

Документация и теория работы парсера регулярных выражений

Супрун Артём Сергеевич

14 апреля 2025 г.

Содержание

1	Введение	2
2	Теоретическая основа	2
2.1	Метод рекурсивного спуска	2
2.2	Обработка регулярных выражений	2
3	Описание структуры данных	2
3.1	Типы узлов	2
3.2	Структура узла	3
4	Функции обработки ошибок и создания узлов	3
4.1	Функция <code>error</code>	3
4.2	Функция <code>create_node</code>	4
5	Рекурсивный спуск: функции синтаксического анализа	4
5.1	<code>parse_expression</code>	4
5.2	<code>parse_term</code>	4
5.3	<code>parse_factor</code>	5
5.4	<code>parse_base</code>	5
6	Функции вывода и освобождения памяти	6
6.1	<code>print_tree</code>	6
6.2	<code>free_tree</code>	7
7	Основная функция программы	7
8	Заключение	8

1 Введение

В данном документе представлена теоретическая основа и реализация парсера регулярных выражений на языке C. Использование метода рекурсивного спуска позволяет поэтапно обрабатывать различные компоненты выражения — от литералов до операторов конкатенации, альтернативы и повторения (оператор `*`).

2 Теоретическая основа

2.1 Метод рекурсивного спуска

Рекурсивный спуск — это метод синтаксического анализа, при котором для каждой нетерминальной грамматики создаётся отдельная функция. Каждая функция принимает указатель на текущую позицию во входной строке и обновляет его по мере анализа. Такой подход позволяет эффективно обрабатывать вложенные структуры и обеспечивать понятное разделение логики.

Грамматика регулярных выражений в рассматриваемом примере выглядит следующим образом:

$$\begin{aligned}\text{Expression} &\rightarrow \text{Term} \{ ' | ' \text{Term} \} \\ \text{Term} &\rightarrow \text{Factor} \{ \text{Factor} \} \\ \text{Factor} &\rightarrow \text{Base} \{ ' * ' \} \\ \text{Base} &\rightarrow ' (' \text{Expression} ') ' \mid \text{Литерал}\end{aligned}$$

Каждая конструкция реализована отдельной функцией, создающей соответствующие узлы дерева разбора.

2.2 Обработка регулярных выражений

Регулярные выражения — мощный инструмент для поиска и обработки строк. В нашем случае синтаксический анализ включает:

- **Литералы** — отдельные символы.
- **Конкатенацию** — последовательное соединение символов и подвыражений.
- **Альтернативу** — выбор между несколькими выражениями с помощью оператора `' | '`.
- **Клини-стар** — оператор повторения `(' * ')`, позволяющий повторять предыдущий элемент неограниченное число раз.

3 Описание структуры данных

Для представления дерева разбора используются структура `Node` и перечисление возможных типов узлов.

3.1 Типы узлов

Определены следующие типы узлов:

- **NODE_LITERAL**: узел, представляющий конкретный символ.

- **NODE_CONCAT**: узел для конкатенации двух подвыражений.
- **NODE_ALTERNATION**: узел для альтернативы (оператор '|').
- **NODE_STAR**: узел для оператора повторения (оператор '*').

3.2 Структура узла

Структура Node включает следующие поля:

- **type** — тип узла.
- **c** — символ, используемый для узлов типа литерал.
- **left** — указатель на левый подузел (например, первый операнд).
- **right** — указатель на правый подузел (например, второй операнд для конкатенации или альтернативы).

Пример определения:

```

1 typedef enum {
2     NODE_LITERAL,
3     NODE_CONCAT,
4     NODE_ALTERNATION,
5     NODE_STAR
6 } NodeType;
7
8 typedef struct Node {
9     NodeType type;
10    char c; //
11    struct Node *left; //
12    struct Node *right; //
13 } Node;
```

Листинг 1: Определения типов узлов и структуры

4 Функции обработки ошибок и создания узлов

4.1 Функция error

Функция `error` выводит сообщение об ошибке. При наличии символа он включается в сообщение, после чего программа завершается с кодом ошибки.

```

1 void error(const char *message, char symbol) {
2     if (symbol != '\0') {
3         fprintf(stderr, ": %s '%c'\n", message, symbol);
4     } else {
5         fprintf(stderr, ": %s\n", message);
6     }
7     exit(EXIT_FAILURE);
8 }
```

Листинг 2: Функция обработки ошибок

4.2 Функция create_node

Функция `create_node` выделяет память для нового узла, инициализирует его поля и возвращает указатель на узел. При неудачном выделении памяти вызывается функция `error`.

```
1 Node* create_node(NodeType type, char c, Node* left, Node* right) {
2     Node *node = (Node *)malloc(sizeof(Node));
3     if (!node) {
4         error(" ", '\0');
5     }
6     node->type = type;
7     node->c = c;
8     node->left = left;
9     node->right = right;
10    return node;
11 }
```

Листинг 3: Функция создания узла

5 Рекурсивный спуск: функции синтаксического анализа

Ниже приведены функции, реализующие парсер с рекурсивным спуском. Каждая функция принимает аргумент `input` — указатель на указатель на текущую позицию во входной строке, что обеспечивает удобное продвижение по тексту.

