

Parsing

Note by Baris Aktemur:

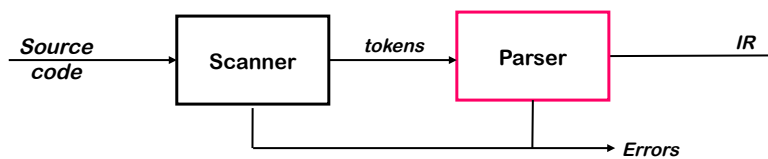
Our slides are adapted from Cooper and Torczon's slides that they prepared for COMP 412 at Rice.

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.

The Front End



Parser

- Checks the stream of words and their parts of speech (produced by the scanner) for grammatical correctness
- Determines if the input is syntactically well formed
- Need a mathematical model of syntax — a grammar G
- Need an algorithm for testing membership in $L(G)$
- Need to keep in mind that our goal is building parsers, not studying the mathematics of arbitrary languages

Limits of Regular Languages

Advantages of Regular Expressions

- Simple & powerful notation for specifying patterns
- Automatic construction of fast recognizers
- Many kinds of syntax can be specified with REs

Example — a regular expression for arithmetic expressions

$Term \rightarrow [a-zA-Z] ([a-zA-Z] | [0-9])^*$

$Op \rightarrow + | - | * | /$

$Expr \rightarrow (Term Op)^* Term$

$([a-zA-Z] ([a-zA-Z] | [0-9])^* (+ | - | * | /)^* [a-zA-Z] ([a-zA-Z] | [0-9])^*$

Limitation: Operator precedence?

2

Grammars

0	$Expr$	\rightarrow	$Expr Op Expr$
1			<u>number</u>
2			<u>id</u>
3	Op	\rightarrow	+
4			-
5			*
6			/

Rule	Sentential Form
—	$Expr$
0	$Expr Op Expr$
2	$\langle id, x \rangle Op Expr$
4	$\langle id, x \rangle - Expr$
0	$\langle id, x \rangle - Expr Op Expr$
1	$\langle id, x \rangle - \langle num, 2 \rangle Op Expr$
5	$\langle id, x \rangle - \langle num, 2 \rangle * Expr$
2	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$

- Such a sequence of rewrites is called a *derivation*
- Process of discovering a derivation is called *parsing*

We denote this derivation: $Expr \Rightarrow^* \underline{id} - \underline{num} * \underline{id}$

3

Derivations

The point of parsing is to construct a derivation

- At each step, we choose a nonterminal to replace
- Different choices can lead to different derivations

Two derivations are of interest

- *Leftmost derivation* – replace leftmost NT at each step
- *Rightmost derivation* – replace rightmost NT at each step

These are the two *systematic* derivations

(We don't care about randomly-ordered derivations!)

4

The Two Derivations for $x - 2 * y$

Rule	Sentential Form
–	<i>Expr</i>
0	<i>Expr Op Expr</i>
2	$\langle id, x \rangle Op Expr$
4	$\langle id, x \rangle - Expr$
0	$\langle id, x \rangle - Expr Op Expr$
1	$\langle id, x \rangle - \langle num, 2 \rangle Op Expr$
5	$\langle id, x \rangle - \langle num, 2 \rangle * Expr$
2	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$

Leftmost derivation

Rule	Sentential Form
–	<i>Expr</i>
0	<i>Expr Op Expr</i>
2	<i>Expr Op</i> $\langle id, y \rangle$
5	<i>Expr</i> * $\langle id, y \rangle$
0	<i>Expr Op Expr</i> * $\langle id, y \rangle$
1	<i>Expr Op</i> $\langle num, 2 \rangle$ * $\langle id, y \rangle$
4	<i>Expr</i> - $\langle num, 2 \rangle$ * $\langle id, y \rangle$
2	$\langle id, x \rangle$ - $\langle num, 2 \rangle$ * $\langle id, y \rangle$

Rightmost derivation

In both cases, $Expr \Rightarrow^* id - num * id$

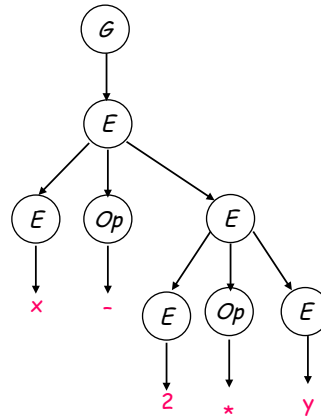
- The two derivations produce different parse trees
- The parse trees imply different evaluation orders!

