



# Parsing — Part II

(Top-down parsing, left-recursion removal)

CS434

Lecture 7

Spring 2005

Department of Computer Science

University of Alabama

Joel Jones

Copyright 2003, Keith D. Cooper, Ken Kennedy & 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.

# 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"  $\Rightarrow$  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



# Top-down Parsing

---

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 the 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?

Version with precedence derived last lecture

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

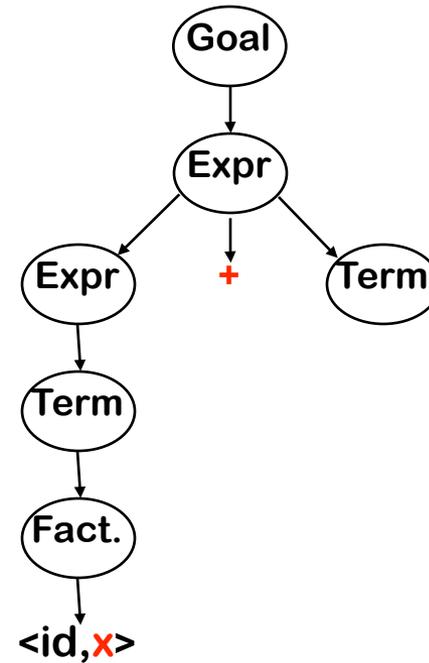
And the input  $x - 2 * y$ .



# Example

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

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



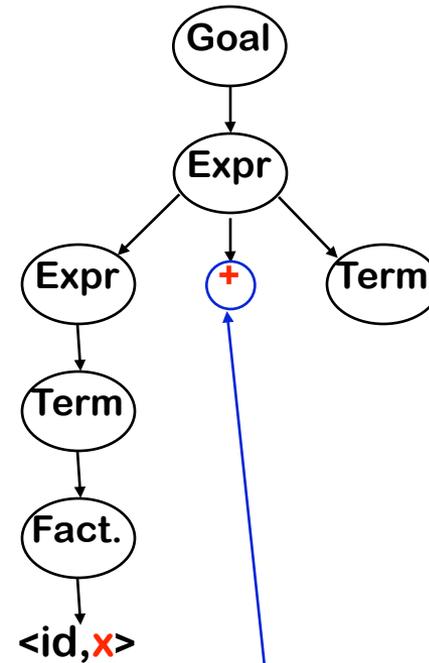
Leftmost derivation, choose productions in an order that exposes problems



# Example

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

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



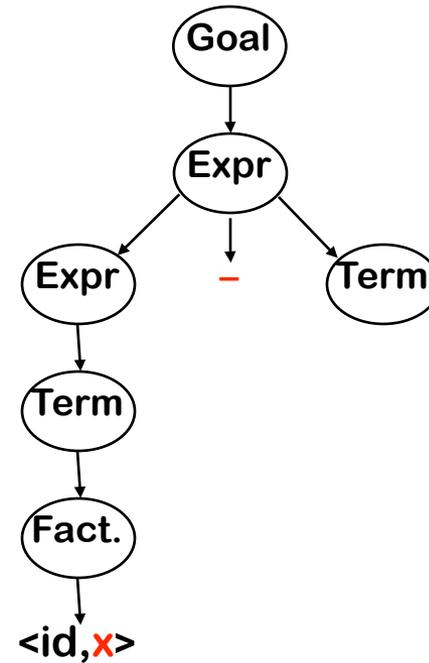
This worked well, except that "-" doesn't match "+"  
The parser must backtrack to here



# Example

Continuing with  $x - 2 * y$ :

<i>Rule</i>	<i>Sentential Form</i>	<i>Input</i>
—	<i>Goal</i>	$\uparrow x - 2 * y$
1	<i>Expr</i>	$\uparrow x - 2 * y$
<hr/>		
3	<i>Expr - Term</i>	$\uparrow x - 2 * y$
4	<i>Term - Term</i>	$\uparrow x - 2 * y$
7	<i>Factor - Term</i>	$\uparrow x - 2 * y$
9	<i>&lt;id,x&gt; - Term</i>	$\uparrow x - 2 * y$
9	<i>&lt;id,x&gt; - Term</i>	$x \uparrow - 2 * y$
—	<i>&lt;id,x&gt; - Term</i>	$x - \uparrow 2 * y$

