

CMSC330: Parsing and Lexing

Chris Kauffman

*Last Updated:
Wed Oct 18 09:25:27 PM EDT 2023*

Logistics

Assignments

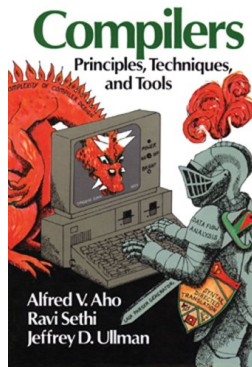
- ▶ Project 5 is almost up
- ▶ NFA to DFA conversion in OCaml
- ▶ P5 due 17-Oct
- ▶ Exam 1 Grades will be posted Wed

Reading

Chapter 4 “Parsing” from *Compilers: Principles, Techniques, and Tools* by Aho et al.
(e.g. the “Dragon Book”)

Goals

Lexing and Parsing input into Trees



ACM named Alfred Vaino Aho and Jeffrey David Ullman recipients of the 2020 ACM A.M. Turing Award for fundamental algorithms and theory underlying programming language implementation and for synthesizing these results and those of others in their highly influential books. . .

– Source: 2020 ACM A.M. Turing Award Laureates

An Opening Example

```
let x = 5 + 10*4 + 7*(3+2) in ...
```

- ▶ Above program fragment is something one would type in and expect compiled/run without trouble
- ▶ The trouble comes when writing the interpreter or compiler
- ▶ Will discuss issues associated with these for the next few sessions
- ▶ Ultimate goal: code to evaluate arithmetic expressions
- ▶ Will require
 - ▶ **Algebraic Types** for symbolic data
 - ▶ **Recursive functions** for lexing/parsing
 - ▶ **List processing** for parsing
 - ▶ **Recursive functions** on trees for evaluation

Good thing that you've mastered all these...

Roadmap for Processing

1. Get input

```
let input = "5 + 10*4 + 7*(3+2)";;
```

2. Lex input to tokens

```
[Int 5; Plus; Int 10; Times; Int 4; Plus; ...
```

3. Parse tokens to expression tree

```
Add(Const(5),  
      Add(Mul(Const(10),  
               Const(4))), ...
```

4. Evaluate tree

80

Getting Input

- ▶ Generally the easy part
- ▶ Obtained via reading from a file or user typing in input
- ▶ Can usually assume it's stored in a string somewhere like

```
let input = "5 + 10*4 + 7*(3+2)";;
```
- ▶ First step in compiler or interpreter is to get input like this

Lexing

lexical analysis \equiv *lexing* \equiv *tokenization*

- ▶ Raw input is just a bunch of characters
- ▶ Lexing is done to ease processing later on
- ▶ Group characters into **tokens** / **lexemes** for operations, numbers, keywords, etc
- ▶ Assign some meaning to tokens via a symbolic name
- ▶ Identify characters that don't belong/not recognized
- ▶ Output of lexing is a stream of such tokens, will use a list of tokens in our work

```
let input = "5 + 10*4 + 7*(3+2)";; (* Lexing: convert this string.*)
let lexed = [Int 5; Plus; Int 10; (* Into this stream of tokens *)
            Times; Int 4; Plus;
            Int 7; Times;
            OParen; Int 3; Plus;
            Int 2; CParen];;
```

Lexing and Token Symbols

- ▶ Typical tokens are **symbolic** and may carry additional data
- ▶ A convenient way to represent them in OCaml is via algebraic types with variants for each token type

```
(* algebraic types for tokens: lexing results *)
type token =
    Plus | Times | OParen | CParen | Int of int;;
```

- ▶ Above tokens are very limited but sufficient for simple arithmetic
- ▶ More extensive arithmetic language processing would include subtraction, division, floating point numbers
- ▶ Fuller programming languages have variable identifiers, keywords like `let/in` and `for/do`

Exercise: Lexing Thought Questions

```
let input = "5 + 10*4 + 7*(3+2)";; (* Lexing: convert this string..*)
let lexed = [Int 5; Plus; Int 10;   (* Into this stream of tokens *)
             Times; Int 4; Plus;
             Int 7; Times;
             OParen; Int 3; Plus;
             Int 2; CParen];;
```

1. Do all characters in the above string appear in the tokens? If not, which are not present and why?
2. Do all the tokens correspond to single characters or are groups of multiple characters associated with a single token?
3. Speculate on how the code for a lexing function will look

