

# Top-down Parsing Recursive Descent & LL(1) Comp 412

The lecture is self consistent, but the order of the three productions for Factor in the expression grammar is different than in 2e

Copyright 2010, Keith D. Cooper & Linda Torczon, all rights reserved.

Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.

Faculty from other educational institutions may use these materials for nonprofit educational purposes, provided this copyright notice is preserved.

### Roadmap (Where are we?)



#### We set out to study parsing

- Specifying syntax
  - Context-free grammars ???
- Top-down parsers
  - Algorithm & its problem with left recursion ?
  - Ambiguity ?
  - Left-recursion removal ?
- Predictive top-down parsing
  - The LL(1) condition ?
  - Simple recursive descent parsers today
  - First and Follow sets today
  - Table-driven LL(1) parsers today

# 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  $\gamma$ 

#### The LL(1) Property

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

$$\mathsf{FIRST}(\alpha) \cap \mathsf{FIRST}(\beta) = \varnothing$$

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  $\varepsilon$ -productions?

⇒ They complicate the definition of LL(1)

If  $A \to \alpha$  and  $A \to \beta$  and  $\epsilon \in \mathsf{FIRST}(\alpha)$ , then we need to ensure that  $\mathsf{FIRST}(\beta)$  is disjoint from  $\mathsf{FOLLOW}(A)$ , 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  $A \to \alpha$  and  $A \to \beta$  implies FIRST<sup>+</sup> $(A \to \alpha) \cap \text{FIRST}^+(A \to \beta) = \emptyset$ 

The intuition is straightforward. If A expands with 2 right hand sides, those rhs' must have distinct first symbols.

We only need FIRST\* sets because of  $\epsilon$ -productions.

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

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

- 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

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

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).

$$A \rightarrow \alpha A'$$
,  $A' \rightarrow \beta_1$  and  $A' \rightarrow \beta_2$ 

Now, if FIRST+(A'  $\rightarrow \beta_1$ )  $\cap$  FIRST+(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 \epsilon then replace all of the productions A \to \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3 \mid ... \mid \alpha \beta_n \mid \gamma with A \to \alpha A' \mid \gamma A' \to \beta_1 \mid \beta_2 \mid \beta_3 \mid ... \mid \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			Term - Expr
3			Term
4	Term	$\rightarrow$	Factor * Term
5			Factor / Term
6			Factor
7	Factor	$\rightarrow$	<u>number</u>
8			<u>id</u>

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.

$$FIRST^{+}(1) = FIRST^{+}(2) = FIRST^{+}(3)$$

and

$$First^+(4) = First^+(5) = First^+(6)$$

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			/ Term
8			ε
9	Factor	$\rightarrow$	<u>number</u>
10			<u>id</u>

```
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 and FOLLOW Sets



#### $FIRST(\alpha)$

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,  $\underline{x} \in \mathsf{FIRST}(\alpha)$  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



```
for each x \in T, FIRST(x) \leftarrow \{x\}
for each A \in NT, FIRST(A) \leftarrow \emptyset
while (FIRST sets are still changing) do
   for each p \in P, of the form A \rightarrow \beta do
      if \beta is B_1B_2...B_k then begin;
          rhs \leftarrow FIRST(B_1) - \{\varepsilon\}
          for i \leftarrow 1 to k-1 by 1 while \varepsilon \in FIRST(B_i) do
              rhs ← rhs \cup (FIRST(B_{i+1}) - {\varepsilon})
              end // for loop
                    // if-then
       end
        if i = k and \varepsilon \in FIRST(B_k)
           then rhs \leftarrow rhs \cup {\varepsilon}
         FIRST(A) \leftarrow FIRST(A) \cup rhs
        end // for loop
              // while loop
    end
```

Outer loop is monotone increasing for FIRST sets

 $\rightarrow$  | T  $\bigcirc$  NT  $\bigcirc$   $\varepsilon$  | is bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

```
For SheepNoise:

FIRST(Goal) = { baa }

FIRST(SN) = { baa }

FIRST(baa) = { baa }
```

#### Computing FIRST Sets



```
for each x \in T, FIRST(x) \leftarrow \{x\}
for each A \in NT, FIRST(A) \leftarrow \emptyset
while (FIRST sets are still changing) do
   for each p \in P, of the form A \rightarrow \beta do
      if \beta is B_1B_2...B_k then begin;
         rhs ← FIRST(B<sub>1</sub>) - {ε}
         for i \leftarrow 1 to k-1 by 1 while \varepsilon \in FIRST(B_i) do
              rhs ← rhs \cup (FIRST(B_{i+1}) - {\varepsilon})
              end // for loop
             // if-then
       end
       if i = k and \varepsilon \in FIRST(B_k)
           then rhs \leftarrow rhs \cup {\varepsilon}
        FIRST(A) \leftarrow FIRST(A) \cup rhs
        end // for loop
              // while loop
    end
```

Outer loop is monotone increasing for FIRST sets

```
\rightarrow | T \cup NT \cup \varepsilon | is nates
```

Inner loop is bounded by the length of the productions in the grammar

```
For SheepNoise:

