

现代编程思想

案例：语法解析器

Hongbo Zhang

语法解析器

- 案例目标

- 解析基于自然数的数学表达式: `"(1+ 5) * 7 / 2"`

- 转化为单词列表

- `LParen Value(1) Plus Value(5) Multiply Value(7) Divide Value(2)`

- 转化为抽象语法树

- `Division(Multiply(Add(Value(1), Value(5)), Value(7)), Value(2))`

- 计算最终结果: 21

- 语法分析

- 对输入文本进行分析并确定其语法结构

- 通常包含词法分析和语法分析

- 本节课均利用语法解析器组合子 (parser combinator) 为例

- 将输入分割为单词
 - 输入：字符串/字节块
 - 输出：单词流
 - 例如： "12 +678" -> [Value(12), Plus, Value(678)]
- 通常可以通过有限状态自动机完成
 - 一般用领域特定语言定义后，由软件自动生成程序
- 算术表达式的词法定义

```
1. Number      = %x30 / (%x31-39) *(%x30-39)
2. LParen      = "("
3. RParen      = ")"
4. Plus        = "+"
5. Minus       = "-"
6. Multiply    = "*"
7. Divide      = "/"
8. Whitespace  = " "
```

词法分析

- 算术表达式的词法定义

```
1. Number = %x30 / (%x31-39) *(%x30-39)
2. Plus   = "+"
```

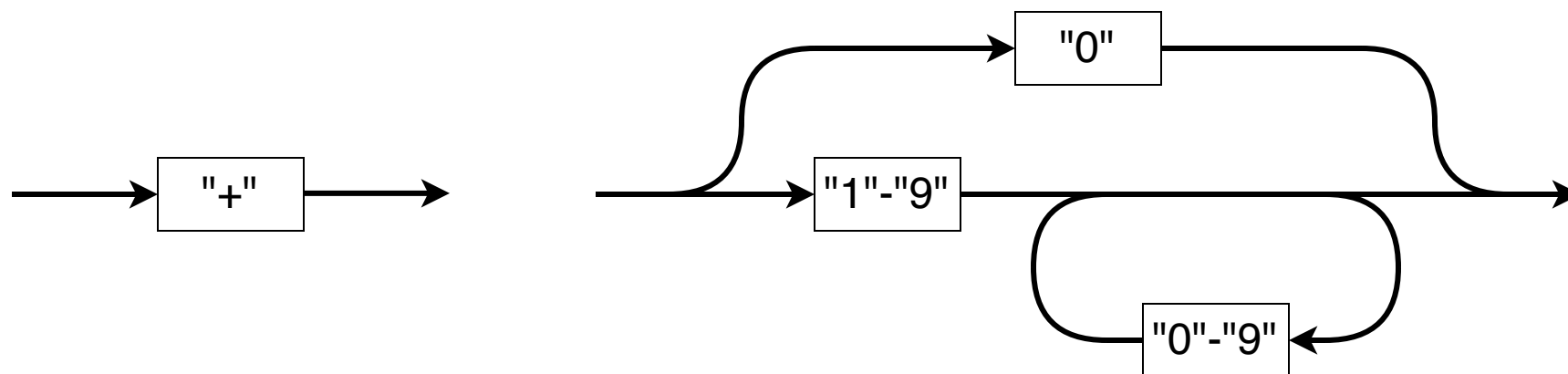
- 每一行对应一个匹配规则

- "xxx" : 匹配内容为xxx的字符串
- a b : 匹配规则a, 成功后匹配规则b
- a / b : 匹配规则a, 匹配失败则匹配规则b
- *a : 重复匹配规则a, 可匹配0或多次
- %x30 : UTF编码十六进制表示为30的字符 ("0")

词法分析

- 算术表达式的词法定义

```
1. Number = %x30 / (%x31-39) * (%x30-39)
2. Plus   = "+"
```



- 单词定义

```
1. enum Token {
2.     Value(Int); LParen; RParen; Plus; Minus; Multiply; Divide
3. } derive(Show)
```

解析器组合子

- 构造可组合的解析器

```
1. // V 代表解析成功后获得的值
2. // @string.View 代表 String 的一部份
3. type Lexer[V] (@string.View) -> (V, @string.View)?
4.
5. fn Lexer::parse[V](self : Lexer[V], str : @string.View) -> (V, @string.View)? {
6.   self.inner()(str)
7. }
```

- 我们简化处理报错信息以及错误位置 (可以使用 `Result[A, B]`)

- 最简单的解析子：判断下一个待读取的字符是否符合条件，符合则读取并前进

```
1. fn pchar(predicate : (Char) -> Bool) -> Lexer[Char] {  
2.   Lexer(input => match input {  
3.     [ch, .. rest] if predicate(ch) => Some((ch, rest))  
4.   } -> None  
5. })  
6. }
```