Answers: Lexing Thought Questions

```
let input = "5 + 10*4 + 7*(3+2)";; (* Lexing: convert this string.. *)
let lexed = [Int 5; Plus; Int 10; (* Into this stream of tokens *)
             Times; Int 4; Plus;
             Int 7; Times;
             OParen; Int 3; Plus;
             Int 2; CParen];;
```

1. Do all characters in the above string appear in the tokens? If not, which are not present and why?
Spaces do not get tokenized; whitespace is commonly ignored.
2. Do all the tokens correspond to single characters or are groups of multiple characters associated with a single token?
In the above example, most tokens are single characters but multi-character numbers like 10 are a single token.
3. Speculate on how the code for a lexing function will look
Lexing functions must scan through the string matching characters and creating tokens. Our version will combine recursion and iteration.

A Simple Lexing Routine

```
1 let lex_string string = (* create a list of tokens *)
2   let len = String.length string in
3   let rec lex pos = (* recursive helper on pos in string *)
4     if pos >= len then (* off end of string ? *)
5       [] (* end of input *)
6     else (* more to lex *)
7       match string.[pos] with (* match a single character *)
8       | ' ' | '\t' | '\n' -> lex (pos+1) (* skip whitespace *)
9       | '+' -> Plus :: (lex (pos+1)) (* single char ops become operators *)
10      | '*' -> Times :: (lex (pos+1)) (* like add and multiply *)
11      | '(' -> OParen :: (lex (pos+1)) (* and open/close parens *)
12      | ')' -> CParen :: (lex (pos+1))
13      | d when is_digit d -> (* see a digit *)
14        let stop = ref pos in (* scan through until a non-digit is found *)
15        while !stop < len && is_digit string.[!stop] do
16          incr stop;
17        done;
18        let numstr = String.sub string pos (!stop - pos) in (* substring is the int *)
19        let num = int_of_string numstr in (* parse the integer *)
20        Int(num) :: (lex !stop) (* and tack onto the stream of tokens *)
21      | _ -> (* any other characters lead to failures *)
22        let msg = sprintf "lex error at char %d, char '%c'" pos string.[pos] in
23        failwith msg
24  in (* end helper *)
25  lex 0 (* call helper *)
26 ;;
```

Sample Calls to `lex_string`

```
1 # lex_string "123";;
2 - : token list = [Int 123]
3
4 # lex_string "*";;                (* NOT VALID SYNTAX BUT THAT IS *)
5 - : token list = [Times]          (* THE BUSINESS OF THE PARSING PHASE *)
6
7 # lex_string "123 +";;
8 - : token list = [Int 123; Plus]
9
10 # lex_string "123 + 19";;
11 - : token list = [Int 123; Plus; Int 19]
12
13 # lex_string "123 + (19)";;
14 - : token list = [Int 123; Plus; OParen; Int 19; CParen]
15
16 # lex_string "123 * (1+19)";;
17 - : token list = [Int 123; Times; OParen; Int 1; Plus; Int 19; CParen]
18
19 # lex_string "/";;                (* UNRECOGNIZED CHARACTERS *)
20 Exception: Failure "lex error at char 0, char '/'".
21
22 # lex_string "123 + 19 - 9";;
23 Exception: Failure "lex error at char 9, char '-'".
```

Notes on Lexing

- ▶ The provided Lexing routine does raw character processing
- ▶ Might benefit from use of **regular expression** matching to make things easier
- ▶ Often tool chains that automatically create lexer code specify tokens via regular expressions

Context Free Grammars and Parsing

- ▶ Have seen that CFGs allow formal specification of the order / syntax of an allowed language
- ▶ Have seen how to derive terminal strings / parse trees from CFGs
- ▶ Parsers solve the inverse problem
 - ▶ Given a terminal string . . .
 - ▶ Determine if it is allowed by the grammar
 - ▶ Construct a parse tree for it
 - ▶ If string cannot be generated by the grammar, reject it with syntax errors
- ▶ Like most **inverse problems**, going backwards (string to grammar) is more difficult than going forwards (grammar to string)