5

Derivations and Parse Trees

Leftmost derivation

Rule	Sentential Form
—	<i>Expr</i>
0	<i>Expr Op Expr</i>
2	$\langle id, \underline{x} \rangle Op Expr$
4	$\langle id, \underline{x} \rangle - Expr$
0	$\langle id, \underline{x} \rangle - Expr Op Expr$
1	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
5	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$



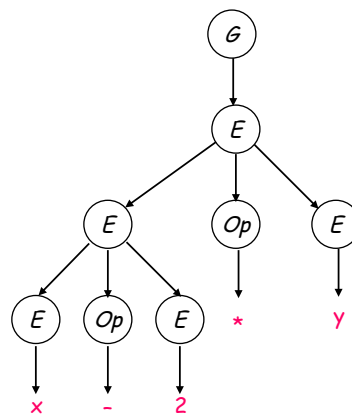
This evaluates as $x - (2 * y)$

6

Derivations and Parse Trees

Rightmost derivation

Rule	Sentential Form
—	<i>Expr</i>
0	<i>Expr Op Expr</i>
2	<i>Expr Op</i> $\langle id, \underline{y} \rangle$
5	<i>Expr</i> * $\langle id, \underline{y} \rangle$
0	<i>Expr Op Expr</i> * $\langle id, \underline{y} \rangle$
1	<i>Expr Op</i> $\langle num, \underline{2} \rangle$ * $\langle id, \underline{y} \rangle$
4	<i>Expr</i> - $\langle num, \underline{2} \rangle$ * $\langle id, \underline{y} \rangle$
2	$\langle id, \underline{x} \rangle$ - $\langle num, \underline{2} \rangle$ * $\langle id, \underline{y} \rangle$



This evaluates as $(x - 2) * y$

This ambiguity is **NOT** good

7

Derivations and Precedence

These two derivations point out a problem with the grammar:
It has no notion of precedence, or implied order of evaluation

To add precedence

- Create a nonterminal for each *level of precedence*
- Isolate the corresponding part of the grammar
- Force the parser to recognize high precedence subexpressions first

For algebraic expressions

- Parentheses first (level 1)
- Multiplication and division, next (level 2)
- Subtraction and addition, last (level 3)

8

Derivations and Precedence

Adding the standard algebraic precedence produces:

	0	Goal	→	Expr
level 3	1	Expr	→	Expr + Term
	2			Expr - Term
	3			Term
level 2	4	Term	→	Term * Factor
	5			Term / Factor
	6			Factor
level 1	7	Factor	→	(Expr)
	8			<u>number</u>
	9			<u>id</u>

This grammar is slightly larger

- Takes more rewriting to reach some of the terminal symbols
- Encodes expected precedence
- Produces same parse tree under leftmost & rightmost derivations
- Correctness trumps the speed of the parser

Let's see how it parses $x - 2 * y$

Cannot handle precedence in an RE for expressions

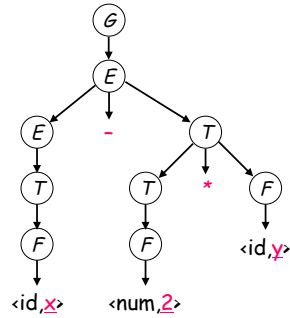
Introduced parentheses, too (beyond power of an RE)

One form of the "classic expression grammar" 9

Derivations and Precedence

Rule	Sentential Form
—	Goal
0	Expr
2	Expr - Term
4	Expr - Term * Factor
9	Expr - Term * <id,y>
6	Expr - Factor * <id,y>
8	Expr - <num,z> * <id,y>
3	Term - <num,z> * <id,y>
6	Factor - <num,z> * <id,y>
9	<id,x> - <num,z> * <id,y>

The rightmost derivation



Its parse tree

It derives $x - (z * y)$, along with an appropriate parse tree.

Both the leftmost and rightmost derivations give the same expression, because the grammar directly and explicitly encodes the desired precedence.

10

Ambiguous Grammars

Let's leap back to our original expression grammar.

It had other problems.