- 例如

```
1. test "pchar" {  
2.   inspect(pchar(ch => ch is 'a').parse("asdf"),  
3.     content=(  
4.       #|Some(('a', "sdf"))  
5.     ))  
6.   inspect(pchar(ch => ch is 'a').parse("sdf"), content="None")  
7. }
```

词法分析

- 单词定义：数字或左右括号或加减乘除

```
1. enum Token {  
2.     Value(Int)  
3.     LParen; RParen; Plus; Minus; Multiply; Divide  
4. } derive(Show)
```

- 分析运算符、括号、空白字符等

```
1. let symbol: Lexer[Char] = pchar(ch => ch  
2.   is ('+' | '-' | '*' | '/' | '(' | ')'))  
3. let whitespace : Lexer[Char] = pchar(ch => ch is ' ')
```


解析器组合子

- 如果解析成功，对解析结果进行转化

```
1. fn[I, O] Lexer::map(self : Lexer[I], f : (I) -> O) -> Lexer[O] {  
2.   Lexer(fn(input) { self.parse(input).map(pair => (f(pair.0), pair.1)) })  
3. }
```

- 分析运算符、括号并映射为对应的枚举值

```
1. let symbol : Lexer[Token] = pchar(ch => ch  
2.   is ('+' | '-' | '*' | '/' | '(' | ')')).map(token => match token {  
3.   '+' => Token::Plus; '-' => Token::Minus  
4.   '*' => Token::Multiply; '/' => Token::Divide  
5.   '(' => Token::LParen; ')' => Token::RParen  
6.   _ => panic()  
7. })
```

解析器组合子

- 解析 `a`，如果成功再解析 `b`，并返回 `(a, b)`

```
1. fn[V1, V2] Lexer::then(self : Lexer[V1], parser2 : Lexer[V2]) -> Lexer[(V1, V2)] {  
2.   Lexer(fn(input) {  
3.     guard self.parse(input) is Some((value, rest)) else { return None }  
4.     guard parser2.parse(rest) is Some((value2, rest2)) else { return None }  
5.     Some(((value, value2), rest2))  
6.   }) }
```

- 解析 `a`，如果失败则解析 `b`

```
1. fn[Value] Lexer::or(self : Lexer[Value], parser2 : Lexer[Value]) -> Lexer[Value] {  
2.   Lexer(fn(input) {  
3.     match self.parse(input) {  
4.       None => parser2.parse(input)  
5.       Some(_) as result => result  
6.     }  
7.   }) }
```

解析器组合子

- 重复解析 `a`，零或多次，直到失败为止

```
1. fn[Value] Lexer::many(self : Lexer[Value]) -> Lexer[@list.T[Value]] {
2.   Lexer(fn(input) {
3.     loop (input, @list.empty()) {
4.       (rest, cumul) =>
5.         match self.parse(rest) {
6.           None => Some((cumul.rev(), rest)) // List 是栈，需要反转
7.           Some((value, rest)) => continue (rest, @list.construct(value, cumul))
8.         }
9.     }
10.  })
11. }
```

词法分析

- 整数分析

```
1. // 通过字符编码将字符转化为数字
2. let zero : Lexer[Int] = pchar(ch => ch is '0').map(_ => 0)
3.
4. let one_to_nine : Lexer[Int] = pchar(ch => ch is ('1'..'9')).map(ch => ch.to_int() - '0'.to_int())
5.
6. let zero_to_nine : Lexer[Int] = pchar(ch => ch is ('0'..'9')).map(ch => ch.to_int() - '0'.to_int())
7.
8. // number = %x30 / (%x31-39) * (%x30-39)
9. let value : Lexer[Token] = zero
10. .or(
11.   one_to_nine
12.   .then(zero_to_nine.many())
13.   .map(pair => {
14.     let (i, ls) = pair
15.     ls.fold((i, j) => i * 10 + j, init=i)
16.   }),
17. )
18. .map(Token::Value(_))
```

