

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

11/1/12

1

Ocamlyacc Input

- File format:

```
%{  
    <header>  
}%  
    <declarations>  
%%  
    <rules>  
%%  
    <trailer>
```

11/1/12

2

Ocamlyacc <header>

- Contains arbitrary Ocaml code
- Typically used to give types and functions needed for the semantic actions of rules and to give specialized error recovery
- May be omitted
- <footer> similar. Possibly used to call parser

11/1/12

3

Ocamlyacc <declarations>

- %token** *symbol ... symbol*
Declare given symbols as tokens
- %token <type>** *symbol ... symbol*
Declare given symbols as token constructors, taking an argument of type <type>
- %start** *symbol ... symbol*
Declare given symbols as entry points; functions of same names in <grammar>.ml

11/1/12

4

Ocamlyacc <declarations>

- %type <type>** *symbol ... symbol*
Specify type of attributes for given symbols. Mandatory for start symbols
- %left** *symbol ... symbol*
- %right** *symbol ... symbol*
- %nonassoc** *symbol ... symbol*
Associate precedence and associativity to given symbols. Same line, same precedence; earlier line, lower precedence (broadest scope)

11/1/12

5

Ocamlyacc <rules>

- nonterminal :**
symbol ... symbol { semantic_action }
| ...
| *symbol ... symbol { semantic_action }*
;
Semantic actions are arbitrary Ocaml expressions
- Must be of same type as declared (or inferred) for *nonterminal*
- Access semantic attributes (values) of symbols by position: \$1 for first symbol, \$2 to second ...

11/1/12

6

Example - Base types

```
(* File: expr.ml *)
type expr =
  Term_as_Expr of term
  | Plus_Expr of (term * expr)
  | Minus_Expr of (term * expr)
and term =
  Factor_as_Term of factor
  | Mult_Term of (factor * term)
  | Div_Term of (factor * term)
and factor =
  Id_as_Factor of string
  | Parenthesized_Expr_as_Factor of expr
```

11/1/12

7

Example - Lexer (exprlex.mll)

```
{ (*open Exprparse*) }
let numeric = ['0' - '9']
let letter = ['a' - 'z' 'A' - 'Z']
rule token = parse
  | "+" {Plus_token}
  | "-" {Minus_token}
  | "*" {Times_token}
  | "/" {Divide_token}
  | "(" {Left_parenthesis}
  | ")" {Right_parenthesis}
  | letter (letter|numeric|"_")* as id {Id_token id}
  | [' ' '\t' '\n'] {token lexbuf}
  | eof {EOL}
```

11/1/12

8

Example - Parser (exprparse.mly)

```
%{ open Expr
%}
%token <string> Id_token
%token Left_parenthesis Right_parenthesis
%token Times_token Divide_token
%token Plus_token Minus_token
%token EOL
%start main
%type <expr> main
%%
```

11/1/12

9

Example - Parser (exprparse.mly)

```
expr:
  term
    { Term_as_Expr $1 }
  | term Plus_token expr
    { Plus_Expr ($1, $3) }
  | term Minus_token expr
    { Minus_Expr ($1, $3) }
```

11/1/12

10

Example - Parser (exprparse.mly)

```
term:
  factor
    { Factor_as_Term $1 }
  | factor Times_token term
    { Mult_Term ($1, $3) }
  | factor Divide_token term
    { Div_Term ($1, $3) }
```

11/1/12

11

Example - Parser (exprparse.mly)

```
factor:
  Id_token
    { Id_as_Factor $1 }
  | Left_parenthesis expr Right_parenthesis
    { Parenthesized_Expr_as_Factor $2 }
main:
  | expr EOL
    { $1 }
```

11/1/12

12

Example - Using Parser

```
# #use "expr.ml";;  
...  
# #use "exprparse.ml";;  
...  
# #use "exprlex.ml";;  
...  
# let test s =  
  let lexbuf = Lexing.from_string (s^"\n") in  
    main token lexbuf;;
```

11/1/12

13

Example - Using Parser

```
# test "a + b";;  
- : expr =  
Plus_Expr  
(Factor_as_Term (Id_as_Factor "a"),  
Term_as_Expr (Factor_as_Term  
(Id_as_Factor "b")))
```