0	Expr	→	Expr Op Expr
1			number
2			id
3	Op	→	+
4			-
5			*
6			/

Rule	Sentential Form
—	Expr
0	Expr Op Expr
②	<id,x> Op Expr
4	<id,x> - Expr
0	<id,x> - Expr Op Expr
1	<id,x> - <num,z> Op Expr
5	<id,x> - <num,z> * Expr
2	<id,x> - <num,z> * <id,y>

- This grammar allows multiple leftmost derivations for $x - z * y$
- Hard to automate derivation if > 1 choice
- The grammar is *ambiguous*

Different choice than the first step

11

Two Leftmost Derivations for $x - 2 * y$

The Difference:

- Different productions chosen on the second step

Rule	Sentential Form
—	<i>Expr</i>
0	<i>Expr Op Expr</i>
②	$\langle id, x \rangle Op Expr$
4	$\langle id, x \rangle - Expr$
0	$\langle id, x \rangle - Expr Op Expr$
1	$\langle id, x \rangle - \langle num, 2 \rangle Op Expr$
5	$\langle id, x \rangle - \langle num, 2 \rangle * Expr$
1	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$

Original choice

Rule	Sentential Form
—	<i>Expr</i>
0	<i>Expr Op Expr</i>
①	<i>Expr Op Expr Op Expr</i>
2	$\langle id, x \rangle Op Expr Op Expr$
4	$\langle id, x \rangle - Expr Op Expr$
1	$\langle id, x \rangle - \langle num, 2 \rangle Op Expr$
5	$\langle id, x \rangle - \langle num, 2 \rangle * Expr$
2	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$

New choice

- Both derivations succeed in producing $x - 2 * y$

12

Ambiguous Grammars

Definitions

- If a grammar has more than one leftmost derivation for a single *sentential form*, the grammar is *ambiguous*
- If a grammar has more than one rightmost derivation for a single *sentential form*, the grammar is *ambiguous*
- The leftmost and rightmost derivations for a *sentential form* may differ, even in an unambiguous grammar
 - However, they must have the same parse tree!

Classic example — the *if-then-else* problem

$Stmt \rightarrow$ if *Expr* then *Stmt*
 | if *Expr* then *Stmt* else *Stmt*
 | ... other stmts ...

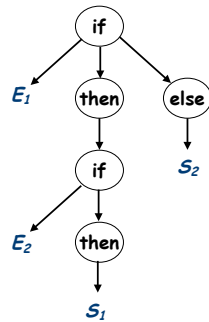
This ambiguity is inherent in the grammar

13

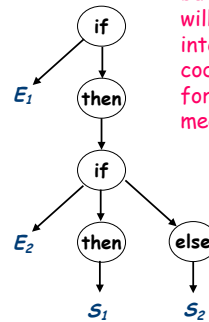
Ambiguity

This sentential form has two derivations

$\text{if } \underline{\text{Expr}}_1 \text{ then if } \underline{\text{Expr}}_2 \text{ then } \underline{\text{Stmt}}_1 \text{ else } \underline{\text{Stmt}}_2$



production 2, then
production 1



production 1, then
production 2

Part of the problem is that the structure built by the parser will determine the interpretation of the code, and these two forms have different meanings!

14

Ambiguity

The grammar forces the structure to match the desired meaning.

Removing the ambiguity

- Must rewrite the grammar to avoid generating the problem
- Match each else to innermost unmatched if (*common sense rule*)

0	$\text{Stmt} \rightarrow \text{if } \underline{\text{Expr}} \text{ then } \underline{\text{Stmt}}$
1	$\text{if } \underline{\text{Expr}} \text{ then } \underline{\text{WithElse}} \text{ else } \underline{\text{Stmt}}$
2	Other Statements
3	$\underline{\text{WithElse}} \rightarrow \text{if } \underline{\text{Expr}} \text{ then } \underline{\text{WithElse}} \text{ else } \underline{\text{WithElse}}$
4	Other Statements

With this grammar, example has only one rightmost derivation

Intuition: once into *WithElse*, we cannot generate an unmatched else
... a final if without an else can only come through rule 2 ...

