# CSCI312 Principles of Programming Languages

LL Parsing

Xu Liu

Derived from Keith Cooper's COMP 412 at Rice University

Copyright © 2006 The McGraw-Hill Companies, Inc.



# Syntactic Analysis

Phase also known as: parser

Purpose is to recognize source structure

Input: tokens

Output: parse tree or abstract syntax tree

A recursive descent parser is one in which each nonterminal in the grammar is converted to a function which recognizes input derivable from the nonterminal.

# Terms

Nullable non-terminals First set

# $T(EBNF) = Code: A \rightarrow W$

- 1 If w is nonterminal, call it.
- 2 If w is terminal, match it against given token.
- 3 If w is  $\{ w' \}$ :

. . .

- while (token in First(w')) T(w')
- 4 If w is: w1 | ... | wn,

switch (token) {
 case First(w1): T(w1); break;

case First(wn): T(wn); break;

5 Switch (cont.): If some wi is empty, use: default: break;

Otherwise

default: error(token);

- 6 If w = [w'], rewrite as (|w'|) and use rule 4.
- 7 If w = X1 ... Xn, T(w) = T(X1); ... T(Xn);

# Outline

See more general problems in a top down parser Backtracking — select appropriate productions Left recursion — revise grammars Predictive parsing — more than recursive descent Table-driven parsing

### Parsing Techniques

Top-down parsers (LL(1), recursive descent)

- Start at the root of the parse tree and grow toward leaves
- Pick a production & try to match the input
- Bad "pick" ⇒ may need to backtrack
- Some grammars are backtrack-free (predictive parsing)

Bottom-up parsers (LR(1), operator precedence)

- Start at the leaves and grow toward root
- As input is consumed, encode possibilities in an internal state
- Start in a state valid for legal first tokens
- Bottom-up parsers handle a large class of grammars

A top-down parser starts with the root of the parse tree The root node is labeled with the goal symbol of the grammar

Top-down parsing algorithm:

Construct the root node of the parse tree

Repeat until lower fringe of the parse tree matches the input string

- 1 At a node labeled A, select a production with A on its lhs and, for each symbol on its rhs, construct the appropriate child
- 2 When a terminal symbol is added to the fringe and it doesn't match the fringe, backtrack
- 3 Find the next node to be expanded

(label  $\in$  NT)

The key is picking the right production in step 1

- That choice should be guided by the input string

### Remember the expression grammar?

We will call this version "the classic expression grammar"

from last lecture

0	Goal	$\rightarrow$	Expr
1	Expr	$\rightarrow$	Expr + Term
2			Expr - Term
3			Term
4	Term	$\rightarrow$	Term * Factor
5		l	Term / Factor
6			Factor
7	Factor	$\rightarrow$	(Expr)
8			number
9		I	id

And the input  $\underline{x} - \underline{2} * \underline{y}$ 

Let's try  $\underline{x} - \underline{2} * \underline{y}$ :

Rule	Sentential Form	Input
_	Goal	<u>↑x - 2 * y</u>



Goal

Let's try  $\underline{x} - \underline{2} * \underline{y}$ : Goal Rule Sentential Form Input Expr 1×-2\*y Goal Expr ↑<u>× - 2</u> \* y 0 Expr Term Expr +Term ↑<u>x</u> - <u>2</u> \* y 1 3 Term + Term  $\uparrow \underline{x} - \underline{2} \times \underline{y}$ (Term 6 Factor + Term  $\uparrow \underline{x} - \underline{2} * \underline{y}$ (Fact.) <id,<u>x</u>> +Term ↑<u>x</u> - <u>2</u> \* <u>y</u> 9  $\rightarrow$ <id,<u>x</u>> +Term <u>x ↑- 2 \* y</u> <id,x>

This worked well, except that "-" doesn't match "+" The parser must backtrack to here

Continuing with  $\underline{x} - \underline{2} * \underline{y}$ :



⇒ Now, we need to expand Term - the last NT on the fringe



- We have more input, but no NTs left to expand
- The expansion terminated too soon
- $\Rightarrow$  Need to backtrack



#### The Point:

The parser must make the right choice when it expands a NT. Wrong choices lead to wasted effort.