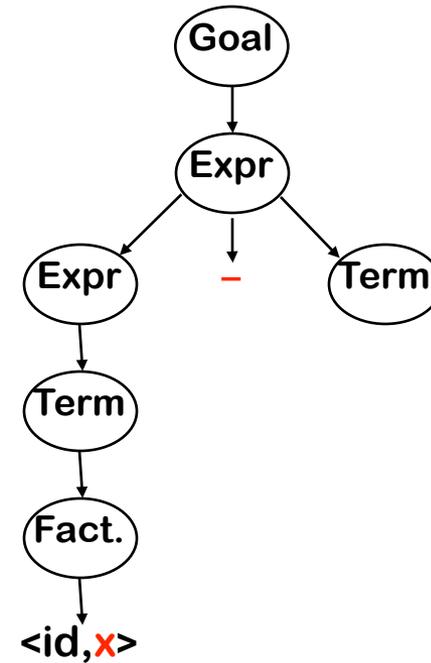




# Example

Continuing with  $x - 2 * y$  :

Rule	Sentential Form	Input
—	<b>Goal</b>	$\uparrow x - 2 * y$
1	<b>Expr</b>	$\uparrow x - 2 * y$
3	<b>Expr - Term</b>	$\uparrow x - 2 * y$
4	<b>Term - Term</b>	$\uparrow x - 2 * y$
7	<b>Factor - Term</b>	$\uparrow x - 2 * y$
9	<b>&lt;id,x&gt; - Term</b>	$\uparrow x - 2 * y$
9	<b>&lt;id,x&gt; - Term</b>	$x \uparrow - 2 * y$
—	<b>&lt;id,x&gt; - Term</b>	$x - \uparrow 2 * y$



This time, “-” and “-” matched

We can advance past “-” to look at “2”

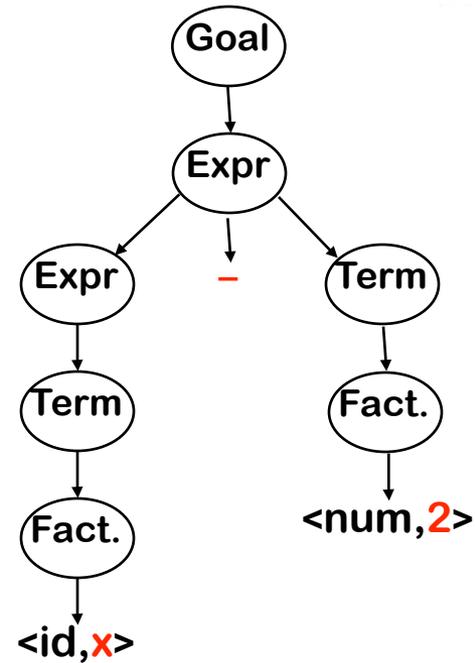
⇒ Now, we need to expand Term - the last NT on the fringe



# Example

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

Rule	Sentential Form	Input
—	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$
7	$\langle id, x \rangle - Factor$	$x - \uparrow 2 * y$
9	$\langle id, x \rangle - \langle num, 2 \rangle$	$x - \uparrow 2 * y$
—	$\langle id, x \rangle - \langle num, 2 \rangle$	$x - 2 \uparrow * y$