15

Ambiguity

if Expr₁ then if Expr₂ then Stmt₁ else Stmt₂

Rule	Sentential Form
—	Stmt
0	if Expr then Stmt
1	if Expr then if Expr then WithElse else Stmt
2	if Expr then if Expr then WithElse else S ₂
4	if Expr then if Expr then S ₁ else S ₂
?	if Expr then if E ₂ then S ₁ else S ₂
?	if E ₁ then if E ₂ then S ₁ else S ₂

Other productions to derive Exprs

This grammar has only one rightmost derivation for the example

16

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

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
- Bottom-up parsers handle a large class of grammars

17

Top-down Parsing

Top-down Parsing

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”

0	<i>Goal</i>	→	<i>Expr</i>	
1	<i>Expr</i>	→	<i>Expr</i> + <i>Term</i>	
2			<i>Expr</i> - <i>Term</i>	
3			<i>Term</i>	
4	<i>Term</i>	→	<i>Term</i> * <i>Factor</i>	<i>And the input <u>x</u> - <u>2</u> * <u>y</u></i>
5			<i>Term</i> / <i>Factor</i>	
6			<i>Factor</i>	
7	<i>Factor</i>	→	(<i>Expr</i>)	
8			<u>number</u>	
9			<u>id</u>	

20

Example

Let's try x - 2 * y:

Goal

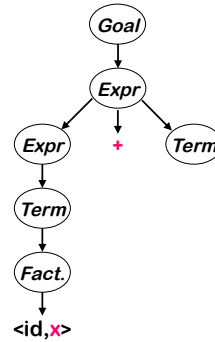
Rule	Sentential Form	Input
—	<i>Goal</i>	↑ <u>x</u> - <u>2</u> * <u>y</u>

↑ is the position in the input buffer 21

Example

Let's try $x - 2 * y$:

Rule	Sentential Form	Input
—	Goal	$\uparrow x - 2 * y$
0	Expr	$\uparrow x - 2 * y$
1	Expr + Term	$\uparrow x - 2 * y$
3	Term + Term	$\uparrow x - 2 * y$
6	Factor + Term	$\uparrow x - 2 * y$
9	$\langle id, x \rangle + Term$	$\uparrow x - 2 * y$
→	$\langle id, x \rangle + Term$	$x \uparrow - 2 * y$



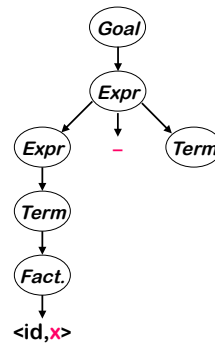
This worked well, except that “-” doesn't match “+”
The parser must backtrack to here

\uparrow is the position in the input buffer 22

Example

Continuing with $x - 2 * y$:

Rule	Sentential Form	Input
—	Goal	$\uparrow x - 2 * y$
0	Expr	$\uparrow x - 2 * y$
2	Expr - Term	$\uparrow x - 2 * y$
3	Term - Term	$\uparrow x - 2 * y$
6	Factor - Term	$\uparrow x - 2 * y$
9	$\langle id, x \rangle - Term$	$\uparrow x - 2 * y$
→	$\langle id, x \rangle \ominus Term$	$x \uparrow \ominus 2 * y$
→	$\langle id, x \rangle - Term$	$x \uparrow \uparrow 2 * y$



Now, “-” and “-” match | Now we can expand Term to match “2”

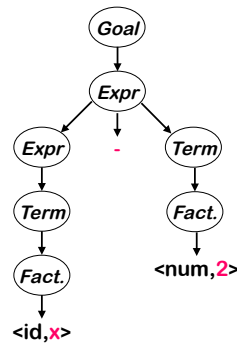
⇒ Now, we need to expand Term - the last NT on the fringe

23

Example

Trying to match the "2" in $x - 2 * y$:

Rule	Sentential Form	Input
→	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$
6	$\langle id, x \rangle - Factor$	$x - \uparrow 2 * y$
8	$\langle id, x \rangle - \langle num, 2 \rangle$	$x - \uparrow 2 * y$
→	$\langle id, x \rangle - \langle num, 2 \rangle$	$x - 2 \uparrow * y$



Where are we?

