

# Monadic Parsing



1

# Schedule

			State.hs	
Tue Nov 12 / Fri Nov 15	Monadic Parsing	Ch. 13	<a href="#">Lecture_Week12</a> <a href="#">Parser.hs</a> , <a href="#">Lecture12.hs</a>	<a href="#">PS8</a> , <a href="#">Data.hs</a> , <a href="#">State.hs</a> , <a href="#">Parser.hs</a> , <a href="#">Lexer.hs</a> , <a href="#">Infer.hs</a> , <a href="#">Ps8.hs</a> <b>Checkpoint #1: attend office hours this week (or earlier)</b>
Tue Nov 19 / Fri Nov 22	More parsing and Parsec		<a href="#">Lecture_Week13</a>	<b>5-8 min presentation in class on Friday</b>
Tue Nov 26	Property Testing; QuickCheck		Quiz 5 on Tue	PS8 due on Tuesday
Tue Dec 3 Fri Dec 6	TBD		Quiz 6 on Fri	<b>Checkpoint #2: attend office hours this week (or earlier)</b>
Tue Dec 10	Project presentations			<b>Project due</b> <b>5-8 min presentation in class</b>

Programming in Haskell, A Milanova

2

# Outline

- Lexing and parsing
- Parser type
- Parser monad
- Lexer primitives
  
- Recursive-descent monadic parsing
- Expression evaluation
  
- Parsec

Programming in Haskell, A Milanova

3

# Recursive-Descent Parsing

- Monadic parsers are recursive-descent parsers
- In recursive-descent there is a procedure for each nonterminal in the grammar, in our example it is *expr*, *term*, and *factor*
- In the monadic parser there is a parsing function for each nonterminal. E.g., *expr* parser corresponds to *expr* nonterminal:

```
expr :: Parser String
expr = do t <- term
        symbol "+"
        e <- expr
        return (t ++ " + " ++ e)
<|> term
```

- To parse a nonterminal, call (i.e., descend) into corresponding routine

4

3

4

## Recursive-Descent Parsing

- Production right-hand-side forms body of procedure. Parser tries one production, then next, and so on
  - E.g.,  $expr ::= term + expr \mid term$   
Parser first tries to parse  $term + expr$ ; if fail, it tries  $term$
- Parsing right-hand-side means calling procedures for nonterminals and consuming terminals, in turn
  - E.g.,  $term + expr$  first calls  $term$ , then (if success) consumes  $+$  then calls  $expr$

```

expr :: Parser String
expr = do t <- term
        symbol "+"
        e <- expr
        return (t ++ " " ++ e)
<|> term
    
```

Programming in Haskell, A Milan

5

5

## Recursive-Descent Parsing

- Parsing an alternative may fail in which case recursive-descent backtracks (i.e., puts consumed input back) and tries next alternative
  - E.g.,  $x1*x2*x3$ . Parser tries  $term + expr$  but fails trying to match  $+$ . It backtracks to the beginning and tries  $term$

```

expr :: Parser String
expr = do t <- term
        symbol "+"
        e <- expr
        return (t ++ " " ++ e)
<|> term
    
```

Programming in Haskell, A Milanova

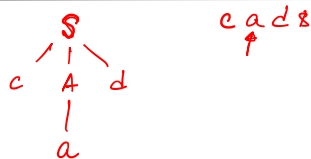
6

6

## Backtracking Example

- Consider grammar
 
$$S ::= c A d$$

$$A ::= a b^x \mid a$$
 and input string cad\$



- Is there backtracking with this grammar?
 
$$expr ::= term + expr \mid term$$

$$term ::= factor * term \mid factor$$

$$factor ::= ( expr ) \mid identifier \mid integer$$

abc\*cd\$

Programming in Haskell, A Milanova

7

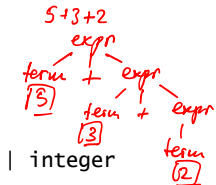
7

- What we can and what we cannot parse with recursive-decent?

