

# Lecture 12 – Parsec



## Outline:

- ▶ Brief review of BNF grammar and parsing
- ▶ Parsing with combinators
- ▶ Creating parsers using Parsec combinators

# Grammar, syntax, and parsing



- ▶ Parsing is the processing of recognizing a string by breaking it down to a set of symbols and analyzing each one against the grammar of the language.
- ▶ A grammar describes how to form strings from a language's alphabet that are valid according to the language's syntax.
- ▶ A parser checks whether a string belongs to a language and transforms the string into a data structure such as a parse tree.
- ▶ A parser is implemented based on a grammar that defines the syntax of a language.



- ▶ BNF grammar is a type of context-free grammar that is commonly used to define programming language syntax.
- ▶ BNF grammar consists of 4 components: tokens (or terminals), non-terminals, production rules, and a start symbol.
  - ▶ Tokens are indivisible unit of syntax such as an identifier, keyword, operators, and literals.
  - ▶ Non-terminals such as expressions and statements are defined with tokens and other non-terminals.
  - ▶ Production rules are equations that define the derivation of non-terminals.
  - ▶ A start symbol is a unique non-terminal that corresponds to the entire program.

# An example of BNF grammar



## ► Arithmetic expressions

```
1 <expr>      ::= <expr> '+' <term>
2           | <expr> '-' <term>
3           | <term>
4
5 <term>       ::= <term> '*' <fact>
6           | <term> '/' <fact>
7           | <fact>
8
9 <fact>       ::= '(' <expr> ')'
10          | <integer>
11
12 <integer>   ::= <digit>
13          | <integer> <digit>
14
15 <digit>     ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
```

<expr>, <term>, <fact>, <integer>, <digit> are non-terminals while +, -, /, \*, (,), and 0 – 9 are tokens.  
<expr> is the root symbol.

# An example of EBNF grammar



## ► Arithmetic expressions in EBNF grammar

```
1 <expr>      ::= <term> {('+' | '-' ) <term>}
2
3 <term>       ::= <fact> {('*' | '/') <fact>}
4
5 <fact>       ::= '(' <expr> ')'
6   | <integer>
7
8 <integer>    ::= <digit> {<digit>}
9
10 <digit>     ::= [0-9]
```

{<digit>} means zero or more repetitions of <digit>

[0-9] means a digit between 0 and 9.

('+' | '-')

means choosing between '+' and '-'.

# An example of parsing

- ▶ Parsing an string can return an abstract syntax tree:

```
1 data Exp = Plus Exp Exp
2           | Minus Exp Exp
3           | Times Exp Exp
4           | Div Exp Exp
5           | Const Integer
```

- ▶ Parsing "(1 + 2) \* 3" results in the following value.

```
1 Times (Plus (Const 1) (Const 2)) (Const 3)
```

- ▶ Parsing "(1 + 2) \* 3" can also directly evaluate it to 9.

## Construct a parser

- ▶ A parser can be constructed from a parser-generator such as Lex, Yacc, Bison, JavaCC, Jison by providing a BNF phrase grammar (with a separate Lexical grammar).
- ▶ A parser can also be written directly using a parser library such as Parsec.
- ▶ Parsec has a collection of combinators that can be used to construct a LL( $k$ ) parser (Left-to-right, Leftmost derivation), where  $k$  is the number tokens of lookahead when parsing a sentence.

## Construct a parser 🔊

- ▶ A LL parser is a top-down parser that can be implemented with a deterministic push-down automaton.

Given  $\Gamma = N \cup \Sigma$ , where  $N$  is the set of non-terminals and  $\Sigma$  is the set of tokens including end-of-file eof.

```
1 stack = [s, eof] -- s is the start symbol
2
3 parse stack
4   | stack == [eof] = True
5   | otherwise =
6
7     let x = head stack
8     in if x `elem` sigma -- x is a token
9         then if readInput == x
10             then parse $ tail stack
11             else False
12         else let a = rules x -- rule x ::= a
13             in parse $ a ++ tail stack
```

The crucial problem is to select the correct production rule to replace non-terminal  $x$  on top of the stack.