5.1 parse_expression

Функция `parse_expression` анализирует выражение согласно грамматике:

$$\text{Expression} \rightarrow \text{Term} \{ '|' \text{Term} \}$$

При обнаружении оператора `'|'` создаётся узел типа `NODE_ALTERNATION`.

```
1 Node* parse_expression(const char **input) {
2     Node *node = parse_term(input);
3     while (**input == '|') {
4         (*input)++; // '|'
5         Node *right = parse_term(input);
6         node = create_node(NODE_ALTERNATION, '\0', node, right);
7     }
8     return node;
9 }
```

Листинг 4: Парсинг выражения

5.2 parse_term

Функция `parse_term` анализирует терм по правилу:

$$\text{Term} \rightarrow \text{Factor} \{ \text{Factor} \}$$

Конкатенация реализована посредством последовательного объединения результатов разбора с созданием узлов типа `NODE_CONCAT`.

```

1 Node* parse_term(const char **input) {
2     Node *node = parse_factor(input);
3     while (**input && **input != ')') {
4         Node *next = parse_factor(input);
5         node = create_node(NODE_CONCAT, '\0', node, next);
6     }
7     return node;
8 }

```

Листинг 5: Парсинг термина

5.3 parse_factor

Функция `parse_factor` обрабатывает фактор по правилу:

$$\text{Factor} \rightarrow \text{Base} \{ '*' \}$$

При каждом обнаружении оператора `'*'` создаётся узел типа `NODE_STAR`.

```

1 Node* parse_factor(const char **input) {
2     Node *node = parse_base(input);
3     while (**input == '*') {
4         (*input)++; // '*'
5         node = create_node(NODE_STAR, '\0', node, NULL);
6     }
7     return node;
8 }

```

Листинг 6: Парсинг фактора

5.4 parse_base

Функция `parse_base` обрабатывает базовые элементы:

- Если встречается символ `'('`, вызывается функция `parse_expression`, после чего ожидается символ `')'`.
- Если символ не является управляющим (то есть не `'('`, `')'`, `'|'`) и не является концом строки, создаётся узел типа `NODE_LITERAL`.

В случае обнаружения ошибки (например, отсутствия закрывающей скобки) вызывается функция `error`.

```

1 Node* parse_base(const char **input) {
2     if (**input == '(') {
3         (*input)++; // '('
4         Node *node = parse_expression(input);
5         if (**input != ')') {
6             error("' )'", **input);
7         }
8         (*input)++; // ')'
9         return node;
10    }
11    if (**input == '\0' || **input == '|' || **input == ')') {

```

```

12     error(" ", **input);
13 }
14 char c = **input;
15 (*input)++; //
16 return create_node(NODE_LITERAL, c, NULL, NULL);
17 }

```

Листинг 7: Парсинг базового элемента

6 Функции вывода и освобождения памяти

6.1 print_tree

Функция `print_tree` рекурсивно выводит дерево разбора в удобном для чтения формате. В зависимости от типа узла выводятся:

- Для литерала — конкретный символ.
- Для конкатенации и альтернативы — соответствующие теги ("CONCAT(" или "ALT(") с рекурсивным выводом поддеревьев.
- Для оператора повторения — тег "STAR(" с выводом содержимого.

```

1 void print_tree(Node *node) {
2     if (!node)
3         return;
4
5     switch (node->type) {
6         case NODE_LITERAL:
7             printf("%c", node->c);
8             break;
9         case NODE_CONCAT:
10            printf("CONCAT(");
11            print_tree(node->left);
12            printf(", ");
13            print_tree(node->right);
14            printf(")");
15            break;
16        case NODE_ALTERNATION:
17            printf("ALT(");
18            print_tree(node->left);
19            printf(", ");
20            print_tree(node->right);
21            printf(")");
22            break;
23        case NODE_STAR:
24            printf("STAR(");
25            print_tree(node->left);
26            printf(")");
27            break;
28        default:
29            break;
30    }

```

Листинг 8: Вывод дерева разбора

6.2 free_tree

После завершения работы с деревом важно освободить выделенную память. Функция `free_tree` рекурсивно проходит по дереву в порядке post-order и освобождает память для каждого узла.

```

1 void free_tree(Node *node) {
2     if (!node)
3         return;
4     free_tree(node->left);
5     free_tree(node->right);
6     free(node);
7 }
```

Листинг 9: Освобождение памяти дерева

7 Основная функция программы

Функция `main` выполняет следующие шаги:

1. Проверяет, передано ли регулярное выражение в качестве аргумента командной строки.
2. Вызывает функцию `parse_expression` для построения дерева разбора.
3. При наличии не обработанных символов выводит сообщение об ошибке.
4. Выводит дерево разбора с помощью `print_tree`.
5. Освобождает память, занятую деревом, с помощью `free_tree`.

Пример кода функции `main`:

```

1 int main(int argc, char *argv[]) {
2     if (argc < 2) {
3         fprintf(stderr, ": %s <_>\n", argv[0]);
4         return EXIT_FAILURE;
5     }
6     const char *input = argv[1];
7     Node *root = parse_expression(&input);
8     if (*input != '\0') {
9         fprintf(stderr, ": : %s\n", input);
10        free_tree(root);
11        return EXIT_FAILURE;
12    }
13    printf(" : \n");
14    print_tree(root);
15    printf("\n");
16    free_tree(root);
17    return EXIT_SUCCESS;
}
```

8 Заключение

В документе подробно рассмотрена реализация парсера регулярных выражений на языке С, основанного на методе рекурсивного спуска. Были рассмотрены:

- Теоретическая основа синтаксического анализа регулярных выражений.
- Структура данных для представления дерева разбора.
- Функции для создания узлов, обработки ошибок, построения и вывода дерева.
- Основной алгоритм разбора, реализованный с помощью функций `parse_expression`, `parse_term`, `parse_factor` и `parse_base`.

Такой подход обеспечивает эффективный синтаксический анализ входной строки и позволяет использовать построенное дерево для дальнейшей интерпретации или компиляции регулярных выражений.