Parsing Approaches

- ▶ Parsing approaches are widely studied and have many flavors
 - ▶ Top Down vs Bottom Up
 - ▶ Generate Leftmost derivation vs Rightmost derivation
 - ▶ Deterministic based on finite look-ahead, sometimes not
 - ▶ Recursive descent vs Pushdown Automata vs Parsing Tables
- ▶ In literature will see acronyms summarizing some aspects of these
 - ▶ LL(1): Left to right parsing, Leftmost derivation, 1-token look-ahead
 - ▶ LALR(1): Left to right parsing, Rightmost derivation, 1-token look-ahead
- ▶ In most practical scenarios, you **shouldn't write your own parser**: parsing is tricky and automated tools exist to aid with creation of fast efficient parses
- ▶ In **educational settings, writing it from scratch helps us learn**
- ▶ Will focus on the traditional Recursive Descent Parser

Goal: Build an Abstract Syntax Tree (AST)

- ▶ `langproc.ml` builds an abstract syntax tree for arithmetic with integers, `+` and `*`
- ▶ Abstract syntax tree for these is comprised of the following

```
1  (* algebraic types for expression tree: parsing results *)
2  type expr =
3      | Add of expr * expr
4      | Mul of expr * expr
5      | Const of int
6      ;;
7
8  let input = "5 + 10*4 + 7*(3+2)";;
9  let parsed =
10     Add(Const(5),
11         Add(Mul(Const(10),
12                 Const(4)),
13             Mul(Const(7),
14                 Add(Const(3),
15                     Const(2))))))
16     ;;
```

CFGs and Recursive Descent Parsers

- ▶ Recall CFGs have Terminal Symbols, Non-terminal Symbols, Production Rules, a Start Symbol
- ▶ A **recursive descent parser** handles parsing via a series of functions, possibly recursive
- ▶ Typically **one function per non-terminal symbol**
- ▶ Initial function to call is the start symbol
- ▶ Functions look for non-terminal symbols in input and attempt to consume them according to production rules in the grammar
- ▶ Recursive Descent has the following properties
 - ▶ Top-Down Parsing: constructs upper parts of Parse Tree / AST building to lower
 - ▶ Non-deterministic: involves search and backtracking
 - ▶ Look-Ahead: will look forward at 1 or more tokens in input to decide how to proceed

Example CFG Recursive Descent Parsing Structure

ParenMath

$$A \rightarrow A + A \quad (1)$$

$$A \rightarrow M \quad (2)$$

$$M \rightarrow M * M \quad (3)$$

$$M \rightarrow N \quad (4)$$

$$N \rightarrow \text{number} \quad (5)$$

$$N \rightarrow (A) \quad (6)$$

Ex: $5 + 10 * 4 + 7 * (3 + 2)$

Recursive Descent Parsers often involve mutually recursive functions so use OCaml's `and` syntax to define a series of such functions

Parser for ParenMath

```
let parse_tokens tokens =  
    (* A prec2: addition only *)  
    let rec prec2 toks =  
        ...  
  
    (* M prec1: multiplication *)  
    and prec1 toks =  
        ...  
  
    (* N prec0: self-evaluating tokens  
       like Int and Paren expressions *)  
    and prec0 toks =  
        ...  
  
    in  
  
    (* end helpers, main code for  
       parse_tokens *)  
    let (expr, rest) = prec2 tokens in  
        ...  
;;
```

Exercise: Highest Precedence Grammar Elements

- ▶ langproc.ml provides a recursive descent parser with 3 precedence levels
- ▶ Below is highest precedence level

```
1 (* prec0: self-evaluating tokens like Int and parenthesized expressions *)
2 let rec prec0 toks =
3   match toks with
4   | [] -> (* out of input *)
5     raise (ParseError {msg="expected an expression"; toks=toks})
6   | Int n :: tail -> (* ints are self-evaluating *)
7     (Const(n),tail)
8   | OParen :: tail -> (* parenthesized expression *)
9     begin
10      let (expr,rest) = prec2 tail in (* start back at lowest precedence *)
11      match rest with
12      | CParen::tail -> (expr,tail)
13      | _ -> raise (ParseError {msg="unclosed parentheses"; toks=rest})
14    end
15   | _ ->
16     raise (ParseError {msg="syntax error"; toks=toks})
```

1. What kind of thing is parameter toks?
2. What types of tokens are handled by prec0
3. What kind of parsing errors can result at this level?