## Example of parsing



- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')'  
2 E ::= 'i'
```

- ▶ To parse the following input,

```
1 (i+i)
```

we apply the following sequence of rules

1	Rule	Stack	Stack-action	Input-action
2	1	[E eof]	pop E, push (E+E)	
3		[ ( E + E ) eof]	pop '(', read '('	
4	2	[E + E ) eof]	pop E, push 'i'	
5		[i + E ) eof]	pop 'i', read 'i'	
6		[+ E ) eof]	pop '+', read '+'	
7	2	[E ) eof]	pop E, push 'i'	
8		[i ) eof]	pop 'i', read 'i'	
9		[ ) eof]	pop ')', read ')'	
10		[eof]		

## Example of parsing – Count i's

- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')'  
2 E ::= 'i'
```

- ▶ Think of a parser as a function from input String to an output that contains the result of the parsing and the remaining input not consumed by the parser.

```
1 type Parser a = String    -- input string  
2                 -> (a,          -- result of parsing  
3                           String) -- remaining string  
4 -- E ::= 'i'  
5 p2 :: Parser Int           -- count the # of i's  
6 p2 ('i':rest) = (1, rest) -- successful and return 1  
7 p2 input      = (0, input) -- unsuccessful
```

# Example of parsing – Count i's 🔊

- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')'  
2 E ::= 'i'
```

- ▶ We quickly run into problems.

```
1 type Parser a = String    -- input string  
2           -> (a,          -- result of parsing  
3                  String) -- remaining string  
4 -- E ::= 'i'  
5 p2 :: Parser Int  
6 p2 ('i':rest) = (1, rest) -- successful and return 1  
7 p2 input       = (0, input) -- unsuccessful  
8  
9 -- E ::= '(' E '+' E ')'  
10 p1 :: Parser Int  
11 p1 ('(': rest) = let (a, rest') = p1 rest'  
12           -- or let (a, rest') = p2 rest'  
13           -- which one do we call?  
14 p1 input = (0, input)
```

# Example of parsing – Count i's

- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')' | 'i'
```

- ▶ Use  $\langle \rangle$  to encode the choice between two parsers.

```
1 type Parser a = String -> (a, String)
2
3 (<|>) :: Parser a -> Parser a -> Parser a
4 p = p1 <|> p2 -- if p1 fails to parse, then try p2
5
6 -- E ::= 'i'
7 p2 :: Parser Int
8 p2 ('i':rest) = (1, rest) -- successful and return 1
9
10 -- E ::= '(' E '+' E ')'
11 p1 :: Parser Int          -- too complicated already
12 p1 ('(': rest) = let (a1, rest1) = p rest
13           in case rest1 of
14             ('+' : rest2) ->
15               let (a2, rest3) = p rest2
16               in case rest3 of
17                 (')', rest4) -> (a1 + a2, rest4)
```

# Example of parsing – Count i's



- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')' | 'i'
```

- ▶ Make Parser a state monad.

```
1 newtype Parser a = Parser {runParser :: String -> (a, String)}
2
3 char :: Char -> Parser Char
4
5 p = p1 <|> p2
6
7 -- E ::= 'i'
8 p2 :: Parser Int
9 p2 = char 'i' >> return 1 -- one i
10
11 -- E ::= '(' E1 '+' E2 ')'
12 p1 :: Parser Int
13 p1 = do char '('
14     a1 <- p    -- # of i's in E1
15     char '+'
16     a2 <- p    -- # of i's in E2
17     char ')'
18     return $ a1 + a2
```

## Example of parsing – Count i's

- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')' | 'i'
```

- ▶ Run  $p$ : unwrap its function, apply it to input, extract result.

```
1 newtype Parser a = Parser {runParser :: String -> (a, String)}
2
3 char :: Char -> Parser Char
4
5 p = p1 <|> p2
6
7 p2 = char 'i' >> return 1
8
9 p1 = do char '('
10           a1 <- p
11           char '+'
12           a2 <- p
13           char ')'
14           return $ a1 + a2
15
16 run p x = fst $ runParser p x
17
18 run p "(i+(i+i))" -- returns 3
```

# Example of parsing – Max depth of parenthesis



- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')' | 'i'
```

- ▶ p parse input to find the max depth of parenthesis.

```
1 newtype Parser a = Parser {runParser :: String -> (a, String)}
2
3 char :: Char -> Parser Char
4
5 p = p1 <|> p2
6
7 p2 = char 'i' >> return 0 -- depth is 0
8
9 -- E ::= '(' E1 '+' E2 ')'
10 p1 = do char '('
11     a1 <- p -- depth of parenthesis in E1
12     char '+'
13     a2 <- p -- depth of parenthesis in E2
14     char ')'
15     return $ (max a1 a2) + 1
16
17 run p "(i+(i+i))" -- returns 2
```

# Example of parsing – Abstract syntax tree



- ▶ Given the BNF grammar below where  $E$  is the start symbol.