11/1/12

14

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*
- Your job: *disambiguate given grammar*
 - Write a new grammar that is **not** ambiguous that generates the **same** language

11/1/12

15

Two Major Sources of Ambiguity

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

11/1/12

16

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

11/1/12

17

Example

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

11/1/12

18

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

11/1/12

19

Precedence in Grammar

- Higher precedence translates to longer derivation chain
- Example:
 $\langle \text{exp} \rangle ::= \langle \text{id} \rangle \mid \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle$
- Becomes
 $\langle \text{exp} \rangle ::= \langle \text{mult_exp} \rangle \mid \langle \text{exp} \rangle + \langle \text{mult_exp} \rangle$
 $\langle \text{mult_exp} \rangle ::= \langle \text{id} \rangle \mid \langle \text{mult_exp} \rangle * \langle \text{id} \rangle$

11/1/12

20

Recursive Descent Parsing

- Recursive descent parsers are a class of parsers derived fairly directly from BNF grammars
- A recursive descent parser traces out a parse tree in top-down order, corresponding to a left-most derivation (LL - left-to-right scanning, leftmost derivation)

11/1/12

21

Recursive Descent Parsing

- Each nonterminal in the grammar has a subprogram associated with it; the subprogram parses all phrases that the nonterminal can generate
- Each nonterminal in right-hand side of a rule corresponds to a recursive call to the associated subprogram

11/1/12

22

Recursive Descent Parsing

- Each subprogram must be able to decide how to begin parsing by looking at the left-most character in the string to be parsed
 - May do so directly, or indirectly by calling another parsing subprogram
- Recursive descent parsers, like other top-down parsers, cannot be built from left-recursive grammars
 - Sometimes can modify grammar to suit

11/1/12

23

Sample Grammar

$$\langle \text{expr} \rangle ::= \langle \text{term} \rangle \mid \langle \text{term} \rangle + \langle \text{expr} \rangle \mid \langle \text{term} \rangle - \langle \text{expr} \rangle$$
$$\langle \text{term} \rangle ::= \langle \text{factor} \rangle \mid \langle \text{factor} \rangle * \langle \text{term} \rangle \mid \langle \text{factor} \rangle / \langle \text{term} \rangle$$
$$\langle \text{factor} \rangle ::= \langle \text{id} \rangle \mid (\langle \text{expr} \rangle)$$

11/1/12

24

Tokens as OCaml Types

- + - * / () <id>
 - Becomes an OCaml datatype
- ```
type token =
 Id_token of string
 | Left_parenthesis | Right_parenthesis
 | Times_token | Divide_token
 | Plus_token | Minus_token
```

11/1/12

25

## Parse Trees as Datatypes

```
<expr> ::= <term> | <term> + <expr>
 | <term> - <expr>
```

```
type expr =
 Term_as_Expr of term
 | Plus_Expr of (term * expr)
 | Minus_Expr of (term * expr)
```

11/1/12

26

## Parse Trees as Datatypes

```
<term> ::= <factor> | <factor> *
 <term>
 | <factor> / <term>
```

```
and term =
 Factor_as_Term of factor
 | Mult_Term of (factor * term)
 | Div_Term of (factor * term)
```

11/1/12

27

## Parse Trees as Datatypes

```
<factor> ::= <id> | (<expr>)
```

```
and factor =
 Id_as_Factor of string
 | Parenthesized_Expr_as_Factor of expr
```

11/1/12

28

## Parsing Lists of Tokens

- Will create three mutually recursive functions:
  - `expr : token list -> (expr * token list)`
  - `term : token list -> (term * token list)`
  - `factor : token list -> (factor * token list)`
- Each parses what it can and gives back parse and remaining tokens

11/1/12

29

## Parsing an Expression

```
<expr> ::= <term> [(+ | -) <expr>]
let rec expr tokens =
 (match term tokens
 with (term_parse , tokens_after_term) ->
 (match tokens_after_term
 with(Plus_token :: tokens_after_plus) ->
```

