = set()
12             for state in current_states:
13                 for c, destinations in self.states_map[state]["
14 transitions"].items():
15                     state_obj = next(iter([s for s in [state] if
16 isinstance(s,
17                     (AsciiState, DotState, CharacterClass,
18                     StarState, PlusState))]), None)
19                     if c == '.' or (state_obj and isinstance(
20 state_obj, DotState)):
```



```

17         next_states.update(destinations)
18         elif (c == char or
19             (state_obj and isinstance(state_obj,
AsciiState)
20             and state_obj.curr_sym == char) or
21             (state_obj and isinstance(state_obj,
CharacterClass)
22             and char in state_obj.chars)):
23             next_states.update(destinations)
24         elif state_obj and isinstance(state_obj,
25             (StarState, PlusState)) and state_obj.
check_self(char):
26             next_states.update(destinations)
27             current_states = self._epsilon_closure(next_states)
28             if not current_states:
29                 break
30             if self.final_state in current_states:
31                 return True
32     return False

```

Цей метод перевіряє, чи містить вхідний рядок підрядок, що відповідає шаблону регулярного виразу. Він використовує алгоритм моделювання роботи NFA:

1. Спочатку перевіряє, чи весь рядок відповідає шаблону
2. Для кожної можливої стартової позиції в тексті:
  - Ініціалізує множину поточних станів початковим станом з епсилон-замиканням
  - Для кожного символу з цієї позиції:
    - Знаходить усі переходи, що відповідають поточному символу
    - Додає цільові стани до множини наступних станів
    - Виконує епсилон-замикання для нових станів
    - Перевіряє, чи досягнуто кінцевого стану

Такий підхід є класичним алгоритмом симуляції NFA і забезпечує коректну роботу з недетермінованими переходами.

## 5.6 Метод is\_full\_match

```

1 def is_full_match(self, text: str) -> bool:
2     """
3     Check if the entire input string matches the regex pattern.
4     """
5     current_states = self._epsilon_closure({self.start_state})
6     for char in text:
7         next_states = set()
8         for state in current_states:
9             for c, destinations in self.states_map[state]["
transitions"].items():
10                 if isinstance(state, DotState):

```

```

11         next_states.update(destinations)
12         elif isinstance(state, AsciiState) and state.
curr_sym == char:
13             next_states.update(destinations)
14             elif isinstance(state, CharacterClass) and char in
state.chars:
15                 next_states.update(destinations)
16                 elif c == char:
17                     next_states.update(destinations)
18                 elif c == '.':
19                     next_states.update(destinations)
20                 for next_state in state.next_states:
21                     if next_state.check_self(char):
22                         next_states.add(next_state)
23             current_states = self._epsilon_closure(next_states)
24             if not current_states:
25                 return False
26             return self.final_state in current_states or any(state.is_final
if
27             hasattr(state, 'is_final') else False for state in
current_states)

```

Цей метод перевіряє, чи повністю відповідає вхідний рядок шаблону регулярного виразу. Він також використовує симуляцію NFA, але з однією стартовою позицією і вимогою, щоб фінальний стан був досягнутий в кінці рядка.

Алгоритм схожий на `check_string`, але не перевіряє всі можливі стартові позиції, а лише проходить по тексту з початку до кінця.

## 6 Алгоритми та методи

### 6.1 Алгоритм побудови НСА (Недетермінованого Скінченно-го Автомата)

Для побудови автомата використовується алгоритм Томпсона (Thompson's construction), що переводить регулярний вираз в еквівалентний NFA. Цей алгоритм обрано через його простоту реалізації та зрозумілість. Основні кроки:

1. Розбиття регулярного виразу на базові компоненти (символи, класи символів, оператори)
2. Побудова станів для кожного компоненту
3. З'єднання станів відповідними переходами
4. Обробка квантифікаторів через епсилон-переходи

Алгоритм має часову складність  $O(n)$  для побудови автомата, де  $n$  – довжина регулярного виразу.

## 6.2 Алгоритм симуляції NFA

Для перевірки відповідності рядка регулярному виразу використовується алгоритм симуляції NFA, який має такі переваги:

1. Прямо відображає недетермінований характер автомата
2. Простіше реалізується, ніж перетворення NFA в DFA
3. Може бути більш ефективним для коротких рядків або простих шаблонів

Однак цей метод має гіршу асимптотичну складність  $O(n \cdot m)$ , де  $n$  – довжина рядка,  $m$  – кількість станів автомата.

## 6.3 Використання епсилон-замикань

Для епсилон-замикань використовується алгоритм пошуку в глибину (DFS), що забезпечує коректну обробку епсилон-переходів. Цей підхід обрано через його простоту і ефективність (складність  $O(V+E)$ , де  $V$  – кількість станів,  $E$  – кількість переходів).

# 7 Аналіз ефективності

## 7.1 Просторова складність

- Побудова автомата:  $O(n)$ , де  $n$  – довжина регулярного виразу
- Зберігання автомата:  $O(m)$ , де  $m$  – кількість станів (пропорційна  $n$ )

## 7.2 Часова складність

- Компіляція шаблону:  $O(n)$
- Перевірка рядка:
  - Найгірший випадок:  $O(n \cdot m)$ , де  $n$  – довжина рядка,  $m$  – кількість станів
  - Це відбувається, коли потрібно перевірити багато шляхів через автомат

# 8 Висновки

Реалізація регулярних виразів через кінцевий автомат є класичним і теоретично обґрунтованим підходом. Обрана архітектура забезпечує:

1. **Розширюваність:** Легко додавати нові типи станів або оператори
2. **Зрозумілість:** Кожен компонент має чітку відповідальність
3. **Коректність:** Реалізація безпосередньо відповідає теоретичному визначенню регулярних виразів

Використання NFA (недетермінованого скінченного автомата) замість DFA зроблено для спрощення реалізації та покращення підтримуваності коду, хоча це впливає на ефективність роботи з великими текстами.

Загалом, дана реалізація представляє збалансований підхід між ефективністю, зрозумілістю, і можливістю подальшого розширення функціоналу.