- 对输入流进行分析
 - 在单词之间可能存在空格

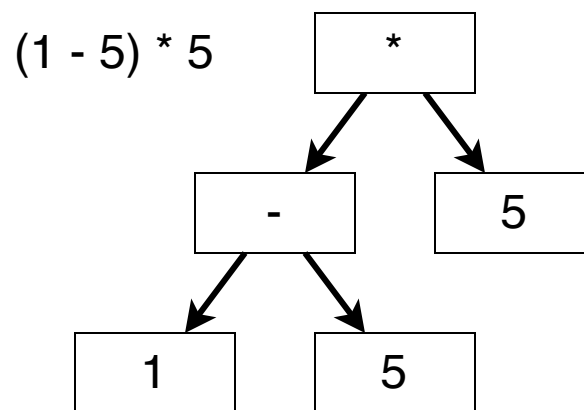
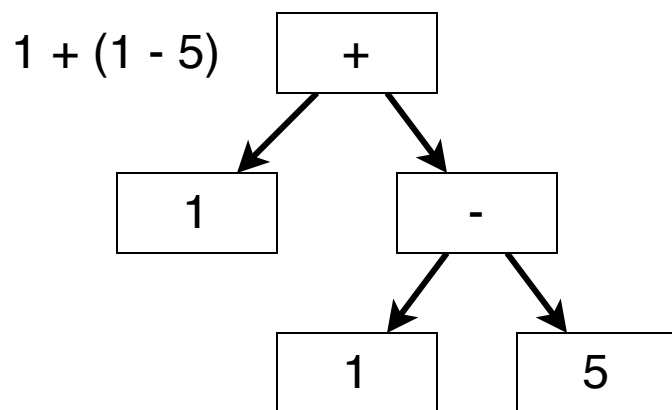
```
1. let tokens : Lexer[@list.T[Token]] = value
2.   .or(symbol)
3.   .then(whitespace.many())
4.   .map(p => p.0)
5.   .many()
6. test {
7.   inspect(
8.     tokens.parse("-10123+--523 103    ( 5 ) ").unwrap(),
9.     content=(
10.      #|(@list.of([Minus, Value(10123), Plus, Minus, Plus, Value(523), Value(103), LParen, Value(5), RParen, RParen]), "")
11.    ),
12.  )
13. }
```

- 我们成功地分割了字符串
 - - 10123 - + - 523 103 (5))
 - 但这不符合数学表达式的语法

语法分析

- 对单词流进行分析，判断是否符合语法
 - 输入：单词流
 - 输出：抽象语法树

```
1. expression = Value / "(" expression ")"  
2. expression =/ expression "+" expression / expression "-" expression  
3. expression =/ expression "*" expression / expression "/" expression
```



语法分析

- 语法定义

```
1. expression = Value / "(" expression ")"  
2. expression =/ expression "+" expression / expression "-" expression  
3. expression =/ expression "*" expression / expression "/" expression
```

- 问题：运算符的优先级、结合性

- 优先级： $a + b \times c \rightarrow a + (b \times c)$
- 结合性： $a + b + c \rightarrow (a + b) + c$
- 当前语法具有二义性

语法分析

- 修改后的语法定义

```
1. atomic      = Value / "(" expression ")"
2. combine     = atomic / combine "*" atomic / combine "/" atomic
3. expression = combine / expression "+" combine / expression "-" combine
```

- 注意到除了简单的组合以外，出现了左递归
 - 左递归会导致我们的解析器进入循环
 - 解析器将尝试匹配运算符左侧的规则而不前进
 - 拓展：自底向上解析器可以处理左递归

语法分析

- 修改后的语法定义

```
1. atomic      = Value / "(" expression ")"
2. combine     = atomic * ( ("*" / "/" ) atomic )
3. expression = combine * ( ("+" / "-") combine )
```

- 数据结构

```
1. enum Expression {
2.     Number(Int)
3.     Plus(Expression, Expression)
4.     Minus(Expression, Expression)
5.     Multiply(Expression, Expression)
6.     Divide(Expression, Expression)
7. }
```

语法解析

- 定义语法解析组合子

```
1. type Parser[V] (@list.T[Token]) -> (V, @list.T[Token])?  
2.  
3. fn[V] Parser::parse(self : Parser[V], tokens : @list.T[Token]) -> (V, @list.T[Token])? { self.inner()(tokens) }
```

- 大部分组合子与 `Lexer[V]` 类似
- 递归组合: `atomic = Value / "(" expression ")"`
 - 延迟定义
 - 递归函数

递归定义

- 延迟定义
 - 利用引用定义 `Ref[Parser[V]]` : `struct Ref[V] { mut val : V }`
 - 在定义其他解析器后更新引用中内容

```
1. fn[Value] Parser::from_ref(ref_ : Ref[Parser[Value]]) -> Parser[Value] {  
2.   Parser(fn(input) { ref_.val.parse(input) })  
3. }
```

- `ref.val` 将在使用时获取，此时已更新完毕

递归定义

- 延迟定义

```
1. fn parser() -> Parser[Expression] {
2.   // 首先定义空引用
3.   let expression_ref : Ref[Parser[Expression]] = { val : Parser(_ => None) }
4.
5.   // atomic = Value / "(" expression ")"
6.   let atom = // 利用引用定义
7.     (lparen.then(Parser::from_ref(expression_ref)).then(rparen).map(expr => expr.0.1).or(number)
8.
9.   // combine = atomic * ( ("*" / "/" ) atomic)
10.  let combine = atom.then(multiply.or(divide).then(atom).many()).map(pair => {
11.    guard pair is (expr, list)
12.    list.fold(init=expr, (expr, op_expr) => match op_expr { ... })
13.  })
14.
15.  // expression = combine * ( "+" / "-" ) combine)
16.  expression_ref.val = combine.then(plus.or(minus).then(combine).many()).map(pair => {
17.    guard pair is (expr, list)
18.    list.fold(init=expr, (expr, op_expr) => match op_expr { ... })
19.  })
20.
21.  expression_ref.val
22. }
```

