

## Programming Languages and Compilers (CS 421)

Elsa L Gunter  
2112 SC, UIUC

<http://courses.engr.illinois.edu/cs421>

Based in part on slides by Mattox Beckman, as updated by Vikram Adve and Gul Agha

10/23/14

1

## General Input

```
{ header }  
let ident = regexp ...  
rule entrypoint [arg1... argn] = parse  
    regexp { action }  
    | ...  
    | regexp { action }  
and entrypoint [arg1... argn] =  
    parse ...and ...  
{ trailer }
```

10/23/14

2

## Ocamllex Input

- *header* and *trailer* contain arbitrary ocaml code put at top and bottom of `<filename>.ml`
- `let ident = regexp ...` Introduces *ident* for use in later regular expressions

10/23/14

3

## Ocamllex Input

- `<filename>.ml` contains one lexing function per *entrypoint*
  - Name of function is name given for *entrypoint*
  - Each entry point becomes an Ocaml function that takes  $n+1$  arguments, the extra implicit last argument being of type `Lexing.lexbuf`
- `arg1... argn` are for use in *action*

10/23/14

4

## Ocamllex Regular Expression

- Single quoted characters for letters: `'a'`
- `_`: (underscore) matches any letter
- `Eof`: special "end\_of\_file" marker
- Concatenation same as usual
- `"string"`: concatenation of sequence of characters
- `e1 / e2`: choice - what was  $e_1 \vee e_2$

10/23/14

5

## Ocamllex Regular Expression

- `[c1 - c2]`: choice of any character between first and second inclusive, as determined by character codes
- `[^c1 - c2]`: choice of any character NOT in set
- `e*`: same as before
- `e+`: same as `e e*`
- `e?`: option - was  $e_1 \vee \epsilon$

10/23/14

6

## Ocamllex Regular Expression

- $e_1 \# e_2$ : the characters in  $e_1$  but not in  $e_2$ ;  $e_1$  and  $e_2$  must describe just sets of characters
- *ident*: abbreviation for earlier reg exp in `let ident = regexp`
- *e<sub>1</sub> as id*: binds the result of  $e_1$  to *id* to be used in the associated *action*

10/23/14

7

## Ocamllex Manual

- More details can be found at

<http://caml.inria.fr/pub/docs/manual-ocaml/manual026.html>

10/23/14

8

## Example : test.mll

```
{ type result = Int of int | Float of float |
  String of string }
let digit = ['0'-'9']
let digits = digit +
let lower_case = ['a'-'z']
let upper_case = ['A'-'Z']
let letter = upper_case | lower_case
let letters = letter +
```

10/23/14

9

## Example : test.mll

```
rule main = parse
  (digits)'.digits as f { Float (float_of_string f) }
  | digits as n          { Int (int_of_string n) }
  | letters as s         { String s}
  | _ { main lexbuf }
{ let newlexbuf = (Lexing.from_channel stdin) in
  print_string "Ready to lex.";
  print_newline ();
  main newlexbuf }
```

10/23/14

10

## Example

```
# #use "test.mll";;
...
val main : Lexing.lexbuf -> result = <fun>
val __ocaml_lex_main_rec : Lexing.lexbuf -> int ->
  result = <fun>
Ready to lex.
hi there 234 5.2
- : result = String "hi"
What happened to the rest?!?
```

10/23/14

11

## Example

```
# let b = Lexing.from_channel stdin;;
# main b;;
hi 673 there
- : result = String "hi"
# main b;;
- : result = Int 673
# main b;;
- : result = String "there"
```

10/23/14

12

## Your Turn

- Work on MP8
  - Add a few keywords
  - Implement booleans and unit
  - Implement Ints and Floats
  - Implement identifiers

10/23/14

13

## Problem

- How to get lexer to look at more than the first token at one time?
- One Answer: *action* tells it to -- recursive calls
- Side Benefit: can add “state” into lexing
- Note: already used this with the `_` case

