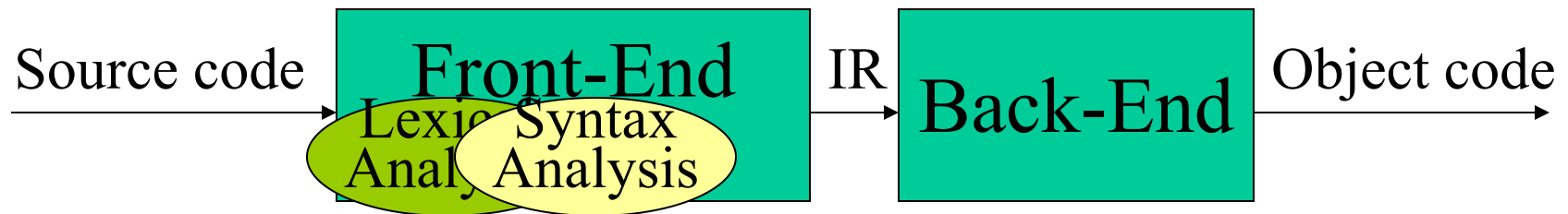


# Lecture 8: Top-Down Parsing



## Parsing:

- Context-free syntax is expressed with a context-free grammar.
- The process of discovering a derivation for some sentence.

## Today's lecture:

Top-down parsing

# Recursive-Descent Parsing

- 1. Construct the root with the starting symbol of the grammar.
- 2. Repeat until the fringe of the parse tree matches the input string:
  - Assuming a node labelled A, select a production with A on its left-hand-side and, for each symbol on its right-hand-side, construct the appropriate child.
  - When a terminal symbol is added to the fringe and it doesn't match the fringe, backtrack.
  - Find the next node to be expanded.

*The key is picking the right production in the first step: that choice should be guided by the input string.*

## Example:

1.  $Goal \rightarrow Expr$

2.  $Expr \rightarrow Expr + Term$

3.       |  $Expr - Term$

4.       |  $Term$

5.  $Term \rightarrow Term * Factor$

6.       |  $Term / Factor$

7.       |  $Factor$

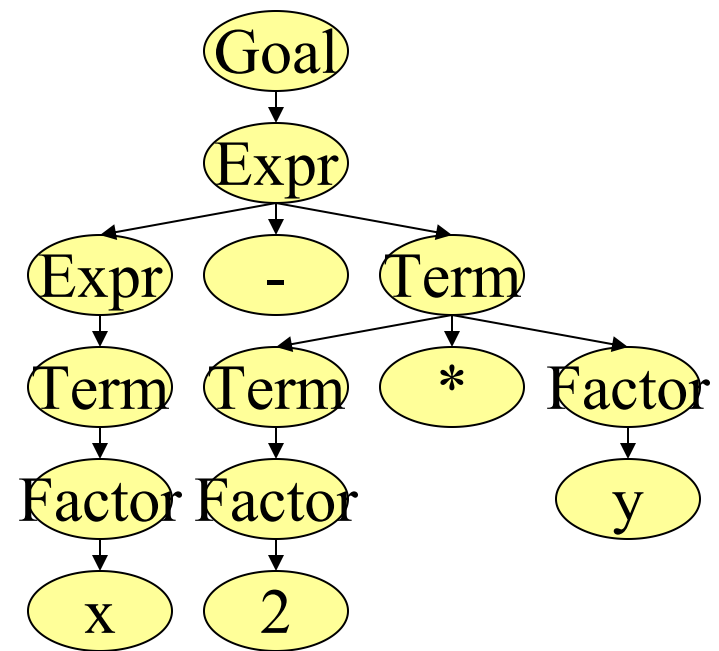
8.  $Factor \rightarrow number$

9.       |  $id$

# Example: Parse $x-2*y$

Steps (one scenario from many)