# Another possible parse

Other choices for expansion are possible

Rule	Sentential Form	Input	
—	Goal	↑ <u>×</u> - <u>2</u> * γ	
0	Expr	1×-2*γ	Consumes no input!
1	Expr +Term	× 2* y	
1	Expr + Term +Term	↑ <u>×</u> - <u>2</u> * y	
1	Expr + Term +Term + Term	↑ <u>× - 2</u> * γ	
1	And so on	<u>x - 2 * y</u>	

#### This expansion doesn't terminate

- Wrong choice of expansion leads to non-termination
- Non-termination is a bad property for a parser to have
- Parser must make the right choice

Top-down parsers cannot handle left-recursive grammars

Formally,

A grammar is left recursive if  $\exists A \in NT$  such that  $\exists a \text{ derivation } A \Rightarrow^{+} A\alpha$ , for some string  $\alpha \in (NT \cup T)^{+}$ 

Our classic expression grammar is left recursive

- This can lead to non-termination in a top-down parser
- In a top-down parser, any recursion must be right recursion
- We would like to convert the left recursion to right recursion

Non-termination is <u>always</u> a bad property in a compiler

# Eliminating Left Recursion

To remove left recursion, we can transform the grammar

Consider a grammar fragment of the form

```
Fee \rightarrow Fee \alpha
```

where neither  $\alpha$  nor  $\beta$  start with Fee



The new grammar defines the same language as the old grammar, using only right recursion.

Added a reference to the empty string

# Eliminating Left Recursion

The expression grammar contains two cases of left recursion

Applying the transformation yields

Expr	→ Term Expr'	Term	$\rightarrow$	Factor Term'
Expr'	→ + Term Expr'	Term'	$\rightarrow$	* Factor Term'
	- Term Expr'			/ Factor Term'
	3			3

These fragments use only right recursion

Right recursion often means right associativity. In this case, the grammar does not display any particular associative bias.

# Picking the "Right" Production

If it picks the wrong production, a top-down parser may backtrack Alternative is to look ahead in input & use context to pick correctly

How much lookahead is needed?

• In general, an arbitrarily large amount

Fortunately,

- Large subclasses of CFGs can be parsed with limited lookahead
- Most programming language constructs fall in those subclasses

Among the interesting subclasses are LL(1) and LR(1) grammars

We will focus, for now, on LL(1) grammars & predictive parsing

### **Predictive Parsing**

#### Basic idea

Given A  $\rightarrow \alpha \mid \beta$  , the parser should be able to choose between  $\alpha$  &  $\beta$ 

#### FIRST sets

For some rhs  $\alpha \in G$ , define FIRST( $\alpha$ ) as the set of tokens that appear as the first symbol in some string that derives from  $\alpha$ That is,  $\underline{x} \in FIRST(\alpha)$  iff  $\alpha \Rightarrow^* \underline{x} \gamma$ , for some  $\gamma$ 

#### Basic idea

Given A  $\rightarrow \alpha \mid \beta$  , the parser should be able to choose between  $\alpha$  &  $\beta$ 

#### FIRST sets

For some rhs  $\alpha \in G$ , define FIRST( $\alpha$ ) as the set of tokens that appear as the first symbol in some string that derives from  $\alpha$ That is,  $\underline{x} \in \text{FIRST}(\alpha)$  iff  $\alpha \Rightarrow^* \underline{x} \gamma$ , for some  $\gamma$ 

#### The LL(1) Property

If  $A \rightarrow \alpha$  and  $A \rightarrow \beta$  both appear in the grammar, we would like

 $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$ 

This would allow the parser to make a correct choice with a lookahead of exactly one symbol !

This is almost correct See the next slide

### **Predictive Parsing**

What about  $\epsilon$ -productions?

- $\Rightarrow$  They complicate the definition of LL(1)
- If  $A \rightarrow \alpha$  and  $A \rightarrow \beta$  and  $\epsilon \in \text{FIRST}(\alpha)$ , then we need to ensure that  $\text{FIRST}(\beta)$  is disjoint from FOLLOW(A), too, where

FOLLOW(A) = the set of terminal symbols that can immediately follow A in a sentential form

Define  $FIRST^{+}(A \rightarrow \alpha)$  as

- FIRST( $\alpha$ )  $\cup$  FOLLOW(A), if  $\varepsilon \in$  FIRST( $\alpha$ )
- FIRST( $\alpha$ ), otherwise

Then, a grammar is LL(1) iff  $A \rightarrow \alpha$  and  $A \rightarrow \beta$  implies

 $\mathsf{FIRST}^{\mathsf{+}}(\mathsf{A} \rightarrow \alpha) \cap \mathsf{FIRST}^{\mathsf{+}}(\mathsf{A} \rightarrow \beta) = \emptyset$ 

# Recursive Descent Parsing

Recall the expression grammar, after transformation

0	Goal	$\rightarrow$	Expr
1	Expr	$\rightarrow$	Term Expr'
2	Expr'	$\rightarrow$	+ Term Expr'
3			- Term Expr'
4			8
5	Term	$\rightarrow$	Factor Term'
6	Term'	$\rightarrow$	* Factor Term
7			/ Factor Term
8			ε
9	Factor	$\rightarrow$	(Expr)
10			<u>number</u>
11			id

This produces a parser with six <u>mutually recursive</u> routines:

- Goal
- Expr
- EPrime
- Term
- TPrime
- Factor

Each recognizes one NT or  $\mathsf{T}$ 

The term <u>descent</u> refers to the direction in which the parse tree is built.