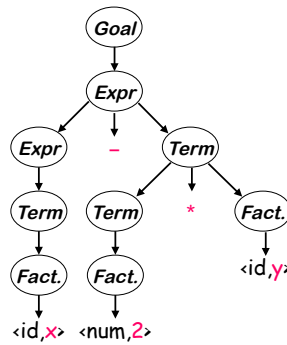
- "2" matches "2"
 - We have more input, but no NTs left to expand
 - The expansion terminated too soon
- ⇒ Need to backtrack

24

Example

Trying again with "2" in $x - 2 * y$:

Rule	Sentential Form	Input
→	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$
4	$\langle id, x \rangle - Term * Factor$	$x - \uparrow 2 * y$
6	$\langle id, x \rangle - Factor * Factor$	$x - \uparrow 2 * y$
8	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$x - \uparrow 2 * y$
→	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$x - 2 \uparrow * y$
→	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$x - 2 * \uparrow y$
9	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$x - 2 * \uparrow y$
→	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$x - 2 * \uparrow y$



The Point:

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

25

Another possible parse

Other choices for expansion are possible

Rule	Sentential Form	Input
–	Goal	$\uparrow x - 2^* y$
0	Expr	$\uparrow x - 2^* y$
1	Expr + Term	$\uparrow x - 2^* y$
1	Expr + Term + Term	$\uparrow x - 2^* y$
1	Expr + Term + Term + Term	$\uparrow x - 2^* y$
1	And so on	$\uparrow x - 2^* y$

Consumes no input!

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

26

Left Recursion

Top-down parsers cannot handle left-recursive grammars

Formally,

A grammar is *left recursive* if $\exists A \in NT$ such that
 \exists a 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 always a bad property in a compiler

27

Eliminating Left Recursion

To remove left recursion, we can transform the grammar

Consider a grammar fragment of the form

$$\begin{array}{l} Fee \rightarrow Fee \alpha \\ | \beta \end{array}$$

where neither α nor β start with Fee

We can rewrite this fragment as

$$\begin{array}{l} Fee \rightarrow \beta Fie \\ Fie \rightarrow \alpha Fie \\ | \varepsilon \end{array}$$

where Fie is a new non-terminal

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

Added a reference to the empty string

28

Eliminating Left Recursion

The expression grammar contains two cases of left recursion

$$\begin{array}{ll} Expr \rightarrow Expr + Term & Term \rightarrow Term * Factor \\ | Expr - Term & | Term * Factor \\ | Term & | Factor \end{array}$$

Applying the transformation yields

$$\begin{array}{ll} Expr \rightarrow Term Expr' & Term \rightarrow Factor Term' \\ Expr' \rightarrow + Term Expr' & Term' \rightarrow * Factor Term' \\ | - Term Expr' & | / Factor Term' \\ | \varepsilon & | \varepsilon \end{array}$$

These fragments use only right recursion

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

Eliminating Left Recursion

Substituting them back into the grammar yields

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

- This grammar is correct, if somewhat non-intuitive.
- It is left associative, as was the original
 - ⇒ The naïve transformation yields a right recursive grammar, which changes the implicit associativity
- A top-down parser will terminate using it.
- A top-down parser may need to backtrack with it.

30

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
- Use the Cocke-Younger, Kasami algorithm or Earley’s algorithm

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

31

Predictive Parsing

Basic idea

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

FIRST sets

For some rhs $\alpha \in \mathcal{G}$, define $\text{FIRST}(\alpha)$ as the set of tokens that appear as the first symbol in some string that derives from α

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

32

Predictive Parsing

Basic idea

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

FIRST sets

For some rhs $\alpha \in \mathcal{G}$, define $\text{FIRST}(\alpha)$ as the set of tokens that appear as the first symbol in some string that derives from α

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

The LL(1) Property

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

$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$$

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

33

Predictive Parsing

What about ϵ -productions?

⇒ 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 $\text{FOLLOW}(A)$, too, where

$\text{FOLLOW}(A)$ = the set of terminal symbols that can immediately follow A in a sentential form

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

- $\text{FIRST}(\alpha) \cup \text{FOLLOW}(A)$, if $\epsilon \in \text{FIRST}(\alpha)$
- $\text{FIRST}(\alpha)$, otherwise