Rule	Sentential Form	Input
-	<i>Goal</i>	$x-2*y$
1	<i>Expr</i>	$x-2*y$
2	<i>Expr + Term</i>	$x-2*y$
4	<i>Term + Term</i>	$x-2*y$
7	<i>Factor + Term</i>	$x-2*y$
9	<i>id + Term</i>	$x-2*y$
Fail	<i>id + Term</i>	$x$   $-2*y$
Back	<i>Expr</i>	$x-2*y$
3	<i>Expr - Term</i>	$x-2*y$
4	<i>Term - Term</i>	$x-2*y$
7	<i>Factor - Term</i>	$x-2*y$
9	<i>id - Term</i>	$x-2*y$
Match	<i>id - Term</i>	$x$ -   $2*y$
7	<i>id - Factor</i>	$x$ -   $2*y$
9	<i>id - num</i>	$x$ -   $2*y$
Fail	<i>id - num</i>	$x-2$   $*y$
Back	<i>id - Term</i>	$x$ -   $2*y$
5	<i>id - Term * Factor</i>	$x$ -   $2*y$
7	<i>id - Factor * Factor</i>	$x$ -   $2*y$
8	<i>id - num * Factor</i>	$x$ -   $2*y$
match	<i>id - num * Factor</i>	$x-2*$   $y$
9	<i>id - num * id</i>	$x-2*$   $y$
match	<i>id - num * id</i>	$x-2*y$



Other choices for expansion are possible:

Rule	Sentential Form	Input
-	<i>Goal</i>	$x-2*y$
1	<i>Expr</i>	$x-2*y$
2	<i>Expr + Term</i>	$x-2*y$
2	<i>Expr + Term + Term</i>	$x-2*y$
2	<i>Expr + Term + Term + Term + Term</i>	$x-2*y$
2	<i>Expr + Term + Term + ... + Term</i>	$x-2*y$

- Wrong choice leads to non-termination!
- This is a bad property for a parser!
- Parser must make the right choice!

# Left-Recursive Grammars

- **Definition**: A grammar is left-recursive if it has a non-terminal symbol  $A$ , such that there is a derivation  $A \Rightarrow Aa$ , for some string  $a$ .
- A left-recursive grammar can cause a recursive-descent parser to go into an infinite loop.
- **Eliminating left-recursion**: In many cases, it is sufficient to replace  $A \rightarrow Aa \mid b$  with  $A \rightarrow bA'$  and  $A' \rightarrow aA' \mid \varepsilon$
- **Example**:

$$Sum \rightarrow Sum + number \mid number$$

would become:

$$Sum \rightarrow number \ Sum'$$

$$Sum' \rightarrow +number \ Sum' \mid \varepsilon$$

# Eliminating Left Recursion

Applying the transformation to the Grammar of the Example in Slide 2 we get:

$$Expr \rightarrow Term\ Expr'$$

$$Expr' \rightarrow +Term\ Expr' \mid -Term\ Expr' \mid \varepsilon$$

$$Term \rightarrow Factor\ Term'$$

$$Term' \rightarrow *Factor\ Term' \mid /Factor\ Term' \mid \varepsilon$$

( $Goal \rightarrow Expr$  and  $Factor \rightarrow number \mid id$  remain unchanged)

Non-intuitive, but it works!

General algorithm: works for non-cyclic, no  $\varepsilon$ -productions grammars

1. Arrange the non-terminal symbols in order:  $A_1, A_2, A_3, \dots, A_n$

2. For  $i=1$  to  $n$  do

for  $j=1$  to  $i-1$  do

I) replace each production of the form  $A_i \rightarrow A_j \gamma$  with

the productions  $A_i \rightarrow \bar{\delta}_1 \gamma \mid \bar{\delta}_2 \gamma \mid \dots \mid \bar{\delta}_k \gamma$

where  $A_j \rightarrow \bar{\delta}_1 \mid \bar{\delta}_2 \mid \dots \mid \bar{\delta}_k$  are all the current  $A_j$  productions

II) eliminate the immediate left recursion among the  $A_i$

# Where are we?

- We can produce a top-down parser, but:
  - if it picks the wrong production rule it has to backtrack.
- **Idea**: look ahead in input and use context to pick correctly.
- How much lookahead is needed?
  - In general, an arbitrarily large amount.
  - Fortunately, most programming language constructs fall into subclasses of context-free grammars that can be parsed with limited lookahead.