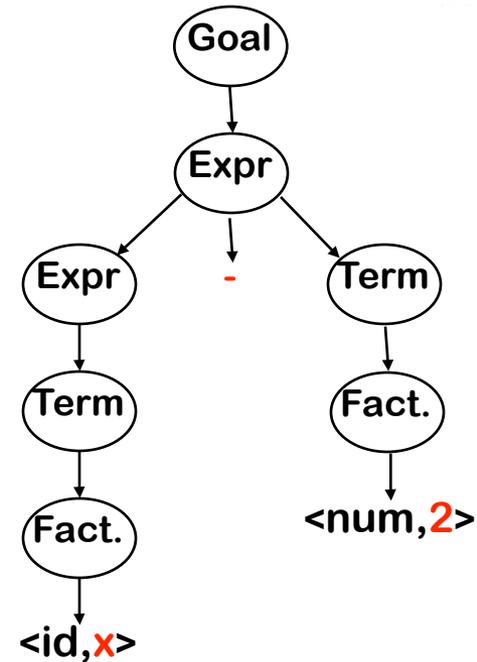




# Example

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

Rule	Sentential Form	Input
—	$\langle \text{id}, x \rangle - \text{Term}$	$x - \uparrow 2 * y$
7	$\langle \text{id}, x \rangle - \text{Factor}$	$x - \uparrow 2 * y$
9	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle$	$x - \uparrow 2 * y$
—	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle$	$x - 2 \uparrow * y$



Where are we?

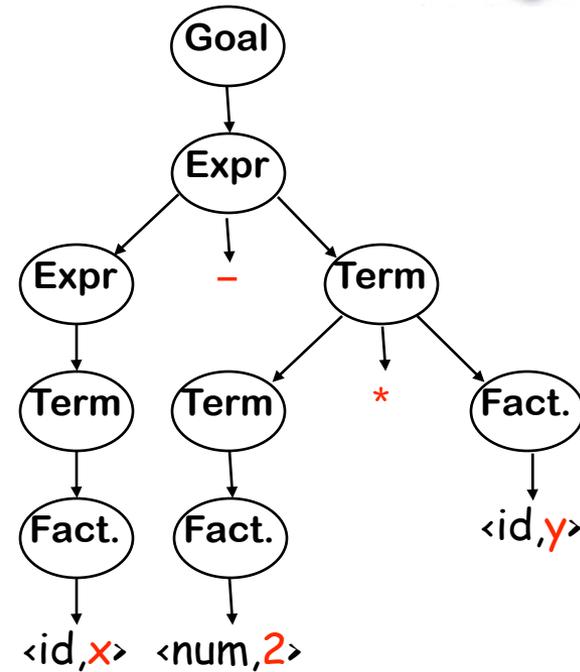
- "2" matches "2"
  - We have more input, but no NTs left to expand
  - The expansion terminated too soon
- ⇒ Need to backtrack



# Example

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

Rule	Sentential Form	Input
—	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$
5	$\langle id, x \rangle - Term * Factor$	$x - \uparrow 2 * y$
7	$\langle id, x \rangle - Factor * Factor$	$x - \uparrow 2 * y$
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	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$x - 2 * \uparrow y$
—	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$x - 2 * y \uparrow$



This time, we matched & consumed all the input

⇒ Success!



# Another possible parse

Other choices for expansion are possible

Rule	Sentential Form	Input
—	Goal	$\uparrow \underline{x} - \underline{2} * \underline{y}$
1	Expr	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	Expr + Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	Expr + Term + Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	Expr + Term + Term + Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	Expr + Term + Term + ... + Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$

consuming no input !

This doesn't terminate

(obviously)

- 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

# 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 expression grammar is left recursive

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

Non-termination is a bad property in any part of a compiler



# Eliminating Left Recursion

---

To remove left recursion, we can transform the grammar

Consider a grammar fragment of the form

$$\begin{aligned} \text{Fee} &\rightarrow \text{Fee } \alpha \\ &| \beta \end{aligned}$$

where neither  $\alpha$  nor  $\beta$  start with  $\text{Fee}$

We can rewrite this as

$$\begin{aligned} \text{Fee} &\rightarrow \beta \text{Fie} \\ \text{Fie} &\rightarrow \alpha \text{Fie} \\ &| \varepsilon \end{aligned}$$

where  $\text{Fie}$  is a new non-terminal

This accepts the same language, but uses only right recursion



# Eliminating Left Recursion

The expression grammar contains two cases of left recursion

$$\begin{array}{l} \text{Expr} \rightarrow \text{Expr} + \text{Term} \\ \quad | \text{Expr} - \text{Term} \\ \quad | \text{Term} \end{array} \qquad \begin{array}{l} \text{Term} \rightarrow \text{Term} * \text{Factor} \\ \quad | \text{Term} / \text{Factor} \\ \quad | \text{Factor} \end{array}$$