递归定义

- 递归函数
 - 解析器本质上是一个函数
 - 定义互递归函数后，将函数装进结构体

```
1. fn recursive_parser() -> Parser[Expression] {
2.   // 定义互递归函数
3.   // atomic = Value / "(" expression ")"
4.   letrec atom = fn(tokens: @list.T[Token]) -> (Expression, @list.T[Token])? {
5.     lparen.then( Parser(expression) ).and(rparen).map(expr => expr.0.1)
6.     .or(number).parse(tokens)
7.   }
8.   and combine = fn(tokens: @list.T[Token]) -> (Expression, @list.T[Token])? { ... }
9.   and expression = fn (tokens: @list.T[Token]) -> (Expression, @list.T[Token])? { ... }
10.
11.   // 返回函数代表的解析器
12.   Parser(expression)
13. }
```

语法树之外：Tagless Final

- 计算表达式，除了生成为抽象语法树再解析，我们还可以有其他的选择
- 我们通过“行为”来进行抽象

```
1. trait Expr : Add + Sub + Mul + Div {  
2.   number(Int) -> Self  
3. }
```

- 接口的不同实现即是对行为的不同语义

语法树之外：Tagless Final

- 我们利用行为的抽象定义解析器

```
1. fn[E : Expr] recursive_parser() -> Parser[E] {
2.   let number : Parser[E] = ptoken(token => token is Value(_)).map(ptoken => {
3.     guard ptoken is Value(i)
4.     E::number(i) // 利用抽象的行为
5.   })
6.
7.   letrec atomic = fn(tokens: @list.T[Token]) -> (E, @list.T[Token])? { ... }
8.   // 转化为 a * b * c * ... 和 a / b / c / ...
9.   and combine = fn(tokens: @list.T[Token]) -> (E, @list.T[Token])? { ... }
10.  // 转化为 a + b + c + ... 和 a - b - c - ...
11.  and expression=fn(tokens: @list.T[Token]) -> (E, @list.T[Token])? { ... }
12.
13.  Parser(expression)
14. }
15.
16. // 结合在一起
17. fn[E : Expr] parse_string(
18. str : @string.View,
19. ) -> (E, @string.View, @list.T[Token])? {
20.   guard tokens.parse(str) is Some((token_list, rest_string)) else { return None }
21.   guard recursive_parser().parse(token_list) is Some((expr, rest_token)) else { return None }
22.   Some((expr, rest_string, rest_token))
23. }
```

语法树之外：Tagless Final

- 我们可以提供不同的实现，获得不同的诠释

```
1. enum Expression { ... } derive(Show) // 语法树实现
2. type BoxedInt Int derive(Show) // 整数实现
3. // 实现接口 (此处省略其他方法)
4. fn BoxedInt::number(i: Int) -> BoxedInt { BoxedInt(i) }
5. fn Expression::number(i: Int) -> Expression { Number(i) }
6. test {
7.   // 获得语法树
8.   inspect(
9.     (
10.      parse_string("1 + 1 * (307 + 7) + 5 - 3 - 2") :
11.      (Expression, @string.View, @list.T[Token])?.unwrap(),
12.      content=(
13.        #|(Minus(Minus(Plus(Plus(Number(1), Multiply(Number(1), Plus(Number(307), Number(7))))), Number(5)), Number(3)), Number(2)), "", @list.of([]))
14.      ),
15.    )
16.    // 获得计算结果
17.    inspect(
18.      (
19.        parse_string("1 + 1 * (307 + 7) + 5 - 3 - 2") :
20.        (BoxedInt, @string.View, @list.T[Token])?.unwrap(),
21.        content=(
22.          #|(BoxedInt(315), "", @list.of([]))
23.        ),
24.      )
25.    }
```


总结

- 本节课展示了一个语法解析器
 - 介绍了词法解析的概念
 - 介绍了语法解析的概念
 - 展示了语法解析组合子的定义与实现
 - Tagless Final的概念与实现
- 拓展阅读
 - 调度场算法
 - 斯坦福CS143 第1-8课 或
 - 《编译原理》前五章 或
 - 《现代编译原理》前三章
- 拓展练习
 - 实现兼容各类“流”的语法解析组合子