

Programming Languages and Compilers (CS 421)

Grigore Rosu
2110 SC, UIUC

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

Slides by Elsa Gunter, based in part on slides by Mattox Beckman, as updated by Vikram Adve and Gul Agha

3/20/2014

1

Objective

- Finish the discussion on lexing
- Ocamllex (takes .mll files and generates .ml)
- Context-Free Grammars and BNF

3/20/2014

2

General Input for .mll Files

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

3/20/2014

3

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

3/20/2014

4

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*

3/20/2014

5

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 *e₁ ∨ e₂*

3/20/2014

6

Ocamllex Regular Expression

- $[c_1 - c_2]$: choice of any character between first and second inclusive, as determined by character codes
- $[\wedge c_1 - c_2]$: choice of any character NOT in set
- e^* : same as before
- e^+ : same as $e e^*$
- $e?$: option - was $e_1 \vee \varepsilon$

3/20/2014

7

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* = *regex*
- e_1 as *id*: binds the result of e_1 to *id* to be used in the associated *action*

3/20/2014

8

Ocamllex Manual

- More details can be found at

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

3/20/2014

9

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 +
```

3/20/2014

10

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 }
```

3/20/2014

11

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???
```

3/20/2014

12

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"
```

3/20/2014

13

Problem

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

3/20/2014

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 }
```

3/20/2014

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

3/20/2014

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 }
```

3/20/2014

17

Dealing with comments

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

3/20/2014

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 }
```

3/20/2014

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 }
```

3/20/2014

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 }
```

3/20/2014

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

3/20/2014

22

Sample Grammar

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

3/20/2014

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*

3/20/2014

24

BNF Grammars

- BNF rules (aka *productions*) have form
$$\mathbf{X} ::= \gamma$$
where \mathbf{X} is any nonterminal and γ 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

3/20/2014

25

Sample Grammar

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

3/20/2014

26

BNF Derivations

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

3/20/2014

27

BNF Derivations

- Start with the start symbol:

<Sum> =>

3/20/2014

28

BNF Derivations

- Pick a non-terminal

<Sum> =>

3/20/2014

29

BNF Derivations

- Pick a rule and substitute:
 - $\text{<Sum>} ::= \text{<Sum>} + \text{<Sum>}$
- <Sum> => <Sum> + <Sum>

3/20/2014

30

BNF Derivations

- Pick a non-terminal:

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

3/20/2014 31

BNF Derivations

- Pick a rule and substitute:
 - $\langle \text{Sum} \rangle ::= (\langle \text{Sum} \rangle)$

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

3/20/2014 32

BNF Derivations

- Pick a non-terminal:

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

3/20/2014 33

BNF Derivations

- Pick a rule and substitute:
 - $\langle \text{Sum} \rangle ::= \langle \text{Sum} \rangle + \langle \text{Sum} \rangle$

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

3/20/2014 34

BNF Derivations

- Pick a non-terminal:

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

3/20/2014 35

BNF Derivations

- Pick a rule and substitute:
 - $\langle \text{Sum} \rangle ::= 1$

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

3/20/2014 36

BNF Derivations

- Pick a non-terminal:

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

3/20/2014 37

BNF Derivations

- Pick a rule and substitute:
 - $\langle \text{Sum} \rangle ::= 0$

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

3/20/2014 38

BNF Derivations

- Pick a non-terminal:

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

3/20/2014 39

BNF Derivations

- Pick a rule and substitute
 - $\langle \text{Sum} \rangle ::= 0$

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

3/20/2014 40

BNF Derivations

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

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

3/20/2014 41

BNF Derivations

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

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

3/20/2014 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

3/20/2014

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)

3/20/2014

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

3/20/2014

45

Extended BNF Grammars

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

3/20/2014

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

3/20/2014

47

Example

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

3/20/2014

48

Example cont.

- 1 * 1 + 0: <exp>

<exp> is the start symbol for this parse tree

3/20/2014 49

Example cont.

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

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

3/20/2014 50

Example cont.

- 1 * 1 + 0: <exp>
 - <factor>
 - <bin>
 - *
 - <exp>

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

3/20/2014 51

Example cont.

- 1 * 1 + 0: <exp>
 - <factor>
 - <bin>
 - 1
 - *
 - <exp>
 - <factor>
 - +
 - <factor>

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

3/20/2014 52

Example cont.

- 1 * 1 + 0: <exp>
 - <factor>
 - <bin>
 - 1
 - *
 - <exp>
 - <factor>
 - <bin>
 - +
 - <factor>
 - <bin>

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

3/20/2014 53

Example cont.

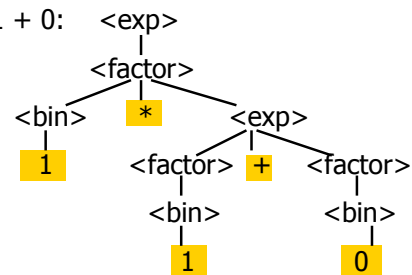
- 1 * 1 + 0: <exp>
 - <factor>
 - <bin>
 - 1
 - *
 - <exp>
 - <factor>
 - <bin>
 - 1
 - +
 - <factor>
 - <bin>
 - 0

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

3/20/2014 54

Example cont.

- 1 * 1 + 0:



Fringe of tree is string generated by grammar

3/20/2014

55

Your Turn: 1 * 0 + 0 * 1

3/20/2014

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

3/20/2014

57

Example

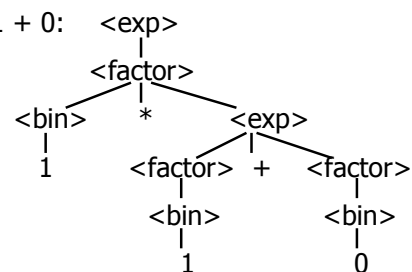
- Recall grammar:
 - $\langle \text{exp} \rangle ::= \langle \text{factor} \rangle \mid \langle \text{factor} \rangle + \langle \text{factor} \rangle$
 - $\langle \text{factor} \rangle ::= \langle \text{bin} \rangle \mid \langle \text{bin} \rangle * \langle \text{exp} \rangle$
 - $\langle \text{bin}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

3/20/2014

58

Example cont.

- 1 * 1 + 0:



3/20/2014

59

Example cont.

- Can be represented as