10/23/14

14

## Example

```
rule main = parse
  (digits) '!' digits as f { Float
  (float_of_string f) :: main lexbuf }
  | digits as n      { Int (int_of_string n) ::
  main lexbuf }
  | letters as s     { String s :: main lexbuf }
  | eof              { [] }
  | _                { main lexbuf }
```

10/23/14

15

## Example Results

Ready to lex.

hi there 234 5.2

- : result list = [String "hi"; String "there"; Int 234; Float 5.2]

#

Used Ctrl-d to send the end-of-file signal

10/23/14

16

## Dealing with comments

### First Attempt

```
let open_comment = "(*"
let close_comment = "*)"
rule main = parse
  (digits) '!' digits as f { Float (float_of_string
  f) :: main lexbuf }
  | digits as n      { Int (int_of_string n) ::
  main lexbuf }
  | letters as s     { String s :: main lexbuf }
```

10/23/14

17

## Dealing with comments

```
| open_comment      { comment lexbuf }
| eof                { [] }
| _ { main lexbuf }
and comment = parse
  close_comment     { main lexbuf }
  | _                { comment lexbuf }
```

10/23/14

18

## Dealing with nested comments

```
rule main = parse ...
| open_comment { comment 1 lexbuf }
| eof { [] }
| _ { main lexbuf }
and comment depth = parse
open_comment { comment (depth+1)
lexbuf }
| close_comment { if depth = 1
then main lexbuf
else comment (depth - 1) lexbuf }
| _ { comment depth lexbuf }
```

10/23/14

19

## Dealing with nested comments

```
rule main = parse
(digits) '.' digits as f { Float (float_of_string f) ::
main lexbuf }
| digits as n { Int (int_of_string n) :: main
lexbuf }
| letters as s { String s :: main lexbuf }
| open_comment { (comment 1 lexbuf) }
| eof { [] }
| _ { main lexbuf }
```

10/23/14

20

## Dealing with nested comments

```
and comment depth = parse
open_comment { comment (depth+1) lexbuf }
| close_comment { if depth = 1
then main lexbuf
else comment (depth - 1) lexbuf }
| _ { comment depth lexbuf }
```

10/23/14

21

## Types of Formal Language Descriptions

- Regular expressions, regular grammars
- Context-free grammars, BNF grammars, syntax diagrams
- Finite state automata
  
- Whole family more of grammars and automata – covered in automata theory

10/23/14

22

## Sample Grammar

- Language: Parenthesized sums of 0's and 1's
- $\langle \text{Sum} \rangle ::= 0$
- $\langle \text{Sum} \rangle ::= 1$
- $\langle \text{Sum} \rangle ::= \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$
- $\langle \text{Sum} \rangle ::= (\langle \text{Sum} \rangle)$

10/23/14

23

## BNF Grammars

- Start with a set of characters, **a,b,c,...**
  - We call these *terminals*
- Add a set of different characters, **X,Y,Z,**  
...
  - We call these *nonterminals*
- One special nonterminal **S** called *start symbol*

10/23/14

24

## BNF Grammars

- BNF rules (aka *productions*) have form
$$\mathbf{X} ::= y$$
where  $\mathbf{X}$  is any nonterminal and  $y$  is a string of terminals and nonterminals
- BNF *grammar* is a set of BNF rules such that every nonterminal appears on the left of some rule

10/23/14

25

## Sample Grammar

- Terminals: 0 1 + ( )
- Nonterminals: <Sum>
- Start symbol = <Sum>
- <Sum> ::= 0
- <Sum> ::= 1
- <Sum> ::= <Sum> + <Sum>
- <Sum> ::= (<Sum>)
- Can be abbreviated as
$$\langle \text{Sum} \rangle ::= 0 \mid 1 \mid \langle \text{Sum} \rangle + \langle \text{Sum} \rangle \mid (\langle \text{Sum} \rangle)$$