```
1 E ::= '(' E '+' E ')' | 'i'
```

- ▶  $p$  returns a Node value –  $p :: \text{Parser Node}$

```
1 newtype Parser a = Parser {runParser :: String -> (a, String)}
2
3 data Node = Plus Node Node | I
4
5 p = p1 <|> p2
6
7 p2 = char 'i' >> return I
8
9 -- E ::= '(' E1 '+' E2 ')'
10 p1 = do char '('
11     a1 <- p -- Node for E1
12     char '+'
13     a2 <- p -- Node for E2
14     char ')'
15     return $ Plus a1 a2
16
17 run p "(i+(i+i))" -- Plus I (Plus I I)
```

## Example of parsing – XML tags

- ▶ XML is the start symbol; {x} means zero or more copies of x.

```
1     XML ::= OpenTag { XML } EndTag
2 OpenTag ::= '<' Name '>'
3 EndTag ::= '<' '/' Name '>'
```

- ▶ parse an XML string with matching open and end tags.

```
1 data XML = Tag String [XML]
2
3 openTag = do {char '<'; name <- word; char '>'; return name}
4
5 endTag name = string $ "</" ++ name ++ ">"
6
7 xml = do name <- openTag      -- name of open tag
8         content <- many xml
9         endTag name          -- matching end tag
10        return $ Tag name content
11
12 -- new combinators:
13 --   string :: String -> Parser ()
14 --   many    :: Parser a -> Parser [a]
15 --   word    :: Parser String
```

# Example of parsing – XML tags



- ▶ parse an XML string with matching open and end tags.

```
1 data XML = Tag String [XML]
2
3 openTag = do {char '<'; name <- word; char '>'; return name}
4
5 endTag name = string $ "</" ++ name ++ "/"
6
7 xml = do name <- openTag      -- name of open tag
8           content <- many xml
9           endTag name        -- matching end tag
10          return $ Tag name content
11
12 -- many    :: Parser a -> Parser [a]
13 -- word    :: Parser String
```

- ▶ 'string' parses a specific string and return () if successful.

```
1 string :: String -> Parser ()
2
3 string ""      = return ()                  -- parse empty string
4 string (c:cs) = char c >> string cs -- parse char one by one
```

## Example of parsing – XML tags 🔊

- ▶ parse an XML string with matching open and end tags.

```
1 data XML = Tag String [XML]
2
3 openTag = do {char '<'; name <- word; char '>'; return name}
4
5 endTag name = string $ "</" ++ name ++ ">"
6
7 xml = do name <- openTag      -- name of open tag
8         content <- many xml
9         endTag name        -- matching end tag
10        return $ Tag name content
11
12 -- word :: Parser String
```

- ▶ 'many' runs  $p$  zero or more times – returns  $p$ 's output in a list.

```
1 many :: Parser a -> Parser [a]
2 many p = many1 p <|> return []
3
4 many1 p = do   -- many1 runs p one or more times
5     x <- p
6     xs <- many1 p <|> return []
7     return $ x:xs
```

# Example of parsing – XML tags

- ▶ parse an XML string with matching open and end tags.

```
1 data XML = Tag String [XML]
2
3 openTag = do {char '<'; name <- word; char '>'; return name}
4
5 endTag name = string $ "</" ++ name ++ ">"
6
7 xml = do name <- openTag      -- name of open tag
8           content <- many xml
9           endTag name          -- matching end tag
10          return $ Tag name content
```

- ▶ 'word' parses any string consisting of one or more letters.

```
1 word :: Parser String
2 word = many1 letter
3
4 letter :: Parser Char   -- parse an alphabet
```

## Example of parsing – XML tags 🔊

- ▶ parse an XML string with matching open and end tags.

```
1 data XML = Tag String [XML]
2
3 openTag = do {char '<'; name <- word; char '>'; return name}
4
5 endTag name = string $ "</" ++ name ++ "/"
6
7 xml = do name <- openTag      -- name of open tag
8         content <- many xml
9         endTag name          -- matching end tag
10        return $ Tag name content
```

- ▶ But this can't even parse<sup>1</sup> "<a></a>".

