

КОМБИНАТОРЫ ПАРСЕРОВ

Михаил Беляев

18 ноября 2015

Комбинаторы парсеров (также называемые «монадными комбинаторами парсеров») впервые появились в 2008 году в статье Frost, Hafiz & Callaghan.

Референсная реализация кода к статье была написана на Haskell.

ЗАЧЕМ ВООБЩЕ ЭТО?

Комбинаторы парсеров позволяют:

- Разбивать разбор грамматик на основе рекурсивного спуска на модули
- Решить при этом (частично) проблему разбора контекстно-зависимых грамматик за полиномиальное время

Кроме того, отдельные части становятся предельно простыми в сравнении с традиционными ad-hoc парсерами и средствами генерации парсеров

На данный момент самой известной является библиотека `Parsec`, самой быстрой — библиотека `Attoparsec`

Существует большое количество подражателей в других языках:

- `JParsec` и `ParsecJ` для Java
- `Benny` для Javascript
- `Parsy`, `Picoparse`, `FuncparserLib` для Python
- `Boost::Spirit` и другие для C++
- Стандартные библиотеки `Scala`, `F#` и некоторых других языков
- и т.д.

- Есть базовый набор простых парсеров, которые умеют разбирать примитивный ввод
- Есть набор простых комбинаторов, позволяющих на основе одних парсеров строить другие

ПРОСТЫЕ ПАРСЕРЫ

```
-- Принимаем 1 символ, возвращаем его же
char :: Char -> Parser Char
-- Принимаем символы из диапазона
range :: Char -> Char -> Parser Char
-- Принимаем символы из списка
oneOf :: [Char] -> Parser Char
```

ПРОСТЫЕ ПАРСЕРЫ

```
-- Принимаем 1 символ, возвращаем его же
char :: Char -> Parser Char
-- Принимаем символы из диапазона
range :: Char -> Char -> Parser Char
-- Принимаем символы из списка
oneOf :: [Char] -> Parser Char
```

Или, более общий вариант

```
conforms :: (Char -> Bool) -> Parser Char
char c = conforms (== c)
range a b = conforms (\c -> c >= a && c <= b)
oneOf cs = conforms (\c -> c `elem` cs)
```

```
-- Не пытаться ничего читать,  
-- сразу вернуть ошибку  
parserError :: Error -> Parser a  
-- Не пытаться ничего читать,  
-- сразу вернуть значение  
parserSuccess :: a -> Parser a
```



```
-- Всё, кроме заданного
parserNot :: (Parser Char) -> Parser Char
-- Попробовать первый, если не сработал, то второй
parserOr :: Parser a ->
           Parser b ->
           Parser (Either a b)
```

```
-- Несколько элементов подряд
parserSeq :: Parser a ->
           Parser b ->
           Parser (a,b)
-- Много (0 или больше) элементов подряд
parserMany :: Parser a -> Parser [a]
```

```
-- Парсер является функтором!  
parserMap :: (a -> b) -> Parser a -> Parser b  
-- Парсер является монадой!  
parserReturn = parserSuccess  
parserFail = parserError  
parserFlatten :: (Parser (Parser a)) -> Parser a
```

```
parserTry :: Parser a -> Parser a
```

Особенность реализации: иногда нужно попробовать разобрать текст и в случае неудачи вернуться туда, откуда был начат разбор

СИНТАКСИЧЕСКИЙ САХАР ДЛЯ УПОТРЕБИМЫХ ОПЕРАЦИЙ

```
space :: Parser ()  
eof  :: Parser ()  
digit :: Parser Char  
alpha :: Parser Char  
upper :: Parser Char  
...
```

ПРИМЕР: РАЗБОР НОМЕРА ГРУППЫ

```
data GroupNumber = GroupNumber Integer Integer
parseNumber = read <$> parserMany1 digit
parseGroupNumber = do{ dept <- parseNumber;
                        char '/';
                        gp <- parseNumber;
                        return $ GroupNumber dept gp }
parseManyGNs = do{ gn <- parseGroupNumber;
                   gns <- parserMany parseGroupNumber';
                   return (gn: gns) }
where
  parseGroupNumber' = do { char ',';
                          parseGroupNumber }
```

- Парсер это не просто монада, а монада-трансформер
- Позволяет в процесс разбора притаскивать своё состояние
- Поддерживает чтение из множества разных типов
- `parser0r` называется `<|>`
- В остальном все парсер-комбинаторы похожи