FIRST(Goal) = \{ \underline{baa} \}

FIRST(SN) = \{ \underline{baa} \}

FIRST(\underline{baa}) = \{ \underline{baa} \}
```

### Computing FIRST Sets



```
for each x \in T, FIRST(x) \leftarrow \{x\}
for each A \in NT, FIRST(A) \leftarrow \emptyset
while (FIRST sets are still changing) do
   for each p \in P, of the form A \rightarrow \beta do
      if \beta is B_1B_2...B_k then begin;
         rhs \leftarrow FIRST(B_1) - \{\varepsilon\}
          for i \leftarrow 1 to k-1 by 1 while \varepsilon \in FIRST(B_i) do
              rhs ← rhs \cup (FIRST(B_{i+1}) - {\varepsilon})
              end // for loop
                    // if-then
       end
        if i = k and \varepsilon \in FIRST(B_k)
           then rhs \leftarrow rhs \cup {\varepsilon}
         FIRST(A) \leftarrow FIRST(A) \cup rhs
        end // for loop
              // while loop
    end
```

Outer loop is monotone increasing for FIRST sets

```
\rightarrow | T \cup NT \cup \varepsilon | is For Goal \rightarrow SN
```

Inner loop is bounded by the length of the productions in the grammar

```
For SheepNoise:

FIRST(Goal) = \{ \underline{baa} \}

FIRST(SN) = \{ \underline{baa} \}

FIRST(\underline{baa}) = \{ \underline{baa} \}
```

### 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_1B_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_i) - {\varepsilon})
                  else TRAILER \leftarrow FIRST(B<sub>i</sub>) // no \varepsilon => no right context
            else TRAILER \leftarrow \{B_i\}
                                       //B_i \in T \Rightarrow only 1 symbol
```

# Classic Expression Grammar

0	Goal	$\rightarrow$	Expr	Symbol	FIRST	FOLLOW
1	Expr		Term Expr'	<u>num</u>	<u>num</u>	Ø
2	Expr'		+ Term Expr'	<u>id</u>	<u>id</u>	Ø
3	CAPI	ı	•	+	+	Ø
			- Term Expr'	-	-	Ø
4		ı	3	*	*	Ø
5	Term	$\rightarrow$	Factor Term'	/	/	Ø
6	Term'	$\rightarrow$	* Factor Term'	<u>(</u>	(	Ø
7			/ Factor Term'	)	)	Ø
8			3	eof	eof	Ø
9	Factor	$\rightarrow$	number	3	3	Ø
10		1	<u>id</u>	Goal	<u>(,id,num</u>	eof
11		1	(Expr)	Expr	<u>(,id,num</u>	), eof
FIRS	FIRST <sup>+</sup> $(A \rightarrow \beta)$ is identical to FIRST $(\beta)$				+, -, ε	), eof
except for productiond 4 and 8			Term	<u>(,id,num</u>	+, -, <u>)</u> , eof	
FIRST+(Expr' $\rightarrow \varepsilon$ ) is $\{\varepsilon, j, eof\}$			Term'	*,/,ε	+,-, <u>)</u> , eof	
FIRST+(Term' $\rightarrow \epsilon$ ) is $\{\epsilon,+,-,\}$ , eof $\}$			Factor	(,id,num	+,-,*,/, <u>)</u> ,eof	

# Classic Expression Grammar

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

Prod'n	FIRST+
0	<u>(,id,num</u>
1	<u>(,id,num</u>
2	+
3	-
4	ε,), eof
5	<u>(,id,num</u>
6	*
7	/
8	ε,+,-,), eof
9	<u>number</u>
10	<u>id</u>
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 <</li>

#### 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

#### Strategy

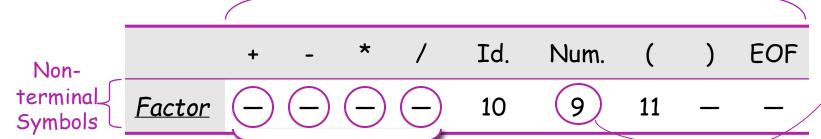
- Encode knowledge in a table
- Use a standard "skeleton" parser to interpret the table

#### Example

- The non-terminal Factor has 3 expansions
  - (Expr) or Identifier or Number
- Table might look like:

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

Terminal Symbols



Comp 412, Fall 2010

Cannot expand Factor into an operator  $\Rightarrow$  error

Expand Factor by rule 9 with input "number"

# 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	Term'	8	8	6	7	_	_	_	8	8
built ear	Factor	_	_	_	_	10	9	11	_	_

Table differs from 2e because order of productions in Factor differs.

# 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 = EOF and word = EOF then
    break & report success // exit on success
  else if TOS is a terminal then
    if TOS matches word then
      pop Stack // recognized TOS
      word ← NextWord()
    else report error looking for TOS // error exit
                            // TOS is a non-terminal
  else
    if TABLE[TOS, word] is A \rightarrow B_1B_2...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  $y \in FIRST^{+}(X \rightarrow \beta)$
- 2. 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

# Grammar and Sets for the LL(1) Construction



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

Prod'n	FIRST+
0	(,id,num
1	<u>(,id,num</u>
2	+
3	-
4	ε,), eof
5	<u>(,id,num</u>
6	*
7	/
8	ε,+,-,), eof
9	<u>number</u>
10	<u>id</u>
11	<u>(</u>

Right-recursive variant of the classic expression grammar
Comp 412, Fall 2010

Augmented FIRST sets for the grammar