Applying the transformation yields

$$\begin{array}{l} \text{Expr} \rightarrow \text{Term Expr}' \\ \text{Expr}' | + \text{Term Expr}' \\ \quad | - \text{Term Expr}' \\ \quad | \epsilon \end{array} \qquad \begin{array}{l} \text{Term} \rightarrow \text{Factor Term}' \\ \text{Term}' | * \text{Factor Term}' \\ \quad | / \text{Factor Term}' \\ \quad | \epsilon \end{array}$$

These fragments use only right recursion

They retain the original left associativity



# Eliminating Left Recursion

Substituting them back into the grammar yields

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

- This grammar is correct, if somewhat non-intuitive.
- It is left associative, as was the original
- A top-down parser will terminate using it.
- A top-down parser may need to backtrack with it.



# Eliminating Left Recursion

The transformation eliminates immediate left recursion  
What about more general, indirect left recursion ?

The general algorithm:

arrange the NTs into some order  $A_1, A_2, \dots, A_n$

for  $i \leftarrow 1$  to  $n$

for  $s \leftarrow 1$  to  $i - 1$

Must start with 1 to ensure that  
 $A_1 \rightarrow A_1 \beta$  is transformed

replace each production  $A_i \rightarrow A_s \gamma$  with  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ ,  
where  $A_s \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all the current productions for  $A_s$   
eliminate any immediate left recursion on  $A_i$   
using the direct transformation

This assumes that the initial grammar has no cycles ( $A_i \Rightarrow^+ A_i$ ),  
and no epsilon productions

And back



# Eliminating Left Recursion

---

How does this algorithm work?

1. Impose arbitrary order on the non-terminals
2. Outer loop cycles through NT in order
3. Inner loop ensures that a production expanding  $A_i$  has no non-terminal  $A_s$  in its rhs, for  $s < i$
4. Last step in outer loop converts any direct recursion on  $A_i$  to right recursion using the transformation showed earlier
5. New non-terminals are added at the end of the order & have no left recursion

At the start of the  $i^{\text{th}}$  outer loop iteration

For all  $k < i$ , no production that expands  $A_k$  contains a non-terminal

$A_s$  in its rhs, for  $s < k$

# Example

---



- Order of symbols:  $G, E, T$

$G \rightarrow E$

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow E \sim T$

$T \rightarrow \underline{\text{id}}$

# Example

---



- Order of symbols:  $G, E, T$

1.  $A_i = G$

$G \rightarrow E$

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow E \sim T$

$T \rightarrow \underline{id}$



# Example

---

- Order of symbols:  $G, E, T$

1.  $A_j = G$

$G \rightarrow E$

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow E \sim T$

$T \rightarrow \underline{id}$

2.  $A_j = E$

$G \rightarrow E$

$E \rightarrow T E'$

$E' \rightarrow + T E'$

$E' \rightarrow \varepsilon$

$T \rightarrow E \sim T$

$T \rightarrow \underline{id}$



# Example

- Order of symbols:  $G, E, T$

1.  $A_j = G$

$G \rightarrow E$

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow E \sim T$

$T \rightarrow \underline{id}$

2.  $A_j = E$

$G \rightarrow E$

$E \rightarrow TE'$

$E' \rightarrow +TE'$

$E' \rightarrow \varepsilon$

$T \rightarrow E \sim T$

$T \rightarrow \underline{id}$

3.  $A_j = T, A_S = E$

$G \rightarrow E$

$E \rightarrow TE'$

$E' \rightarrow +TE'$

$E' \rightarrow \varepsilon$

$T \rightarrow TE' \sim T$

$T \rightarrow \underline{id}$

Go to  
Algorithm



# Example

- Order of symbols:  $G, E, T$

1.  $A_i = G$

$G \rightarrow E$

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow E \sim T$

$T \rightarrow \underline{\text{id}}$

2.  $A_i = E$

$G \rightarrow E$

$E \rightarrow T E'$

$E' \rightarrow + T E'$

$E' \rightarrow \varepsilon$

$T \rightarrow E \sim T$

$T \rightarrow \underline{\text{id}}$

3.  $A_i = T, A_S = E$

$G \rightarrow E$

$E \rightarrow T E'$

$E' \rightarrow + T E'$

$E' \rightarrow \varepsilon$

$T \rightarrow T E' \sim T$

$T \rightarrow \underline{\text{id}}$

4.  $A_i = T$

$G \rightarrow E$

$E \rightarrow T E'$

$E' \rightarrow + T E'$

$E' \rightarrow \varepsilon$

$T \rightarrow \underline{\text{id}} T'$

$T' \rightarrow E' \sim T T'$

$T' \rightarrow \varepsilon$



# Roadmap (Where are we?)

---

We set out to study parsing

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

# 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



# 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  $\text{FIRST}(\alpha)$  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  $\gamma$