```
1 run xml "<a></a>"
2 -- parse error at (line 1, column 5):
3 -- unexpected "/"
4 -- expecting letter
```

Open tag and end tag share <. After <a>, the parser expects another open tag or </a>. After openTag fails, it can't try endTag since < is consumed – no backtracking.

---

<sup>1</sup>Error handling is not discussed yet.

## Example of parsing – XML tags 🔊

- ▶ A parser without backtracking is LL(1).

Use *try p* to enable backtracking for *p* to allow LL(*k*),  $k \geq 2$ .

```
1 try :: Parser a -> Parser a
```

- ▶ *try p* has effect only if it runs on the left of some ⟨|⟩.

```
1 data XML = Tag String [XML]
2
3 openTag = do {char '<'; name <- word; char '>'; return name}
4
5 endTag name = string $ "</" ++ name ++ ">"
6
7 xml = do name <- try openTag -- try parse open tag
8     content <- many xml
9     endTag name           -- matching end tag
10    return $ Tag name content
11
12 run xml "<a></a>"
13 -- Tag "a" []
14
15 run xml "<a><b><c></c><d></d></b></a>"
16 -- Tag "a" [Tag "b" [Tag "c" [], Tag "d" []]]
```

## Example of parsing – Sentence

- ▶ A sentence is a sequence of words separated by separators such as space and comma and terminated by tokens such as dot and question mark.

```
1 sentence :: Parser [String]
2 sentence = do words <- sepBy1 word separator
3                     terminator
4                     return words
5
6 separator :: Parser ()
7 separator = skipMany1 (space <|> char ',')
8
9 terminator = oneOf ".?!"
10 space = oneOf "\t"
11
12 --      sepBy1 :: Parser a1 -> Parser a2 -> Parser [a1]
13 --      oneOf :: [Char] -> Parser Char
14 --      skipMany1 :: Parser a -> Parser ()
15
16 run sentence "Coding is easy, but debugging is hard!"
17 -- ["Coding","is","easy","but","debugging","is","hard"]
18
19 run sentence "a, b, , c."
20 -- ["a","b","c"]
```

## Example of parsing – Sentence



```
▶ 1 sentence :: Parser [String]
  2 sentence = do words <- sepBy1 word separator
  3             terminator
  4             return words
  5
  6 separator :: Parser ()
  7 separator = skipMany1 (space <|> char ',')
  8
  9 terminator = oneOf ".?!"
10 space = oneOf " \t"
11
12 --      oneOf :: [Char] -> Parser Char
13 --      skipMany1 :: Parser a -> Parser ()
```

- ▶ *sepBy1 p sep* parses one or more occurrences of *p* separated by *sep* and return a list of values returned by *p*.

```
1 sepBy1 p sep = do
  2   x <- p
  3   xs <- many (sep >> p)
  4   return $ x:xs
  5
  6 -- zero or more occurrences of p separated by sep
  7 sepBy p sep = sepBy1 p sep <|> return []
```

## Example of parsing – Sentence



```
▶ 1 sentence :: Parser [String]
  2 sentence = do words <- sepBy1 word separator
                  terminator
                  return words
  3
  4
  5
  6 separator :: Parser ()
  7 separator = skipMany1 (space <|> char ',')
  8
  9 terminator = oneOf ".?!"
10 space = oneOf "\t"
```

▶ *oneOf str* parses one of the char in *str*.

*skipMany1 p* parses one or more occurrences of *p* and throws away the results.

```
1 oneOf :: [Char] -> Parser Char
2
3 skipMany1 :: Parser a -> Parser ()
4 skipMany1 p = many1 p >> return ()
5
6 -- parse zero or more occurrence of p and skip results.
7 skipMany p = many p >> return ()
```

# Example of parsing – Arithmetic expression

- ▶ Parse an arithmetic expression and return abstract syntax tree.