10/23/14

26

## BNF Derivations

- Given rules
$$\mathbf{X} ::= y\mathbf{Z}w \text{ and } \mathbf{Z} ::= v$$
we may replace  $\mathbf{Z}$  by  $v$  to say
$$\mathbf{X} \Rightarrow y\mathbf{Z}w \Rightarrow yvw$$
- Sequence of such replacements called *derivation*
- Derivation called *right-most* if always replace the right-most non-terminal

10/23/14

27

## BNF Derivations

- Start with the start symbol:

<Sum> =>

10/23/14

28

## BNF Derivations

- Pick a non-terminal

<Sum> =>

10/23/14

29

## BNF Derivations

- Pick a rule and substitute:
  - <Sum> ::= <Sum> + <Sum>

<Sum> => <Sum> + <Sum>

10/23/14

30

## BNF Derivations

- Pick a non-terminal:

$\langle \text{Sum} \rangle \Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$

10/23/14

31

## BNF Derivations

- Pick a rule and substitute:

■  $\langle \text{Sum} \rangle ::= ( \langle \text{Sum} \rangle )$   
 $\langle \text{Sum} \rangle \Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$

10/23/14

32

## BNF Derivations

- Pick a non-terminal:

$\langle \text{Sum} \rangle \Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$

10/23/14

33

## BNF Derivations

- Pick a rule and substitute:

■  $\langle \text{Sum} \rangle ::= \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$   
 $\langle \text{Sum} \rangle \Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$

10/23/14

34

## BNF Derivations

- Pick a non-terminal:

$\langle \text{Sum} \rangle \Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$

10/23/14

35

## BNF Derivations

- Pick a rule and substitute:

■  $\langle \text{Sum} \rangle ::= 1$   
 $\langle \text{Sum} \rangle \Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle$   
 $\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + \langle \text{Sum} \rangle$

10/23/14

36

## BNF Derivations

- Pick a non-terminal:

$$\begin{aligned} \langle \text{Sum} \rangle &\Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + \langle \text{Sum} \rangle \end{aligned}$$

10/23/14

37

## BNF Derivations

- Pick a rule and substitute:

$$\begin{aligned} &\bullet \langle \text{Sum} \rangle ::= 0 \\ \langle \text{Sum} \rangle &\Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + 0 \end{aligned}$$

10/23/14

38

## BNF Derivations

- Pick a non-terminal:

$$\begin{aligned} \langle \text{Sum} \rangle &\Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + 0 \end{aligned}$$

10/23/14

39

## BNF Derivations

- Pick a rule and substitute

$$\begin{aligned} &\bullet \langle \text{Sum} \rangle ::= 0 \\ \langle \text{Sum} \rangle &\Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) 0 \\ &\Rightarrow ( 0 + 1 ) + 0 \end{aligned}$$

10/23/14

40

## BNF Derivations

- $( 0 + 1 ) + 0$  is generated by grammar

$$\begin{aligned} \langle \text{Sum} \rangle &\Rightarrow \langle \text{Sum} \rangle + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + \langle \text{Sum} \rangle ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + \langle \text{Sum} \rangle \\ &\Rightarrow ( \langle \text{Sum} \rangle + 1 ) + 0 \\ &\Rightarrow ( 0 + 1 ) + 0 \end{aligned}$$

10/23/14

41

$\langle \text{Sum} \rangle ::= 0 \mid 1 \mid \langle \text{Sum} \rangle + \langle \text{Sum} \rangle \mid ( \langle \text{Sum} \rangle )$

$\langle \text{Sum} \rangle \Rightarrow$

10/23/14

42

## BNF Semantics

- The meaning of a BNF grammar is the set of all strings consisting only of terminals that can be derived from the Start symbol