We will defer the problem of how to compute FIRST sets until we look at the LR(1) table construction algorithm



# 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  $\text{FIRST}(\alpha)$  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

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

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



# Predictive Parsing

---

What about  $\epsilon$ -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  $\text{FOLLOW}(\alpha)$ , too

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

- $\text{FIRST}(\alpha) \cup \text{FOLLOW}(\alpha)$ , 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

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

$\text{FOLLOW}(\alpha)$  is the set of all words in the grammar that can legally appear immediately after an  $\alpha$



# 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

$$\text{FIRST}^+(\beta_1) \cap \text{FIRST}^+(\beta_2) \cap \text{FIRST}^+(\beta_3) = \emptyset$$

```
/* find an A */  
if (current_word ∈ FIRST(β1))  
    find a β1 and return true  
else if (current_word ∈ FIRST(β2))  
    find a β2 and return true  
else if (current_word ∈ FIRST(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  $\beta_i$ ” (§ 3.3.4 in EAC)



# Recursive Descent Parsing

Recall the expression grammar, after transformation

1	<i>Goal</i>	→	<i>Expr</i>
2	<i>Expr</i>	→	<i>Term Expr'</i>
3	<i>Expr'</i>	→	+ <i>Term Expr'</i>
4			- <i>Term Expr'</i>
5			$\epsilon$
6	<i>Term</i>	→	<i>Factor Term'</i>
7	<i>Term'</i>	→	* <i>Factor Term'</i>
8			/ <i>Factor Term'</i>
9			$\epsilon$
10	<i>Factor</i>	→	<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.

# Recursive Descent Parsing

(Procedural)



A couple of routines from the expression parser

## Goal()

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

## Expr()

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

looking for EOF,  
found token

## Factor()

```
if (token = Number) then  
  token ← next_token();  
  return true;  
else if (token = Identifier) then  
  token ← next_token();  
  return true;  
else  
  report syntax error;  
  return false;
```

*EPrime*, *Term*, & *TPrime* follow the same basic lines (Figure 3.7, EAC)

looking for Number or Identifier,  
found token instead



# Recursive Descent Parsing

To build a parse tree:

- Augment parsing routines to build nodes
- Pass nodes between routines using a stack
- Node for each symbol on rhs
- Action is to pop rhs nodes, make them children of lhs node, and push this subtree

To build an abstract syntax tree

- Build fewer nodes
- Put them together in a different order