```
1 data Exp = Plus Exp Exp
2           | Minus Exp Exp
3           | Times Exp Exp
4           | Div Exp Exp
5           | Const Integer deriving Show
6
7 expr = term `chainl1` addop -- expr ::= term {addop term}
8 term = fact `chainl1` mulop -- term ::= fact {mulop fact}
9 -- fact ::= '(' expr ')' | integer
10 fact = parens expr <|> (Const <$> integer)
11
12 addop = (char '+' >> return Plus) -- addop ::= '+' | '-'
13     <|> (char '-' >> return Minus)
14 mulop = (char '*' >> return Times) -- mulop ::= '*' | '/'
15     <|> (char '/' >> return Div)
16
17 parens p = do { char '('; x <- p; char ')'; return x }
18 integer = read <$> many1 digit -- integer ::= digit {digit}
19
20 digit :: Parser Char -- digit ::= [0--9]
21 -- chainl1 :: Parser t -> Parser (t -> t -> t) -> Parser t
```

## Example of parsing – Arithmetic expression

- ▶ Parse an arithmetic expression and return abstract syntax tree.

```
1 expr = term `chainl1` addop -- expr ::= term {addop term}
2 term = fact `chainl1` mulop -- term ::= fact {mulop fact}
3 -- fact ::= '(' expr ')' | integer
4 fact = parens expr <|> (Const <$> integer)
5
6 addop = (char '+' >> return Plus) -- addop ::= '+' | '-'
7     <|> (char '-' >> return Minus)
8 mulop = (char '*' >> return Times) -- mulop ::= '*' | '/'
9     <|> (char '/' >> return Div)
10
11 parens p = do { char '('; x <- p; char ')'; return x }
12 integer = read <$> many1 digit -- integer ::= digit {digit}
13
14 digit :: Parser Char -- digit ::= [0--9]
```

- ▶ *chainl1 p op* folds *p* one or more times with *op* as an left associative operator.

```
1 chainl1 p op = p >>= rest
2   where rest x = do { f <- op; y <- p; rest $ f x y }
3           <|> return x
4
5 chainl p op x = chainl1 p op <|> return x
```

# Example of parsing – Arithmetic expression

- ▶ Parse an arithmetic expression and return abstract syntax tree.

```
1 expr = term `chainl1` addop -- expr ::= term {addop term}
2 term = fact `chainl1` mulop -- term ::= fact {mulop fact}
3 -- fact ::= '(' expr ')' | integer
4 fact = parens expr <|> (Const <$> integer)
5
6 addop = (char '+' >> return Plus) -- addop ::= '+' | '-'
7     <|> (char '-' >> return Minus)
8 mulop = (char '*' >> return Times) -- mulop ::= '*' | '/'
9     <|> (char '/' >> return Div)
10
11 parens p = do { char '('; x <- p; char ')'; return x }
12 integer = read <$> many1 digit -- integer ::= digit {digit}
13
14 digit :: Parser Char -- digit ::= [0--9]
```

- ▶ Test it

```
1 run expr "1+(2-3)*4"
2 --Plus (Const 1) (Times (Minus (Const 2) (Const 3)) (Const 4))
```

# Example of parsing – Arithmetic expression



- ▶ Parse an arithmetic expression and evaluate it

```
1 expr = term `chainl1` addop -- expr ::= term {addop term}
2 term = fact `chainl1` mulop -- term ::= fact {mulop fact}
3 -- fact ::= '(' expr ')' | integer
4 fact = parens expr <|> integer
5
6 addop = (char '+' >> return (+)) -- addop ::= '+' | '-'
7     <|> (char '-' >> return (-))
8 mulop = (char '*' >> return (*)) -- mulop ::= '*' | '/'
9     <|> (char '/' >> return div)
10
11 parens p = do { char '('; x <- p; char ')'; return x }
12 integer = read <$> many1 digit -- integer ::= digit {digit}
13
14 digit :: Parser Char -- digit ::= [0--9]
```

- ▶ Test it

```
1 run expr "1+(2-3)*4"
2 -- evaluates -3
```

# Parsec



- ▶ Example adopted from Chapter 16 of *Real World Haskell*.

```
1 import Text.Parsec
2
3 type Parser = Parsec [Char] ()
4
5 csvFile :: Parser [[String]]
6 csvFile = endBy line newline
7
8 line = sepBy cell (char ',')
9 cell = many (noneOf ",\n\r")
10 noneOf :: [Char] -> Parser Char
```

- ▶ Runner

```
1 run :: Show a => Parser a -> String -> IO()
2 run p input = putStrLn $ case parse p "" input of
3     Left error -> show error
4     Right x      -> show x
5
6 -- parse parser source-name input -> either error or result
7 parse :: Parser a -> String -> String -> Either ParseError a
8
9 run csvFile "10,one\n20,two\n"
10 -- [[10,"one"],[20,"two"]]
```