Answers: Highest Precedence Grammar Elements

1. What kind of thing is parameter toks?
It is a list of token types. The head element is analyzed in the match/with.
2. What types of tokens are handled by prec0
Int tokens which are converted to Const expressions and OParen/CParen tokens which continue parsing again.
3. What kind of parsing errors can result at this level?
Running out of input, failure to close a parenthesis, and general syntax errors.

Exercise: Arithmetic Operations

```
1 (* prec1: multiplication *)
2 and prec1 toks =
3   let (lexpr, rest) = prec0 toks in   (* try higher prec first *)
4   match rest with
5   | Times :: tail ->                  (* * is first *)
6     let (rexpr, rest) = prec1 tail in (* recurse to get right-hand expr *)
7     (Mul(lexpr, rexpr), rest)         (* multiply left/right expr *)
8   | _ -> (lexpr, rest)                (* not a multiply *)
9
10 (* prec2: addition only *)
11 and prec2 toks =
12   let (lexpr, rest) = prec1 toks in   (* try higher prec first *)
13   match rest with
14   | Plus :: tail ->                   (* + is first *)
15     let (rexpr, rest) = prec2 tail in (* recurse to get right-hand expr *)
16     (Add(lexpr, rexpr), rest)         (* add left/right expr *)
17   | _ -> (lexpr, rest)                (* not an addition *)
```

1. prec1 calls prec0 and prec2 calls prec1. In this scheme, does a large number indicate high or low precedence?
2. When prec0 is called, what is the result? Is this the same for when prec1 is called? how about prec2?
3. What happens if prec1 and prec2 cannot match a token at the beginning of the list of tokens?

Answers: Arithmetic Operations

1. `prec1` calls `prec0` and `prec2` calls `prec1`. In this scheme, does a large number indicate high or low precedence?
Large numbers indicate lower precedence as `prec0` are self-evaluating or parenthesized expressions. This is the opposite of the OCaml documentation where higher numbers indicate higher precedence.
2. When `prec0` is called, what is the result? Is this the same for when `prec1` is called? how about `prec2`?
A pair of an expression (AST) and the remaining list of tokens is returned in both cases. This is also what `prec0` returns.
3. What happens if `prec1` and `prec2` cannot match a token at the beginning of the list of tokens?
They return the expression given by the `prec0` or `prec1` paired with the entire list of tokens with nothing removed. This backtracks to a previous function.

The Full(-ish) Parser

3 precedence levels, 3 functions to handle them

```
1 (* parse tokens list of arithmetic
2   language to an AST *)
3 let parse_tokens tokens =
4
5   (* prec0: self-evaluating tokens like
6     Int and parenthesized expressions *)
7   let rec prec0 toks =
8     match toks with
9     | [] ->
10       raise (Expect Expression)
11     | Int n :: tail ->
12       (Const(n),tail)
13     | OParen :: tail ->
14       begin
15         let (expr,rest) = prec2 tail in
16         match rest with
17         | CParen::tail -> (expr,tail)
18         | _ -> raise (Unclosed Paren)
19       end
20     | _ ->
21       raise (Syntax Error)
22
23   (* prec1: multiplication *)
24   and prec1 toks =
25     let (lexpr, rest) = prec0 toks in
26     match rest with
27     | Times :: tail ->
28       let (rexpr,rest) = prec1 tail in
29       (Mul(lexpr,rexpr), rest)
30     | _ -> (lexpr, rest)
31
32   (* prec2: addition only *)
33   and prec2 toks =
34     let (lexpr, rest) = prec1 toks in
35     match rest with
36     | Plus :: tail ->
37       let (rexpr,rest) = prec2 tail in
38       (Add(lexpr,rexpr), rest)
39     | _ -> (lexpr, rest)
40
41   in
42
43   (* end helpers, main code for
44     parse_tokens *)
45   let (expr, rest) = prec2 tokens in
46   match rest with
47   | [] -> expr
48   | _ -> raise (Tokens Remain)
49 ;;
```