10/23/14

43

## Regular Grammars

- Subclass of BNF
- Only rules of form  
 $\langle \text{nonterminal} \rangle ::= \langle \text{terminal} \rangle \langle \text{nonterminal} \rangle$  or  
 $\langle \text{nonterminal} \rangle ::= \langle \text{terminal} \rangle$  or  
 $\langle \text{nonterminal} \rangle ::= \epsilon$
- Defines same class of languages as regular expressions
- Important for writing lexers (programs that convert strings of characters into strings of tokens)
- Close connection to nondeterministic finite state automata – nonterminals  $\cong$  states; rule  $\cong$  edge

10/23/14

44

## Example

- Regular grammar:  
 $\langle \text{Balanced} \rangle ::= \epsilon$   
 $\langle \text{Balanced} \rangle ::= 0 \langle \text{OneAndMore} \rangle$   
 $\langle \text{Balanced} \rangle ::= 1 \langle \text{ZeroAndMore} \rangle$   
 $\langle \text{OneAndMore} \rangle ::= 1 \langle \text{Balanced} \rangle$   
 $\langle \text{ZeroAndMore} \rangle ::= 0 \langle \text{Balanced} \rangle$
- Generates even length strings where every initial substring of even length has same number of 0's as 1's

10/23/14

45

## Extended BNF Grammars

- Alternatives: allow rules of form  $X ::= y/z$ 
  - Abbreviates  $X ::= y, X ::= z$
- Options:  $X ::= y[v]$ 
  - Abbreviates  $X ::= yvz, X ::= yz$
- Repetition:  $X ::= y\{v\}^*z$ 
  - Can be eliminated by adding new nonterminal  $V$  and rules  $X ::= yz, X ::= yVz, V ::= v, V ::= vV$

10/23/14

46

## Parse Trees

- Graphical representation of derivation
- Each node labeled with either non-terminal or terminal
- If node is labeled with a terminal, then it is a leaf (no sub-trees)
- If node is labeled with a non-terminal, then it has one branch for each character in the right-hand side of rule used to substitute for it

10/23/14

47

## Example

- Consider grammar:  
 $\langle \text{exp} \rangle ::= \langle \text{factor} \rangle$   
 $\quad \quad \quad | \langle \text{factor} \rangle + \langle \text{factor} \rangle$   
 $\langle \text{factor} \rangle ::= \langle \text{bin} \rangle$   
 $\quad \quad \quad | \langle \text{bin} \rangle * \langle \text{exp} \rangle$   
 $\langle \text{bin} \rangle ::= 0 \quad | \quad 1$
- Problem: Build parse tree for  $1 * 1 + 0$  as an  $\langle \text{exp} \rangle$

10/23/14

48



### Example cont.

- 1 \* 1 + 0: <exp>

<exp> is the start symbol for this parse tree

### Example cont.

- 1 \* 1 + 0: <exp>  
                  |  
                  <factor>

Use rule: <exp> ::= <factor>

### Example cont.

- 1 \* 1 + 0: <exp>  
                  |  
                  <factor>  
                  / | \  
                  <bin> \* <exp>

Use rule: <factor> ::= <bin> \* <exp>

### Example cont.

- 1 \* 1 + 0: <exp>  
                  |  
                  <factor>  
                  / | \  
                  <bin> \* <exp>  
                  |    / | \  
                  1    <factor> + <factor>

Use rules: <bin> ::= 1 and  
          <exp> ::= <factor> +  
                  <factor>

### Example cont.

- 1 \* 1 + 0: <exp>  
                  |  
                  <factor>  
                  / | \  
                  <bin> \* <exp>  
                  |    / | \  
                  1    <factor> + <factor>  
                          |    |    |  
                          <bin> <bin>

Use rule: <factor> ::= <bin>

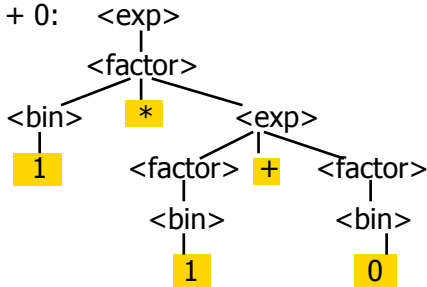
### Example cont.