```
Expr()  
  result ← true;  
  if (Term() = false)  
    then return false;  
  else if (EPrime() = false)  
    then result ← false;  
  else  
    build an Expr node  
    pop EPrime node  
    pop Term node  
    make EPrime & Term  
      children of Expr  
    push Expr node  
  return result;
```

Success ⇒ build a piece of the parse tree



# Left Factoring

What if my grammar does not have the LL(1) property?

⇒ Sometimes, we can transform the grammar

## The Algorithm

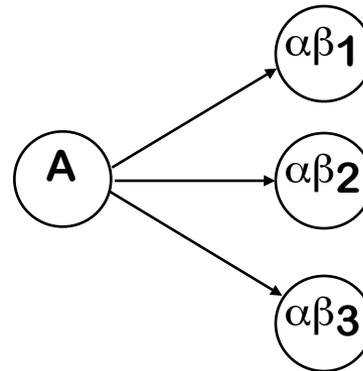
$\forall A \in NT$ ,  
find the longest prefix  $\alpha$  that occurs in two  
or more right-hand sides of  $A$   
if  $\alpha \neq \epsilon$  then replace all of the  $A$  productions,  
 $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma$ ,  
with  
 $A \rightarrow \alpha Z \mid \gamma$   
 $Z \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$   
where  $Z$  is a new element of  $NT$   
Repeat until no common prefixes remain



# Left Factoring

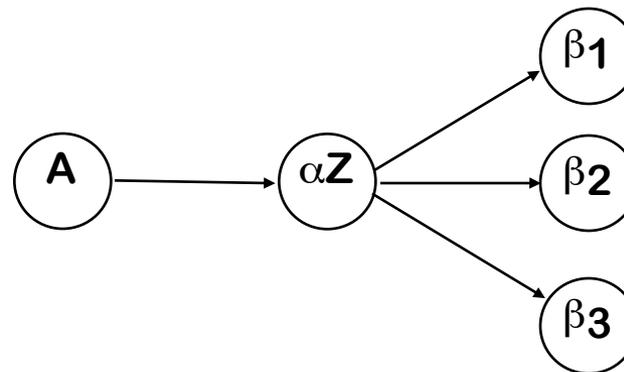
A graphical explanation for the same idea

$A \rightarrow \alpha\beta_1$   
|  $\alpha\beta_2$   
|  $\alpha\beta_3$



becomes ...

$A \rightarrow \alpha Z$   
 $Z \rightarrow \beta_1$   
|  $\beta_2$   
|  $\beta_n$



# Left Factoring

(An example)



Consider the following fragment of the expression grammar

*Factor* → Identifier  
| Identifier [ *ExprList* ]  
| Identifier ( *ExprList* )

FIRST(rhs<sub>1</sub>) = { Identifier }  
FIRST(rhs<sub>2</sub>) = { Identifier }  
FIRST(rhs<sub>3</sub>) = { Identifier }

After left factoring, it becomes

*Factor* → Identifier *Arguments*  
*Arguments* → [ *ExprList* ]  
| ( *ExprList* )  
| ε

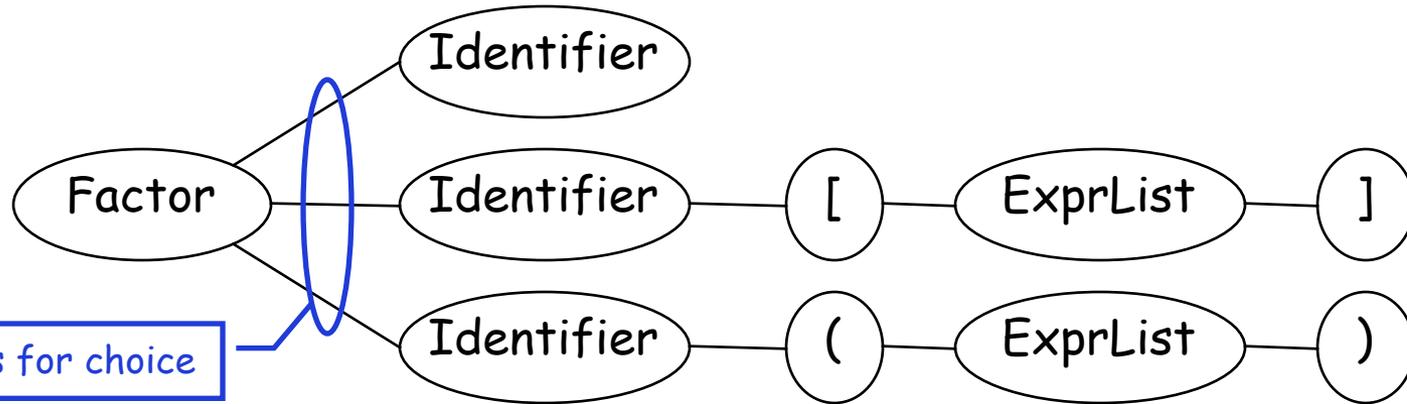
FIRST(rhs<sub>1</sub>) = { Identifier }  
FIRST(rhs<sub>2</sub>) = { [ }  
FIRST(rhs<sub>3</sub>) = { ( }  
FIRST(rhs<sub>4</sub>) = FOLLOW(*Factor*)  
⇒ It has the LL(1) property

This form has the same syntax, with the LL(1) property



# Left Factoring

Graphically



becomes ...