Sample Calls to parse_tokens

```
1 # parse_tokens (lex_string "11+5");;          (* Simple tree with Const leaves *)
2 - : expr = Add (Const 11,
3               Const 5)
4
5 # parse_tokens (lex_string "11+5+2+9");;      (* Chain-esque tree of repeated Add *)
6 - : expr = Add (Const 11,
7               Add (Const 5,
8                   Add (Const 2,
9                       Const 9)))
9
10
11 # parse_tokens (lex_string "11+5*2+9");;      (* Mult has higher precedence *)
12 - : expr = Add (Const 11,
13               Add (Mul (Const 5,
14                       Const 2),
15                   Const 9))
15
16
17 # parse_tokens (lex_string "2*(11+5)*2");;    (* Parens change precedence *)
18 - : expr = Mul (Const 2,
19               Mul (Add (Const 11,
20                       Const 5),
21                   Const 2))
21
22
23 # parse_tokens (lex_string "2 +");;           (* Parse Error Cases *)
24 Exception: ParseError {msg = "expected an expression"; toks = []}.
25
26 # parse_tokens (lex_string "2 * (3+5 " );;
27 Exception: ParseError {msg = "unclosed parentheses"; toks = []}.
28
29 # parse_tokens (lex_string "2 * 3 + * 5");;
30 Exception: ParseError {msg = "syntax error"; toks = [Times; Int 5]}.
```

Exercise: Show Abstract Syntax

- Show the AST or report errors for the following examples
- First few examples are already done for reference

```
# parse_tokens (lex_string "11+5+2+9");;      (* EXAMPLE A *)  
- : expr = Add (Const 11,  
               Add (Const 5,  
                   Add (Const 2,  
                       Const 9)))
```

```
# parse_tokens (lex_string "2*(11+5)*2");;      (* EXAMPLE B *)  
- : expr = Mul (Const 2,  
               Mul (Add (Const 11,  
                       Const 5),  
                   Const 2))
```

```
# parse_tokens (lex_string "(11 + 2)*5");;      (* PROBLEM 1 *)
```

```
# parse_tokens (lex_string "(11 + )*5");;      (* PROBLEM 2 *)
```

```
# parse_tokens (lex_string "11*5*2*9");;      (* PROBLEM 3 *)
```

```
# parse_tokens (lex_string "(11+2)+(5+9)");;    (* PROBLEM 4 *)
```

```
# parse_tokens (lex_string "11*5*2+9");;      (* PROBLEM 5 *)
```


Answers: Show Parse Tree

```
1 # parse_tokens (lex_string "(11 + 2)*5");;      (* PROBLEM 1 *)
2 - : expr = Mul (Add (Const 11,                  (* parens give addition higher *)
3                  Const 2),                      (* precedence, later multiply *)
4                  Const 5)
5
6 # parse_tokens (lex_string "(11 + )*5");;      (* PROBLEM 2 *)
7 Exception: ParseError {msg = "syntax error"; toks = [CParen; Times; Int 5]}.
8
9 # parse_tokens (lex_string "11*5*2*9");;      (* PROBLEM 3 *)
10 - : expr = Mul (Const 11,                      (* chain-like tree of repeated *)
11               Mul (Const 5,                  (* multiplications *)
12                   Mul (Const 2,
13                       Const 9)))
14
15 # parse_tokens (lex_string "(11+2)+(5+9)");;  (* PROBLEM 4 *)
16 - : expr = Add (Add (Const 11,                (* Parens give first/last add *)
17                   Const 2),                  (* higher precedence than middle *)
18               Add (Const 5,
19                   Const 9))
20
21 # parse_tokens (lex_string "11*5*2+9");;      (* PROBLEM 5 *)
22 - : expr = Add (Mul (Const 11,                (* Multiply everything first *)
23                 Mul (Const 5,                (* then add *)
24                     Const 2)),
25                 Const 9)
```

These make for easy exam problems...

Recursive Descent Parsing is a Search Process

- ▶ It is difficult but very worthwhile to reason about how the recursive descent parser works in total
- ▶ Parsing is really a **search process** in which each function tries to consume tokens with some failures allowing backtracking
- ▶ Gets even trickier to understand with parentheses involved which circle back to top of parsing functions