- 1 \* 1 + 0: <exp>  
                  |  
                  <factor>  
                  / | \  
                  <bin> \* <exp>  
                  |    / | \  
                  1    <factor> + <factor>  
                          |    |    |  
                          <bin> <bin>  
                                  |    |  
                                  1    0

Use rules: <bin> ::= 1 | 0

## Example cont.

- 1 \* 1 + 0:



Fringe of tree is string generated by grammar

10/23/14

55

## Your Turn: 1 \* 0 + 0 \* 1

10/23/14

56

## Parse Tree Data Structures

- Parse trees may be represented by OCaml datatypes
- One datatype for each nonterminal
- One constructor for each rule
- Defined as mutually recursive collection of datatype declarations

10/23/14

57

## Example

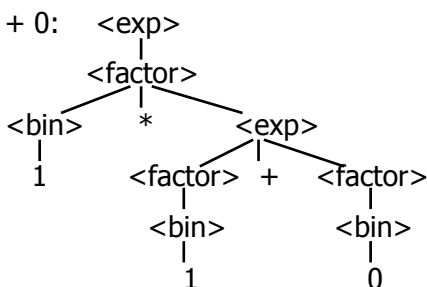
- Recall grammar:  
 $\langle \text{exp} \rangle ::= \langle \text{factor} \rangle \mid \langle \text{factor} \rangle + \langle \text{factor} \rangle$   
 $\langle \text{factor} \rangle ::= \langle \text{bin} \rangle \mid \langle \text{bin} \rangle * \langle \text{exp} \rangle$   
 $\langle \text{bin} \rangle ::= 0 \mid 1$
- type exp = Factor2Exp of factor  
| Plus of factor \* factor  
and factor = Bin2Factor of bin  
| Mult of bin \* exp  
and bin = Zero | One

10/23/14

58

## Example cont.

- 1 \* 1 + 0:



10/23/14

59

## Example cont.

- Can be represented as

```
Factor2Exp
(Mult(One,
      Plus(Bin2Factor One,
            Bin2Factor Zero)))
```

10/23/14

60

## Ambiguous Grammars and Languages

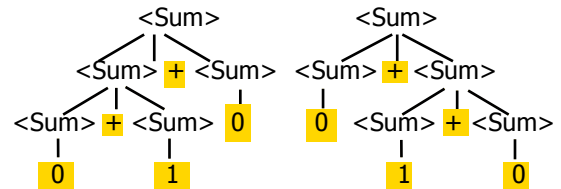
- A BNF grammar is *ambiguous* if its language contains strings for which there is more than one parse tree
- If all BNF's for a language are ambiguous then the language is *inherently ambiguous*

10/23/14

61

## Example: Ambiguous Grammar

- $0 + 1 + 0$



10/23/14

62

## Example

- What is the result for:

$$3 + 4 * 5 + 6$$

10/23/14

63

## Example

- What is the result for:

$$3 + 4 * 5 + 6$$

- Possible answers:

- $41 = ((3 + 4) * 5) + 6$
- $47 = 3 + (4 * (5 + 6))$
- $29 = (3 + (4 * 5)) + 6 = 3 + ((4 * 5) + 6)$
- $77 = (3 + 4) * (5 + 6)$

10/23/14

64

## Example

- What is the value of:

$$7 - 5 - 2$$

10/23/14

65

## Example

- What is the value of:

$$7 - 5 - 2$$

- Possible answers:

- In Pascal, C++, SML assoc. left  
 $7 - 5 - 2 = (7 - 5) - 2 = 0$
- In APL, associate to right  
 $7 - 5 - 2 = 7 - (5 - 2) = 4$

10/23/14

66



## Two Major Sources of Ambiguity

- Lack of determination of operator precedence
- Lack of determination of operator associativity
  
- Not the only sources of ambiguity