- This grammar can be parsed:
 
$$expr ::= term + expr \mid term$$

$$term ::= factor * term \mid factor$$

$$factor ::= ( expr ) \mid identifier \mid integer$$



- How about this "better" disambiguation of original  $expr$  grammar:
 
$$expr ::= expr + term \mid term$$

$$term ::= term * factor \mid factor$$

$$factor ::= ( expr ) \mid identifier \mid integer$$

LEFT RECURSION IS BAD.

Programming in Haskell, A Milanova

8

8

*Expr ~ term + term + term    id + id \* id*

- For some grammars recursive-descent does not backtrack
- We can rewrite our grammar like this:
 

```

      expr ::= term term_tail
      term_tail ::= + term term_tail | ε
      term ::= factor factor_tail
      factor_tail ::= * factor factor_tail | ε
      factor ::= ( expr ) | identifier | integer
      
```
- A classic LL(1) grammar meaning parser can predict the production by looking at just one token of lookahead
- The monadic parser does not consume and backtrack; symbol does a "lookahead" of at most one token and decides. Thus, a lot more efficient than previous parser!

9

9

- We can rewrite our grammar one more time:
 

```

      expr ::= term term_tail
      term_tail ::= + term term_tail | ε
      
```

 becomes
 

```

      expr ::= term term_tail
      term_tail ::= + expr | ε
      
```

 becomes
 

```

      expr ::= term (+ expr | ε)
      
```

*TEXTBOOK.*

```

expr :: Parser String
expr = do t <- term
      do symbol "+"
      e <- expr
      return (t ++ " + " ++ e)
<|> return t
  
```

10

10

## Evaluation

*1 + 2 - 3*

- Look at simplified version of grammar:
 

```

      expr ::= integer ttail
      ttail ::= + integer ttail | - integer ttail | ε
      
```
- Let us write a parser for this grammar. It parses into strings:

```

expr :: Parser String
expr =
  do i <- integer
      tt <- ttail
      return ((show i)++tt)

ttail :: Parser String
ttail =
  do s <- symbol "+ "<|> symbol "-"
      i <- integer
      tt <- ttail
      return (s ++ (show i) ++ tt)
  <|> return ""
  
```

*> runParser expr "1+2-3"*

11

11

- Look at simplified version of grammar:
 

```

      expr ::= integer ttail
      ttail ::= + integer ttail | - integer ttail | ε
      
```
- Now let's look at parse tree of expression 5 - 3 - 2

```

      expr
      /  |  \
Integer  -  ttail
 [5]      /  |  \
          Integer - ttail
          [3]   /  |  \
               Integer - ttail
               [2]   |
                   ε
  
```

12

12

```

expr ::= integer ttail
ttail ::= + integer ttail | - integer ttail | ε

```

■ We want to build a parser that also evaluates expression. First try:

```

ttail :: Parser Int
ttail = do symbol "+"
          i <- integer
          tt <- ttail
          return (i + tt)
        <|> do symbol "-"
          i <- integer
          tt <- ttail
          return (i - tt)
        <|> return 0

```

```

expr :: Parser Int
expr =
do i <- integer
  tt <- ttail
  return (i + tt)

```

Programming in Haskell, A Milanova 13

13

```

expr ::= integer ttail
ttail ::= + integer ttail | - integer ttail | ε

```

■ It's wrong! Parser evaluates 5 - 3 - 2 incorrectly. Here is how:

■ Problem of course is that grammar is right-recursive and tree tilts to the right while we are doing a bottom-up evaluation

Programming in Haskell, A Milanova 14

14

```

expr ::= integer ttail
ttail ::= (+ | -) integer ttail | ε

```

■ Parse of our running example 5 - 3 - 2 with lists of tuples:

Programming in Haskell, A Milanova 15

15

■ One way to fix

```

expr ::= integer ttail
ttail ::= (+ | -) integer ttail | ε

```

```

type IntOp = Int -> Int -> Int

expr :: Parser Int
expr = do i <- integer
          -- tt :: [(IntOp,Int)] i.e., flattens tree
          tt <- ttail
          let comb :: Int -> (IntOp, Int) -> Int
              comb x (op,y) = op x y
              -- foldl tt list around i
              return (foldl comb i tt)

```

5 [(-), 3], (-), 2]