11/1/12

30

## Parsing an Expression

```

<expr> ::= <term> [(+ | -) <expr>]
let rec expr tokens =
 (match term tokens
 with (term_parse , tokens_after_term) ->
 (match tokens_after_term
 with (Plus_token :: tokens_after_plus) ->

```

11/1/12

31

## Parsing a Plus Expression

```

<expr> ::= <term> [(+ | -) <expr>]
let rec expr tokens =
 (match term tokens
 with (term_parse , tokens_after_term) ->
 (match tokens_after_term
 with (Plus_token :: tokens_after_plus) ->

```

11/1/12

32

## Parsing a Plus Expression

```

<expr> ::= <term> [(+ | -) <expr>]
let rec expr tokens =
 (match term tokens
 with (term_parse , tokens_after_term) ->
 (match tokens_after_term
 with (Plus_token :: tokens_after_plus) ->

```

11/1/12

33

## Parsing a Plus Expression

```

<expr> ::= <term> [(+ | -) <expr>]
let rec expr tokens =
 (match term tokens
 with (term_parse , tokens_after_term) ->
 (match tokens_after_term
 with (Plus_token :: tokens_after_plus) ->

```

11/1/12

34

## Parsing a Plus Expression

```

<expr> ::= <term> + <expr>
(match expr tokens_after_plus
 with (expr_parse , tokens_after_expr) ->
 (Plus_Expr (term_parse , expr_parse),
 tokens_after_expr))

```

11/1/12

35

## Parsing a Plus Expression

```

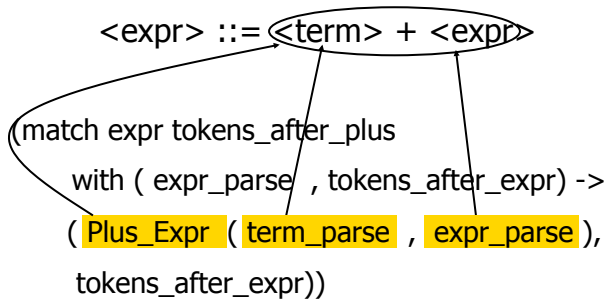
<expr> ::= <term> + <expr>
(match expr tokens_after_plus
 with (expr_parse , tokens_after_expr) ->
 (Plus_Expr (term_parse , expr_parse),
 tokens_after_expr))

```

11/1/12

36

## Building Plus Expression Parse Tree



11/1/12

37

## Parsing a Minus Expression

$\langle \text{expr} \rangle ::= \langle \text{term} \rangle - \langle \text{expr} \rangle$

```

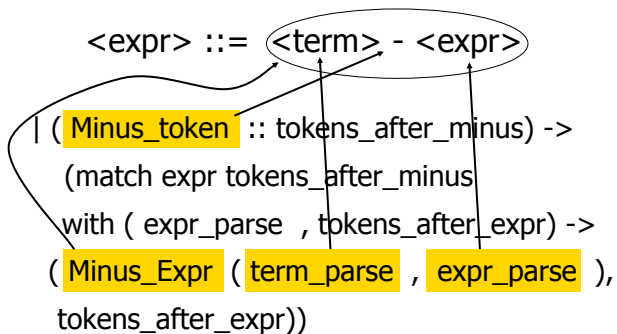
| (Minus_token :: tokens_after_minus) ->
 (match expr tokens_after_minus
 with (expr_parse , tokens_after_expr) ->
 (Minus_Expr (term_parse , expr_parse),
 tokens_after_expr))

```

11/1/12

38

## Parsing a Minus Expression



11/1/12

39

## Parsing an Expression as a Term

$\langle \text{expr} \rangle ::= \langle \text{term} \rangle$

```

| _ -> (Term_as_Expr term_parse ,
tokens_after_term)))

```

- Code for **term** is same except for replacing addition with multiplication and subtraction with division

11/1/12

40

## Parsing Factor as Id

$\langle \text{factor} \rangle ::= \langle \text{id} \rangle$

```

and factor tokens =
 (match tokens
 with (Id_token id_name :: tokens_after_id) =
 (Id_as_Factor id_name, tokens_after_id))