TOKEN POSITION	FUNCTION CALL STACK	EXPRESSON TREE
5 + 8 * 4		
<5> + 8 * 4	prec2	
<5> + 8 * 4	prec2,prec1	
<5> + 8 * 4	prec2,prec1,prec0	
5 <+> 8 * 4	prec2,prec1	Const(5)
5 <+> 8 * 4	prec2	Const(5)
5 + <8> * 4	prec2,prec2	Add(Const(5),..)
5 + <8> * 4	prec2,prec2,prec1	Add(Const(5),..)
5 + <8> * 4	prec2,prec2,prec1,prec0	Add(Const(5),..)
5 + 8 <*> 4	prec2,prec2,prec1	Add(Const(5),..) Const(8)
5 + 8 * <4>	prec2,prec2,prec1,prec1	Add(Const(5), Mul(Const(8), ...))
5 + 8 * <4>	prec2,prec2,prec1,prec1,prec0	Add(Const(5), Mul(Const(8), ...))
5 + 8 * 4<>	prec2,prec2,prec1,prec1,prec0	Add(Const(5), Mul(Const(8), Const(4)))

Aside: Tracing Function Execution

- ▶ A **call trace** conveys a sequence of function calls with parameter and return values
- ▶ Useful to understand code flow within a group of functions, now the stack looks
- ▶ Like many programming environments with a REPL, OCaml can **automatically generate call traces**
- ▶ `#trace funcname;;` turns on tracing of function, shows calls with params and return values
- ▶ Look for trace features in every language either directly or via debugger tools

```
# let doub x = 2 * x;;
val doub : int -> int = <fun>

# let octosum (x,y) =
  (doub (doub x)) + (doub (doub y));;
val octosum : int * int -> int = <fun>

# #trace octosum;;
octosum is now traced.
# #trace doub;;
doub is now traced.

# octosum (2,7);;
octosum <-- (2, 7)
double <-- 7
double --> 14
double <-- 14
double --> 28
double <-- 2
double --> 4
double <-- 4
double --> 8
octosum --> 36
- : int = 36
```

Tracing Parsing Functions

```
1 #use "langproc_trace.ml";;          (* defines parsing funcs at top level *)
2 ...                                (* rather than inside parse_tokens *)
3 val prec2 : token list -> token list * expr = <fun>
4 val prec1 : token list -> token list * expr = <fun>
5 val prec0 : token list -> token list * expr = <fun>
6 ...
7 # #trace prec2;;                    (* trace each of the parsing funcs *)
8 prec2 is now traced.
9 # #trace prec1;;
10 prec1 is now traced.
11 # #trace prec0;;
12 prec0 is now traced.
13
14 # let parsed = parse_tokens (lex_string "7");;
15 prec2 <-- [Int 7]
16 prec1 <-- [Int 7]
17 prec0 <-- [Int 7]
18 prec0 --> ([], Const 7)
19 prec1 --> ([], Const 7)
20 prec2 --> ([], Const 7)
21 val parsed : expr = Const 7
```

A More Intense Trace Example

- ▶ Indented (by hand) to show matching of call/return
- ▶ Token list param/return on the left and expr return on right

```
1 # parse_tokens (lex_string "7*(3+2)");;
2 prec2 <-- [Int 7; Times; OParen; Int 3; Plus; Int 2; CParen]
3   prec1 <-- [Int 7; Times; OParen; Int 3; Plus; Int 2; CParen]
4     prec0 <-- [Int 7; Times; OParen; Int 3; Plus; Int 2; CParen]
5       prec0 --> ([Times; OParen; Int 3; Plus; Int 2; CParen], Const 7)
6         prec1 <-- [OParen; Int 3; Plus; Int 2; CParen]
7           prec0 <-- [OParen; Int 3; Plus; Int 2; CParen]
8             prec2 <-- [Int 3; Plus; Int 2; CParen]
9               prec1 <-- [Int 3; Plus; Int 2; CParen]
10                 prec0 <-- [Int 3; Plus; Int 2; CParen]
11                   prec0 --> ([Plus; Int 2; CParen], Const 3)
12                     prec1 --> ([Plus; Int 2; CParen], Const 3)
13                       prec2 <-- [Int 2; CParen]
14                         prec1 <-- [Int 2; CParen]
15                           prec0 <-- [Int 2; CParen]
16                             prec0 --> ([CParen], Const 2)
17                               prec1 --> ([CParen], Const 2)
18                                 prec2 --> ([CParen], Const 2)
19                                   prec2 --> ([CParen], Add (Const 3, Const 2))
20                                     prec0 --> ([], Add (Const 3, Const 2))
21                                       prec1 --> ([], Add (Const 3, Const 2))
22                                         prec1 --> ([], Mul (Const 7, Add (Const 3, Const 2)))
23                                           prec2 --> ([], Mul (Const 7, Add (Const 3, Const 2)))
24 - : expr = Mul (Const 7, Add (Const 3, Const 2))
```