Programming in Haskell, A Milanova 16

16

```

expr ::= integer ttail
ttail ::= (+ | -) integer ttail | ε

```

- Parser `ttail` parses into a more complex structure, a list of tuples:

```

symbol' :: Parser IntOp
symbol' = do symbol "+"
           return (+)
         <|> do symbol "-"
           return (-)

ttail :: Parser [(IntOp,Int)]
ttail = do s ← symbol'
           i ← integer
           tt ← ttail
           return (s,i):tt
         <|> return []

```

Programming in Haskell, A Milanova 17

17

```

expr ::= integer ttail
ttail ::= (+ | -) integer ttail | ε

```

- Another way to fix, using a State transformer

```

expr :: Parser Int
expr = do i <- integer
         -- tt :: Parser (S.State Int ())
         -- encodes transform of rest-of-tree
         tt <- ttail
         return (S.runState tt i)

```

Programming in Haskell, A Milanova 18

18

```

expr ::= integer ttail
ttail ::= (+ | -) integer ttail | ε

```

- Parse of our running example `5 - 3 - 2` with state transformer:

Programming in Haskell, A Milanova 19

19

```

expr ::= integer ttail
ttail ::= (+ | -) integer ttail | ε

```

- Another way to fix, using a State transformer


```

expr :: Parser Int
expr = do i <- integer
         -- tt :: Parser (S.State int ())
         -- encodes transform of rest-of-tree
         tt <- ttail
         let (_,i') = S.runState tt i
         return i'

```

Programming in Haskell, A Milanova 20

20




$expr ::= integer\ ttail$   
 $ttail ::= (+ | -)\ integer\ ttail | \epsilon$

- Parser `ttail` parses into a state transformer:

```
ttail :: Parser (S.State Int ())
ttail = do
  op <- symbol
  i <- integer
  tt <- ttail
  -- encode new st based on state transformer tt
  return (do { pr <- S.get; S.put (op pr i); tt; return () })
  -- return (S.get >>= (\pr ->
  --   let (_,pr') = S.runState tt (op pr i) in S.put pr'))
  -- if epsilon, return transformer \pr -> ((),pr)
  <|> return (do { pr <- S.get; return () })
  -- return (S.get >>= (\pr -> S.put pr))
```

Programming in Haskell, A Milanova 21

21




## Exercise

- Expand grammar with terms with multiplication and integer division (as opposed to just integer terms)
- Should be able to generalize to an arbitrary L-attributed interpretation of the tree and possibly shift-reduce parsing

Programming in Haskell, A Milanova 22

22




## Outline

- Lexing and parsing
- Parser type
- Parser monad
- Lexer primitives
- Recursive-descent monadic parsing
- Expression evaluation
- Parsec**

Programming in Haskell, A Milanova 23

23



## Parsec

- We wrote parsing primitives and combinators from “first principles” in Haskell
- We then built token parsers (i.e., the lexer)
- We then combined token parsers to build more complex parsers, e.g., for CFGs (i.e., known as a parser in compilers)
- We then added interpretation of structure
- You can use Parsec! To install: `cabal install parsec`

Programming in Haskell, A Milanova 24

24

## Parser Types

Using Text.Parsec

- ParsecT is most general:

```
type ParsecT s u m a
-- Input stream s, typically String
-- User state u wrapped in a monad, e.g., Int
-- Monad to wrap output m, e.g., [], Maybe
-- Output type a, e.g., Char
```

- Parsec instantiates ParsecT with Identity monad, i.e., no wrapping:

```
type Parsec s u = ParsecT s u Identity
```

Programming in Haskell, A Milanova

25

25

- To run a Parsec parser:

```
> :t runParser
... =>
Parsec s u a -> u -> SourceName -> s -> (Either ParseError a)
> runParser digit () "" "9abc"
Right '9'
> runParser digit () "" "abc"
Left (line 1, column 1):
unexpected "a"
expecting digit
```

- An even simpler parse, no user state:

```
> parse digit "" "9abc"
Right '9'
> parse digit "" "abc"
Left ... -- Error
```

Pr

26

26

- Remember our Parser a type:

```
newtype Parser a = P (String -> [(a,String)])
runParser (P f) = f
```