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

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

34

Classic Expression Grammar

0	<i>Goal</i> → <i>Expr</i>	Symbol	FIRST	FOLLOW
1	<i>Expr</i> → <i>Term Expr'</i>	num	num	∅
2	<i>Expr'</i> → + <i>Term Expr'</i>	id	id	∅
3	- <i>Term Expr'</i>	+	+	∅
4	ε	-	-	∅
5	<i>Term</i> → <i>Factor Term'</i>	*	*	∅
6	<i>Term'</i> → * <i>Factor Term'</i>	/	/	∅
7	/ <i>Factor Term'</i>	((∅
8	ε))	∅
9	<i>Factor</i> → <u>number</u>	eof	eof	∅
10	<u>id</u>	ε	ε	∅
11	(<i>Expr</i>)	<i>Goal</i>	(,id,num	eof
		<i>Expr</i>	(,id,num),eof
		<i>Expr'</i>	+, -, ε),eof
		<i>Term</i>	(,id,num	+, -,),eof
		<i>Term'</i>	*, /, ε	+, -,),eof
		<i>Factor</i>	(,id,num	+, -, *, /,),eof

$\text{FIRST}^+(A \rightarrow \beta)$ is identical to $\text{FIRST}(\beta)$ except for production 4 and 8

$\text{FIRST}^+(\text{Expr}' \rightarrow \epsilon)$ is $\{\epsilon, \text{eof}\}$

$\text{FIRST}^+(\text{Term}' \rightarrow \epsilon)$ is $\{\epsilon, +, -, \text{eof}\}$

Classic Expression Grammar

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

36

Example

0	<i>Goal</i> → <i>Expr</i>
1	<i>Expr</i> → <i>Term Expr'</i>
2	<i>Expr'</i> → + <i>Expr</i>
3	- <i>Expr</i>
4	ε
5	<i>Term</i> → <i>Factor Term'</i>
6	<i>Term'</i> → * <i>Term</i>
7	/ <i>Term</i>
8	ε
9	<i>Factor</i> → <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

37

Predictive Parsing

Given a grammar that has the $LL(1)$ property

- Can write a simple routine to recognize each *lhs*
- Code is both simple & fast

Consider $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with
 $FIRST^+(A \rightarrow \beta_i) \cap FIRST^+(A \rightarrow \beta_j) = \emptyset$ if $i \neq j$

```

/* find an A */
if (current_word ∈ FIRST(A→β1))
  find a β1 and return true
else if (current_word ∈ FIRST(A→β2))
  find a β2 and return true
else if (current_word ∈ FIRST(A→β3))
  find a β3 and return true
else
  report an error and return false
  
```

Grammars with the $LL(1)$ property are called *predictive grammars* because the parser can “predict” the correct expansion at each point in the parse.

Parsers that capitalize on the $LL(1)$ property are called *predictive parsers*.

One kind of predictive parser is the *recursive descent* parser.

Of course, there is more detail to “find a β_i ” (p. 103 in EAC, 1st Ed.)

38

Recursive Descent Parsing

Recall the expression grammar, after transformation

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

This produces a parser with six *mutually recursive* routines:

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

Each recognizes one *NT* or *T*

The term *descent* refers to the direction in which the parse tree is built.

39

Recursive Descent Parsing (Procedural)

A couple of routines from the expression parser

Goal()

```
token ← next_token();  
if (Expr() = true & token = EOF)  
  then report success;  
else  
  report syntax error;  
  return false;
```

Expr()

```
if (Term() = false)  
  then return false;  
else return Eprime();
```

looking for Number, Identifier, or
"(", found token instead, or failed
to find Expr or ")" after "("

Factor()

```
if (token = Number) then  
  token ← next_token();  
  return true;  
else if (token = Identifier) then  
  token ← next_token();  
  return true;  
else if (token = Lparen)  
  token ← next_token();  
  if (Expr() = true & token = Rparen) then  
    token ← next_token();  
    return true;  
// fall out of if statement  
report syntax error;  
return false;
```

*EPrime, Term, & TPrime follow the same
basic lines (Figure 3.10, EaC2)*

40

Bottom-up Parsing

Parsing Techniques

Bottom-up parsers LR(1)

- 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
- Bottom-up parsers build a rightmost derivation in reverse
- Parsers can be auto-generated from grammars

- We will skip this topic...

42

Summary

	<i>Advantages</i>	<i>Disadvantages</i>
<i>Top-down Recursive descent, LL(1)</i>	Fast Good locality Simplicity Good error detection	Hand-coded High maintenance Right associativity
<i>LR(1)</i>	Fast Deterministic langs. Automatable Left associativity	Large working sets Poor error messages Large table sizes

43