

# 解析器组合子的原理与实现

叶家炜

Apr 22, 2025

在计算机科学领域，解析器（Parser）是将原始数据转换为结构化表示的核心工具。无论是编译器处理源代码、解释器执行脚本，还是配置文件读取，解析器都扮演着至关重要的角色。传统解析技术如正则表达式或自动生成工具（如 Yacc/Bison）往往面临可维护性差或灵活性不足的问题。解析器组合子（Parser Combinators）通过函数式编程的组合思想，提供了一种优雅解决方案。

## 1 解析器组合子基础

解析器组合子的核心在于「组合」二字。它将解析器视为一等公民（First-class Parser），允许通过高阶函数将简单解析器组合成复杂解析器。例如一个匹配数字的解析器可以与匹配运算符的解析器组合，最终形成算术表达式解析器。这种方法的优势在于代码可读性强、扩展灵活，并且天然适配函数式编程范式。

## 2 解析器组合子的原理

### 2.1 解析器的类型定义

解析器的本质是一个函数：接收输入字符串，返回解析结果。在 TypeScript 中可以定义为：

```
1 type Parser<T> = (input: string) => ParseResult<T>;  
  
3 interface ParseResult<T> {  
    success: boolean;  
5    value?: T;  
    remaining?: string;  
7    error?: string;  
}
```

这里 `Parser<T>` 表示生成类型 `T` 的解析器。`ParseResult` 包含解析是否成功、解析值、剩余字符串和错误信息。例如解析数字 `123` 后，`value` 可能为 `123`，`remaining` 为后续字符串。

### 2.2 基本解析器构建

字符解析器是最基础的构建单元。以下是一个匹配特定字符的解析器实现：

```
const char = (c: string): Parser<string> => (input) => {
```

```

2  if (input[0] === c) {
    return {
4      success: true,
      value: c,
6      remaining: input.slice(1)
    };
8  }
  return {
10    success: false,
    error: `Expected '${c}', got '${input[0]}'`
12  };
};

```

该函数接收目标字符 `c`，返回一个解析器。当输入字符串首字符匹配时，返回成功结果；否则返回错误信息。类似地，可以构建 `string` 解析器来匹配完整字符串。

## 2.3 组合子操作

组合子的威力在于将原子解析器组合成复杂结构。以「交替组合子」为例：

```

1  const or = <T>(p1: Parser<T>, p2: Parser<T>): Parser<T> => (input) => {
    const result1 = p1(input);
3   if (result1.success) return result1;
    return p2(input);
5  };

```

此组合子尝试用第一个解析器 `p1` 解析输入，若失败则尝试 `p2`。例如 `or(char('a'), char('b'))` 将匹配 `'a'` 或 `'b'`。类似地，`and` 组合子用于串联解析器：

```

1  const and = <T, U>(p1: Parser<T>, p2: Parser<U>): Parser<[T, U]> => (input) => {
    const result1 = p1(input);
3   if (!result1.success) return { success: false, error: result1.error };

    const result2 = p2(result1.remaining || '');
    if (!result2.success) return { success: false, error: result2.error };
7
    return {
9      success: true,
      value: [result1.value, result2.value],
11     remaining: result2.remaining
    };
13 };

```

此实现中, `and` 依次应用 `p1` 和 `p2`, 并将两者的结果合并为元组。例如 `and(char('a'), char('b'))` 将匹配序列 `ab`。

## 3 实现一个简单的解析器组合子库

### 3.1 重复组合子的实现

处理重复结构是解析常见需求。以下 `many` 组合子实现了零次或多次匹配：

```
1 const many = <T>(p: Parser<T>): Parser<T[]> => (input) => {
  const values: T[] = [];
3  let remaining = input;

5  while (true) {
    const result = p(remaining);
7    if (!result.success) break;
    values.push(result.value!);
9    remaining = result.remaining || '';
  }

11  return { success: true, value: values, remaining };
13 };
```

该组合子循环应用解析器 `p` 直到失败, 收集所有成功结果。例如 `many(char('a'))` 可以匹配 `aaa` 或空字符串。

### 3.2 错误处理增强

精确的错误定位对调试至关重要。通过扩展解析结果类型记录位置信息：

```
1 interface ParseResult<T> {
  // ... 原有字段
3  position?: number;
}

5
const withPosition = <T>(p: Parser<T>): Parser<T> => (input) => {
7  const result = p(input);
  if (!result.success) {
9    return {
      ...result,
11     position: input.length - (result.remaining?.length || 0)
    };
  }
13 }

return result;
```

```
15 };
```

此时错误信息可以提示具体出错位置，例如 Expected 'a' at position 5。

## 4 实战：用解析器组合子解析 JSON

### 4.1 JSON 值解析器

JSON 值的解析需要处理多种可能类型。通过 or 组合子实现分发逻辑：

```
1 const jsonValue: Parser<JsonValue> = (input) =>
  or(jsonString,
3    or(jsonNumber,
      or(jsonObject,
5        or(jsonArray,
          or(jsonBoolean,
7            jsonNull
          )
        )
      )
    )
  )(input);
```

此处 JsonValue 是联合类型，包含字符串、数字、对象等可能性。每个分支对应具体类型的解析器。

### 4.2 对象解析器实现

JSON 对象由键值对组成，需处理花括号和逗号分隔符：

```
const jsonObject: Parser<JsonObject> = (input) => {
2  const parser = and(
    and(
4      skipWhitespace(char('{')),
      many(
6          and(
              jsonString,
8              and(
                  skipWhitespace(char(':')),
10             jsonValue
            )
          )
        )
      )
    ),
14 }
```

```
    skipWhitespace(char('}'))
16 );

18 const result = parser(input);
    if (!result.success) return result;

20

22 const entries = result.value[0][1];
    const obj = Object.fromEntries(entries.map(([, v]) => [k, v]));
    return { success: true, value: obj, remaining: result.remaining };
24 };
```

此处 `skipWhitespace` 用于忽略空格，`many` 处理多个键值对，`map` 将结果转换为字典对象。

## 5 优化与进阶话题

### 5.1 左递归处理

传统递归下降解析器难以处理左递归文法，如  $\text{Expr} \rightarrow \text{Expr} + \text{Term}$ 。解析器组合子可通过惰性求值解决：

```
const expr: Parser<Expr> = lazy(() =>
2   or(
    and(expr, and(char('+'), term), (left, [op, right]) => new Add(left, right)),
4   term
  )
6 );
```

`lazy` 包装器延迟解析器的初始化，避免立即执行导致的栈溢出。

### 5.2 记忆化优化

通过缓存解析结果避免重复计算，提升性能：

```
const memoize = <T>(p: Parser<T>): Parser<T> => {
2   const cache = new Map<string, ParseResult<T>>();
    return (input) => {
4     const key = input;
        if (cache.has(key)) return cache.get(key)!;
6     const result = p(input);
        cache.set(key, result);
8     return result;
    };
10 };
```

此技术特别适用于复杂文法的解析，可将时间复杂度从指数级降为线性。

解析器组合子通过函数组合的抽象方式，提供了一种高表达力的解析方案。其优势在于代码的可读性和可维护性，但在处理大规模数据时需谨慎性能优化。现代库如 Haskell 的 Parsec 或 Rust 的 nom 已展示了该技术的工业级应用潜力。未来随着类型系统的发展，结合依赖类型或线性类型可能进一步提升解析器的安全性与效率。