```

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

3/20/2014

60

Ambiguous Grammars and Languages

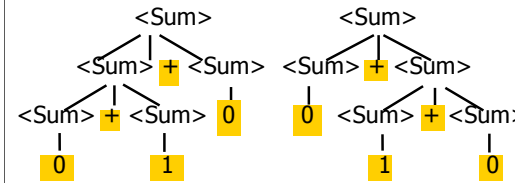
- A BNF grammar is *ambiguous* if its language contains strings for which there is more than one parse tree
- If all BNF's for a language are ambiguous then the language is *inherently ambiguous*

3/20/2014

61

Example: Ambiguous Grammar

- $0 + 1 + 0$



3/20/2014

62

Example

- What is the result for:
 $3 + 4 * 5 + 6$

3/20/2014

63

Example

- What is the result for:
 $3 + 4 * 5 + 6$
- Possible answers:
 - $41 = ((3 + 4) * 5) + 6$
 - $47 = 3 + (4 * (5 + 6))$
 - $29 = (3 + (4 * 5)) + 6 = 3 + ((4 * 5) + 6)$
 - $77 = (3 + 4) * (5 + 6)$

3/20/2014

64

Example

- What is the value of:
 $7 - 5 - 2$

3/20/2014

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
 $7 - 5 - 2 = 7 - (5 - 2) = 4$

3/20/2014

66

Two Major Sources of Ambiguity

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

3/20/2014

67

Disambiguating a Grammar

- Given ambiguous grammar G , with start symbol S , find a grammar G' with same start symbol, such that
language of G = language of G'
- Not always possible
- No algorithm in general

3/20/2014

68

Disambiguating a Grammar

- Idea: Each non-terminal represents all strings having some property
- Identify these properties (often in terms of things that can't happen)
- Use these properties to inductively guarantee every string in language has a unique parse

3/20/2014

69

Example

- Ambiguous grammar:
$$\begin{aligned} \langle \text{exp} \rangle &::= 0 \mid 1 \mid \langle \text{exp} \rangle + \langle \text{exp} \rangle \\ &\quad \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle \end{aligned}$$
- String with more than one parse:
$$\begin{aligned} 0 + 1 + 0 \\ 1 * 1 + 1 \end{aligned}$$
- Source of ambiguity: associativity and precedence

3/20/2014

70

How to Enforce Associativity

- Have at most one recursive call per production
- When two or more recursive calls would be natural leave right-most one for right associativity, left-most one for left associativity

10/4/07

71

Example

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

10/4/07

72

Operator Precedence

- Operators of highest precedence evaluated first (bind more tightly).
- Precedence for infix binary operators given in following table
- Needs to be reflected in grammar

10/4/07

73

Precedence Table - Sample

	Fortan	Pascal	C/C++	Ada	SML
highest	**	*, /, div, mod	++, --	**	div, mod, /, *
	*, /	+, -	*, /, %	*, /, mod	+, -, ^
	+, -		+, -	+, -	::

10/4/07

74

First Example Again

- In any above language,

$$3 + 4 * 5 + 6 = 29$$
- In APL, all infix operators have same precedence
 - Thus we still don't know what the value is (handled by associativity)
- How do we handle precedence in grammar?

10/4/07

75

Precedence in Grammar

- Higher precedence translates to longer derivation chain
- Example:

$$\langle \text{exp} \rangle ::= 0 \mid 1 \mid \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle$$
- Becomes

$$\begin{aligned} \langle \text{exp} \rangle &::= \langle \text{mult_exp} \rangle \mid \langle \text{exp} \rangle + \langle \text{mult_exp} \rangle \\ \langle \text{mult_exp} \rangle &::= \langle \text{id} \rangle \mid \langle \text{mult_exp} \rangle * \langle \text{id} \rangle \\ \langle \text{id} \rangle &::= 0 \mid 1 \end{aligned}$$

10/4/07

76