### What If My Grammar Is Not LL(1)?

Can we transform a non-LL(1) grammar into an LL(1) grammar?

- In general, the answer is no
- In some cases, however, the answer is yes

Assume a grammar G with productions  $A \rightarrow \alpha \beta_1$  and  $A \rightarrow \alpha \beta_2$ 

• If  $\alpha$  derives anything other than  $\epsilon$ , then

$$\mathsf{FIRST}^{\mathsf{+}}(\mathsf{A} \to \alpha \,\beta_1) \cap \mathsf{FIRST}^{\mathsf{+}}(\mathsf{A} \to \alpha \,\beta_2) \neq \emptyset$$

- And the grammar is not LL(1)
- If we pull the common prefix,  $\alpha$ , into a separate production, we may make the grammar LL(1).

$$A \rightarrow \alpha A', A' \rightarrow \beta_1 \text{ and } A' \rightarrow \beta_2$$

Now, if FIRST<sup>+</sup>( $A' \rightarrow \beta_1$ )  $\cap$  FIRST<sup>+</sup>( $A' \rightarrow \beta_2$ ) =  $\emptyset$ , G may be LL(1)

# What If My Grammar Is Not LL(1)?

Left Factoring

```
For each nonterminal A
find the longest prefix \alpha common to 2 or more alternatives for A
if \alpha \neq \varepsilon then
replace all of the productions
A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \alpha \beta_3 | \dots | \alpha \beta_n | \gamma
with
A \rightarrow \alpha A' | \gamma
A' \rightarrow \beta_1 | \beta_2 | \beta_3 | \dots | \beta_n
```

Repeat until no nonterminal has alternative rhs' with a common prefix

This transformation makes some grammars into LL(1) grammars There are languages for which no LL(1) grammar exists

# Left Factoring Example

Consider a simple right-recursive expression grammar

0	Goal	$\rightarrow$	Expr
1	Expr	$\rightarrow$	Term + Expr
2		Ι	Term – Expr
3		Ι	Term
4	Term	$\rightarrow$	Factor * Term
5		Ι	Factor / Term
6		Ι	Factor
7	Factor	$\rightarrow$	number
8		Ι	id

To choose between 1, 2, & 3, an LL(1) parser must look past the <u>number</u> or <u>id</u> to see the operator. FIRST<sup>+</sup>(1) = FIRST<sup>+</sup>(2) = FIRST<sup>+</sup>(3) and FIRST<sup>+</sup>(4) = FIRST<sup>+</sup>(5) = FIRST<sup>+</sup>(6) Let's left factor this grammar.

# Left Factoring Example

After Left Factoring, we have

0	Goal	$\rightarrow$	Expr
1	Expr	$\rightarrow$	Term Expr'
2	Expr'	$\rightarrow$	+ Expr
3			- Expr
4			ε
5	Term	$\rightarrow$	Factor Term'
6	Term'	$\rightarrow$	* Term
7		Ι	/ Term
8		Ι	ε
9	Factor	$\rightarrow$	<u>number</u>
10		Ι	id

Clearly, FIRST\*(2), FIRST\*(3), & FIRST\*(4) are disjoint, as are FIRST\*(6), FIRST\*(7), & FIRST\*(8)

The grammar now has the LL(1) property

This transformation makes some grammars into LL(1) grammars. There are languages for which no LL(1) grammar exists. FIRST( $\alpha$ )

For some  $\alpha \in (T \cup NT)^*$ , define FIRST( $\alpha$ ) as the set of tokens that appear as the first symbol in some string that derives from  $\alpha$ 

That is,  $\underline{x} \in \text{FIRST}(\alpha)$  iff  $\alpha \Rightarrow^* \underline{x} \gamma$ , for some  $\gamma$ 

Follow(A)

- For some  $A \in NT$ , define FOLLOW(A) as the set of symbols that can occur immediately after A in a valid sentential form
- FOLLOW(S) = {EOF}, where S is the start symbol

To build FOLLOW sets, we need FIRST sets ...