```

11/1/12

41

## Parsing Factor as Parenthesized Expression

$\langle \text{factor} \rangle ::= ( \langle \text{expr} \rangle )$

```

| factor (Left_parenthesis :: tokens) =
 (match expr tokens
 with (expr_parse , tokens_after_expr) ->
 (Expr_Paren (expr_parse), tokens_after_expr))

```

11/1/12

42

## Parsing Factor as Parenthesized Expression

$\langle \text{factor} \rangle ::= ( \langle \text{expr} \rangle )$

(match tokens\_after\_expr  
with Right\_parenthesis :: tokens\_after\_rparen ->  
( Parenthesized\_Expr\_as\_Factor expr\_parse ,  
tokens\_after\_rparen)

11/1/12

## Error Cases

- What if no matching right parenthesis?  
| \_ -> raise (Failure "No matching rparen" ) )
- What if no leading id or left parenthesis?  
| \_ -> raise (Failure "No id or lparen" ) );;

11/1/12

44

$(a + b) * c - d$

expr [Left\_parenthesis; Id\_token "a";  
Plus\_token; Id\_token "b";  
Right\_parenthesis; Times\_token;  
Id\_token "c"; Minus\_token;  
Id\_token "d"];;

11/1/12

45

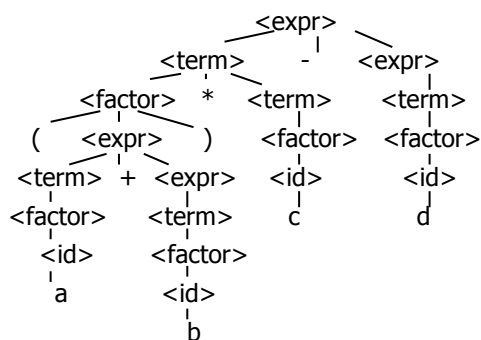
$(a + b) * c - d$

```
- : expr * token list =
(Minus_Expr
 (Mult_Term
 (Parenthesized_Expr_as_Factor
 (Plus_Expr
 (Factor_as_Term (Id_as_Factor "a"),
 Term_as_Expr (Factor_as_Term
 (Id_as_Factor "b")))),
 Factor_as_Term (Id_as_Factor "c")),
 Term_as_Expr (Factor_as_Term (Id_as_Factor
 "d")))),
 [])
```

11/1/12

46

$(a + b) * c - d$



11/1/12

47

$a + b * c - d$

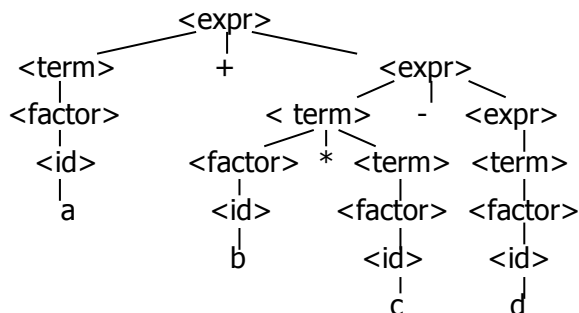
```
expr [Id_token "a"; Plus_token; Id_token "b";
Times_token; Id_token "c"; Minus_token;
Id_token "d"];;
- : expr * token list =
(Plus_Expr
 (Factor_as_Term (Id_as_Factor "a"),
 Minus_Expr
 (Mult_Term (Id_as_Factor "b", Factor_as_Term
 (Id_as_Factor "c")),
 Term_as_Expr (Factor_as_Term (Id_as_Factor
 "d")))),
 [])
```

11/1/12

48



 a + b \* c - d



11/1/12

49

 ( a + b \* c - d


```
expr [Left_parenthesis; Id_token "a";
Plus_token; Id_token "b"; Times_token;
Id_token "c"; Minus_token; Id_token "d"];;
```

Exception: Failure "No matching rparen".

