

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

10/25/16

1

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*

10/25/16

2

BNF Grammars

- BNF rules (aka *productions*) have form $X ::= y$ where **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

10/25/16

3

Sample Grammar

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

10/25/16

4

BNF Derivations

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

10/25/16

5

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

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

10/25/16

6

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

10/25/16

7

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)

10/25/16

8

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

10/25/16

9

Extended BNF Grammars

- Alternatives: allow rules of form $X ::= y/z$
 - 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 $X ::= yz, X ::= yVz, V ::= v, V ::= vV$

10/25/16

10

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

10/25/16

11

Example

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

10/25/16

12

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>

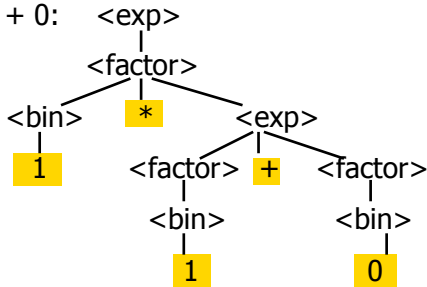
Example cont.

- 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

10/25/16

19

Your Turn: 1 * 0 + 0 * 1

10/25/16

20

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

10/25/16

21

Example

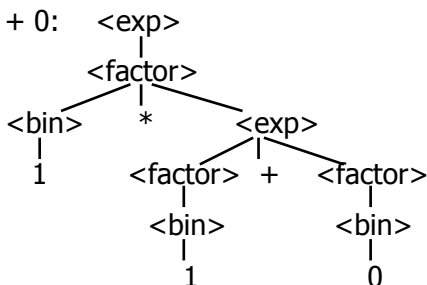
- Recall grammar:
 $\langle \text{exp} \rangle ::= \langle \text{factor} \rangle \mid \langle \text{factor} \rangle + \langle \text{factor} \rangle$
 $\langle \text{factor} \rangle ::= \langle \text{bin} \rangle \mid \langle \text{bin} \rangle * \langle \text{exp} \rangle$
 $\langle \text{bin} \rangle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

10/25/16

22

Example cont.

- 1 * 1 + 0:



10/25/16

23

Example cont.

- Can be represented as

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

10/25/16

24

Ambiguous Grammars and Languages

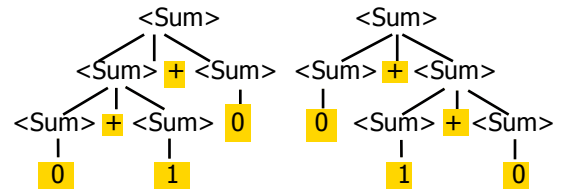
- A BNF grammar is *ambiguous* if its language contains strings for which there is more than one parse tree
- If all BNF's for a language are ambiguous then the language is *inherently ambiguous*

10/25/16

25

Example: Ambiguous Grammar

- $0 + 1 + 0$



10/25/16

26

Example

- What is the result for:

$$3 + 4 * 5 + 6$$

10/25/16

27

Example

- What is the result for:

$$3 + 4 * 5 + 6$$

- Possible answers:

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

10/25/16

28

Example

- What is the value of:

$$7 - 5 - 2$$

10/25/16

29

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

10/25/16

30

Two Major Sources of Ambiguity

- Lack of determination of operator precedence
- Lack of determination of operator associativity

- Not the only sources of ambiguity

10/25/16

31

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

10/25/16

32

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

10/25/16

33

Steps to Grammar Disambiguation

- Identify the rules and a smallest use that display ambiguity
- Decide which parse to keep; why should others be thrown out?
- What syntactic restrictions on subexpressions are needed to throw out the bad (while keeping the good)?
- Add a new non-terminal and rules to describe this set of restricted subexpressions (called stratifying, or refactoring)
- Replace old rules to use new non-terminals
- Rinse and repeat

10/25/16

34

Example

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

10/25/16

35

Two Major Sources of Ambiguity

- Lack of determination of operator precedence
- Lack of determination of operator associativity

- Not the only sources of ambiguity

10/4/07

36

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

37

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

38

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

39

Precedence Table - Sample

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

10/4/07

40

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

41

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

10/4/07

42

Ocamlyacc Input

- File format:

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

10/25/16

43

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

10/25/16

44

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

10/25/16

45

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)

10/25/16

46

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 ...

10/25/16

47

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
```

10/25/16

48

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}
```

10/25/16

49

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
%%
```

10/25/16

50

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) }
```

10/25/16

51

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) }
```

10/25/16

52

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 }
```

10/25/16

53

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;;
```

10/25/16

54



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")))
```