### Computing FIRST Sets

Already studied in previous lectures

Computing FOLLOW Sets

```
for each A \in NT, FOLLOW(A) \leftarrow \emptyset
FOLLOW(S) \leftarrow \{EOF\}
while (FOLLOW sets are still changing)
    for each p \in P, of the form A \rightarrow B_1 B_2 \dots B_k
        TRAILER \leftarrow FOLLOW(A)
        for i \leftarrow k down to 1
            if B_i \in NT then
                                                           // domain check
                 FOLLOW(B_i) \leftarrow FOLLOW(B_i) \cup TRAILER
                 if \varepsilon \in FIRST(B_i)
                                                         // add right context
                   then TRAILER \leftarrow TRAILER \cup (FIRST(B<sub>i</sub>) - {\varepsilon})
                   else TRAILER \leftarrow FIRST(B<sub>i</sub>) // no \varepsilon \Rightarrow no right context
            else TRAILER \leftarrow \{B_i\}
                                           // B_i \in T \Rightarrow only 1 symbol
```

# Classic Expression Grammar

	0	Goal → Expr	Symbol	FIRST	FOLLOW
	1	Expr → Term Expr'	num	num	Ø
	2	Expr' → + Term Expr'	id	id	Ø
	3	I - Term Expr'	+	+	Ø
	4	8	-	-	Ø
	_		*	*	Ø
	5	Term 7 Factor Term'	/	/	Ø
	6	Term' → * Factor	(	(	Ø
	7	/ Factor	)	)	Ø
	8	3	eof	eof	Ø
	9	Factor → <u>number</u>	3	3	Ø
	10	<u>id</u>	Goal	(.id.num	eof
	11	(Expr)	Expr	(.id.num	) eof
F	IRS	$T^{+}(A \rightarrow \beta)$ is identical to	Expr'	<b>+</b> . <b>-</b> . E	) eof
F	IRS	$T(\beta)$ except for productiond 4	Term	(.id.num	+ ). eof
a	nd 8		Term'	*./.ε	+). eof
(			Factor	(,id,num	+,-,*,/,),eof
F	IRS	$T^{+}(Expr' \rightarrow \varepsilon)$ is $\{\varepsilon, j\}$ , eof $\}$		<u> </u>	

# Classic Expression Grammar

0	Goal	→ Expr	Prod'n	FIRST+
1	Expr	→ Term Expr'	0	(.id.num
2	Expr'	→ + Term Expr'	1	(.id.num
3		- Term Expr'	2	+
4		ε	3	-
5	Term	→ Factor Term'	4	ε. <b>). eof</b>
6	Term	→ * Factor	5	(.id.num
7		/ Factor	6	*
8		1 <i>E</i>	7	/
0	Factor	i → number	8	ε. <b>+). eof</b>
7	1 actor		9	number
10			10	id
11		(Expr)	11	(

# Building Top-down Parsers for LL(1) Grammars

Given an LL(1) grammar, and its FIRST & FOLLOW sets ...

- Emit a routine for each non-terminal
  - Nest of if-then-else statements to check alternate rhs's
  - Each returns true on success and throws an error on false
  - Simple, working (perhaps ugly) code
- This automatically constructs a recursive-descent parser

Improving matters

- Nest of if-then-else statements may be slow
  - Good case statement implementation would be better
- What about a table to encode the options?
  - Interpret the table with a skeleton, as we did in scanning

I don't know of a system that does this

...

# **Building Top-down Parsers**



Building Top-down Parsers

Building the complete table

Need a row for every NT & a column for every T

# LL(1) Expression Parsing Table

		+	-	*	/	Id	Num	(	)	EOF
	Goal	_	—	_	_	0	0	0	—	_
	Expr	_	_	_	_	1	1	1	—	—
	Expr'	2	3	_	_	—	—	—	4	4
	Term	_	_	_	_	5	5	5	—	—
Row we b	Term'	8	8	6	7	_	_	_	8	8
earlier	Factor	_	_	_	_	10	9	11	_	_

Building Top-down Parsers

Building the complete table

- Need a row for every NT & a column for every T
- Need an interpreter for the table (skeleton parser)

# LL(1) Skeleton Parser

```
word 

NextWord() 
// Initial conditions, including
push EOF onto Stack // a stack to track local goals
push the start symbol, S, onto Stack
TOS \leftarrow top of Stack
loop forever
 if TOS = FOF and word = FOF then
    break & report success // exit on success
  else if TOS is a terminal then
    if TOS matches word then
       pop Stack // recognized TOS
      word \leftarrow NextWord()
    else report error looking for TOS // error exit
  else
                      // TOS is a non-terminal
    if TABLE[TOS,word] is A \rightarrow B_1 B_2 \dots B_k then
       pop Stack // get rid of A
       push B_k, B_{k-1}, ..., B_1 // in that order
    else break & report error expanding TOS
 TOS \leftarrow top of Stack
```

# **Building Top-down Parsers**

Building the complete table

- Need a row for every NT & a column for every T
- Need a table-driven interpreter for the table
- Need an algorithm to build the table

#### Filling in TABLE[X,y], $X \in NT$ , $y \in T$

- 1. entry is the rule  $X \rightarrow \beta$ , if  $y \in FIRST^{+}(X \rightarrow \beta)$
- entry is error if rule 1 does not define

If any entry has more than one rule, G is not LL(1)

We call this algorithm the LL(1) table construction algorithm