Can't parse because it was expecting a right parenthesis but it got to the end without finding one

11/1/12

50

 a + b ) \* c - d \*)

```
expr [Id_token "a"; Plus_token; Id_token "b";
Right_parenthesis; Times_token; Id_token "c";
Minus_token; Id_token "d"];;
```

- : expr \* token list =

(Plus\_Expr

(Factor\_as\_Term (Id\_as\_Factor "a"),

Term\_as\_Expr (Factor\_as\_Term (Id\_as\_Factor "b"))),

[Right\_parenthesis; Times\_token; Id\_token "c"; Minus\_token; Id\_token "d"])

11/1/12

51

 Parsing Whole String

- Q: How to guarantee whole string parses?
- A: Check returned tokens empty

let parse tokens =

match **expr** tokens

with (expr\_parse, []) -> expr\_parse

| \_ -> raise (Failure "No parse");;

- Fixes <expr> as start symbol

11/1/12

52

 Streams in Place of Lists

- More realistically, we don't want to create the entire list of tokens before we can start parsing
- We want to generate one token at a time and use it to make one step in parsing
- Will use (token \* (unit -> token)) or (token \* (unit -> token option)) in place of token list

11/1/12

53

 Problems for Recursive-Descent Parsing

- Left Recursion:  
A ::= Aw  
translates to a subroutine that loops forever
- Indirect Left Recursion:  
A ::= Bw  
B ::= Av  
causes the same problem

11/1/12

54

## Problems for Recursive-Descent Parsing

- Parser must always be able to choose the next action based only on the very next token
- Pairwise Disjointedness Test: Can we always determine which rule (in the non-extended BNF) to choose based on just the first token

11/1/12

55

## Pairwise Disjointedness Test

- For each rule  
 $A ::= y$   
 Calculate  
 $\text{FIRST}(y) = \{a \mid y \Rightarrow^* aw\} \cup \{\epsilon \mid \text{if } y \Rightarrow^* \epsilon\}$
- For each pair of rules  $A ::= y$  and  $A ::= z$ , require  $\text{FIRST}(y) \cap \text{FIRST}(z) = \{\}$

11/1/12

56

## Example

Grammar:

$\langle S \rangle ::= \langle A \rangle a \langle B \rangle b$

$\langle A \rangle ::= \langle A \rangle b \mid b$

$\langle B \rangle ::= a \langle B \rangle \mid a$

$\text{FIRST}(\langle A \rangle b) = \{b\}$

$\text{FIRST}(b) = \{b\}$

Rules for  $\langle A \rangle$  not pairwise disjoint

11/1/12

57

## Eliminating Left Recursion

- Rewrite grammar to shift left recursion to right recursion
  - Changes associativity
- Given  
 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle$  and  
 $\langle \text{expr} \rangle ::= \langle \text{term} \rangle$
- Add new non-terminal  $\langle e \rangle$  and replace above rules with  
 $\langle \text{expr} \rangle ::= \langle \text{term} \rangle \langle e \rangle$   
 $\langle e \rangle ::= + \langle \text{term} \rangle \langle e \rangle \mid \epsilon$

11/1/12

58

## Factoring Grammar

- Test too strong: Can't handle  
 $\langle \text{expr} \rangle ::= \langle \text{term} \rangle [ ( + \mid - ) \langle \text{expr} \rangle ]$
- Answer: Add new non-terminal and replace above rules by  
 $\langle \text{expr} \rangle ::= \langle \text{term} \rangle \langle e \rangle$   
 $\langle e \rangle ::= + \langle \text{term} \rangle \langle e \rangle$   
 $\langle e \rangle ::= - \langle \text{term} \rangle \langle e \rangle$   
 $\langle e \rangle ::= \epsilon$
- You are delaying the decision point

11/1/12

59

## Example

Both  $\langle A \rangle$  and  $\langle B \rangle$   
have problems:

Transform grammar  
to:

$\langle S \rangle ::= \langle A \rangle a \langle B \rangle b$      $\langle S \rangle ::= \langle A \rangle a \langle B \rangle b$

$\langle A \rangle ::= \langle A \rangle b \mid b$

$\langle A \rangle ::= b \langle A1 \rangle$

$\langle B \rangle ::= a \langle B \rangle \mid a$

$\langle A1 \rangle ::= b \langle A1 \rangle \mid \epsilon$

$\langle B \rangle ::= a \langle B1 \rangle$

$\langle B1 \rangle ::= a \langle B1 \rangle \mid \epsilon$

11/1/12

60