- We can define a similar specialized parser in terms of Parsec:

```
type Parser a = Parsec String () a
runParser' :: Parser a -> String -> Either ParseError a
runParser' p = runParser p () ""
> runParser' digit "987"
Right '9'
```

Programming in Haskell, A Milanova

27

27

## Character Primitives

- Character primitives are defined in terms of tokenPrim since we don't have access to the stream (it is enclosed into a monad)

- The item primitive:

```
item :: Parser Char
item = tokenPrim showChar nextPos Just
  where
    showChar x = "'" ++ [x] ++ "'"
    nextPos pos x xs = updatePosChar pos x

> runParser' item "abc"
Right 'a' -- just as item we defined last week
```

Programming in Haskell, A Milanova

28

28

- The `sat` primitive. Note: there is `satisfy` in `Text.Parsec`.

```

sat :: (Char -> Bool) -> Parser Char
sat p
  = tokenPrim showChar nextPos testChar
  where
    showChar x = "'" ++ [x] ++ "'"
    nextPos pos x xs = updatePosChar pos x
    testChar x = if p x then Just x else Nothing

> runParser' (sat isDigit) "9abc"
Right '9' -- just as sat we defined last week

```

- We can now derive all lexer functions: `natural`, `integer`, `symbol`, and `identifier`; just like we did last week with the `Parser` type. Then we can derive more complex parsers

Programming in Haskell, A Milanova 29

29

## Parsec Lexer

- But we don't have to! We can use the `Text.Parsec.Token`
- `emptyDef` is one of the default language definitions; another one is `haske11Def`; these are records
- We can redefine some features, e.g., identifier starting character and remaining identifier characters:

```

import Text.Parsec
import Text.Parsec.Pos (updatePosChar, updatePosString)
import qualified Text.Parsec.Token as P
import Text.Parsec.Language (emptyDef)

langDef = emptyDef { P.identStart = lower
                    }
                    { P.identLetter = alphaNum

```

Programming in Haskell, A Milanova 30

30

- Then we can instantiate the lexer:

```

-- The lexer
lexer = P.makeTokenParser langDef

identifier :: Parsec String u String
identifier = P.identifier lexer
natural = P.natural lexer
integer = P.integer lexer
symbol = P.symbol lexer

```

- Typically, we just reuse one of the language definitions; if we need a special lexer, we define one from scratch using `tokenPrim`

Programming in Haskell, A Milanova 31

31

## Parsec Expression Parser

- So, let's build a parser for our LL(1) grammar using those lexer primitives. We'll parse back into a `String`

```

expr ::= integer ttail
ttail ::= (+ | -) integer ttail | ε

```

```

expr :: Parser String
expr = do i <- integer
         let tail ← ttail
             return ((show i) ++ tail)

```

SAME AS BEFORE

```

ttail :: Parser String
ttail = do ...

```

32

32



## Evaluation

- Evaluation takes advantage of user state. Need the Parsec type

```
expr ::= integer ttail getState and putState
ttail ::= (+ | -) integer ttail | ε
```

```
expr :: Parsec String Integer Integer
expr = do i <- integer
        putState i
        tt <- ttail
        total <- getState
        return total
        -- or just ttail
```

```
ttail :: Parsec String Integer Integer
ttail = do op <- symbol
          ... i <- integer
          pr <- getState
          putState (op pr i)
          tt <- ttail
          total <- getState
          return total
          <|> getState
```

```
> runParser expr 0 "" "5-3-2"
```

33

## Exercise

*With Parsec*

*MORE COMPLEX STATE.*

- Let's evaluate expression using this slightly different grammar:

```
expr ::= integer + expr | integer - expr | integer
```

```
expr' :: Parsec String (Integer, IntOp) Integer
expr' = try (do i <- integer
              op <- symbol'
              (v, op') <- getState
              putState ((op' v i), op)
              e <- expr
              return e)
        <|> do i <- integer
              (v, op) <- getState
              return (op v i)
```

34

## Some Notes

- Important: Parsec's combinator **do not** backtrack by construction (unlike Parser a we wrote which did backtrack)
- Need to force backtracking with try
- Extend the parsers with more productions and handling of eof

35