Adding more to the Language

- ▶ A good exercise would be to grab the `langproc.ml` file from the codepack and add obvious stuff to the arithmetic expression language like subtraction, division, variable IDs
- ▶ More interesting stuff is programmatic
 - ▶ Name/value binding (including functions)
 - ▶ Control structures like `if/then/else`
 - ▶ Function application and Lambda Expressions
- ▶ All have the same flavor
 - ▶ Add a `prec_ifthenelse` toks function to parse family
 - ▶ Looks for a pattern of tokens that fits
`if <expr> then <expr> else <expr>`
 - ▶ Creates an associated parse tree associated with these like
`Cond(ifexpr, thenexpr, elseexpr)`

But first, there's the small matter of **evaluating** a parsed expression in the simple case of arithmetic

Exercise: Evaluating an AST

Consider example parse tree below

```
# parse_tokens (lex_string "11*5*2+9");;  
- : expr = Add (Mul (Const 11,  
                    Mul (Const 5,  
                        Const 2))),  
            Const 9)
```

1. How does one evaluate such arithmetic expressions by hand?
2. What features do you expect from code that evaluate's the expression tree and produces an integer answer?
3. Write a version of evaluate

```
val evaluate : expr -> int = <fun>
```

Answers: Evaluating an AST

1. How does one evaluate such arithmetic expressions by hand?
Usually ad-hoc but definitely multiply first then add
2. What features would code that would evaluate the expression tree and produce an integer answer?
It will operate on the Abstract Syntax Tree so will likely be recursive. The tree is comprised of algebraic types so pattern matching is expected.
3. Write a version of evaluate

```
1  (* Evaluate AST of expressions to produce an integer result *)
2  let rec evaluate expr =
3      match expr with
4      | Const i -> i
5      | Add(lexpr,rexpr) ->
6          let lans = evaluate lexpr in
7          let rans = evaluate rexr in
8          lans + rans
9      | Mul(lexpr,rexpr) ->
10         let lans = evaluate lexpr in
11         let rans = evaluate rexr in
12         lans * rans
13
14  # evaluate (parse_tokens (lex_string "11*5*2+9"));
15  - : int = 119
```


Aside: OCaml's “Threading” Operator |>

- ▶ At times have a series of functions to apply to data, tedious and backwards looking to write

```
let ans = f3 (f2 (f1 x)) in ...
```

- ▶ A reverse application operator is helpful, referred to in some contexts as a “threading” operator¹ as it threads data through functions in the order they appear
- ▶ OCaml has threading operator called |>

```
let ans = x |> f1 |> f2 |> f3 ...
```

```
(* full cycle of lex, parse, evaluate *)
```

```
# let result = evaluate (parse_tokens (lex_string "2*3+5"));;
```

```
val result : int = 11
```

```
(* with |> threading operator (a.k.a reverse application) *)
```

```
# let result = "2*3+5" |> lex_string |> parse_tokens |> evaluate;;
```

```
val result : int = 11
```

¹This kind of threading has nothing to do with multi-threaded programs; it merely makes certain kinds of function applications easier on the eyes.

Parsing and Evaluation in the General Case

```
# "let x=5 in \
  let z=x+2 in \
  if x > z then \
  then  x*z    \
  else  x+z;"
|> lexcalc_string |> parsecalc_tokens;;
```

Parse tree:

```
Letin( x )
  .var_expr:
    IntExp(5)
  .in_expr:
    Letin( z )
      .var_expr:
        Add
          Varname(x)
          IntExp(2)
      .in_expr:
        Cond
          .if_expr:
            Greater
              Varname(x)
              Varname(z)
          .then_expr:
            Mul
              Varname(x)
              Varname(z)
          .else_expr:
            Add
              Varname(x)
              Varname(z)
```

- ▶ Beyond simple arithmetic, language processing systems add more programmatic features
- ▶ Parsing usually puts these in an AST initially but may convert to other data structures depending on the goals
- ▶ Informally, evaluation “does something” with these elements during a tree walk such as “executing” the tree, producing assembly or virtual machine instructions, drawing something, etc.
- ▶ Formalisms such as **operational semantics** exist which emphasize a mathematically rigorous description of what should be done on encountering each syntactic element of the language