# Predictive Parsing

- Basic idea:
  - For any production  $A \rightarrow a \mid b$  we would like to have a distinct way of choosing the correct production to expand.
- *FIRST* sets:
  - For any symbol  $A$ ,  $FIRST(A)$  is defined as the set of terminal symbols that appear as the first symbol of one or more strings derived from  $A$ .  
E.g. (grammar in Slide 5):  $FIRST(Expr') = \{+, -, \varepsilon\}$ ,  $FIRST(Term') = \{*, /, \varepsilon\}$ ,  
 $FIRST(Factor) = \{number, id\}$
- The LL(1) property:
  - If  $A \rightarrow a$  and  $A \rightarrow b$  both appear in the grammar, we would like to have:  $FIRST(a) \cap FIRST(b) = \emptyset$ . This would allow the parser to make a correct choice with a lookahead of exactly one symbol!

*The Grammar of Slide 5 has this property!*

# Recursive Descent Predictive Parsing

(a practical implementation of the Grammar in Slide 5)

```
Main()
  token=next token();
  if (Expr() != false)
    then <next compilation_step>
  else return False;
```

```
Expr()
  if (Term() == false)
    then result=false
  else if (EPrime() == false)
    then result=false
  else result=true
  return result
```

```
EPrime()
  if (token == '+' or '-') then
    token=next token()
    if (Term() == false)
      then result=false
    elseif (EPrime() == false)
      then result=false
    else result=true
  else result=true /* ε */
  return result
```

```
Term()
  if (Factor() == false)
    then result=false
  else if (TPrime() == false)
    then result=false
  else result=true
  return result
```

```
TPrime()
  if (token == '*' or '/') then
    token=next token()
    if (Factor() == false)
      then result=false
    else if (TPrime() == false)
      then result=false
    else result=true
  else result=true
  return result
```

```
Factor()
  if (token == 'number' or 'id') then
    token=next token()
    result=true
  else
    report syntax error
    result=false
  return result
```

**No backtracking is needed!**  
check :-)



# Left Factoring

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

Sometimes, we can transform a grammar to have this property.

## Algorithm:

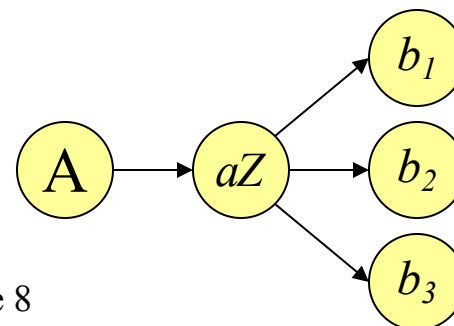
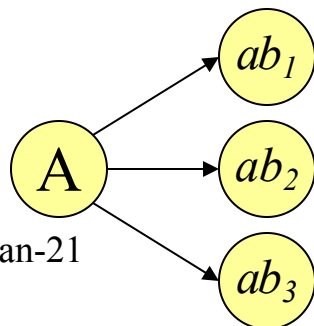
1. For each non-terminal  $A$ , find the longest prefix, say  $a$ , common to two or more of its alternatives

2. if  $a \neq \epsilon$  then replace all the  $A$  productions,  $A \rightarrow ab_1 | ab_2 | ab_3 | \dots | ab_n | \gamma$ , where  $\gamma$  is anything that does not begin with  $a$ , with  $A \rightarrow aZ | \gamma$  and  $Z \rightarrow b_1 | b_2 | b_3 | \dots | b_n$

Repeat the above until no common prefixes remain

**Example:**  $A \rightarrow ab_1 | ab_2 | ab_3$  would become  $A \rightarrow aZ$  and  $Z \rightarrow b_1 | b_2 | b_3$

Note the graphical representation:



# Example

(NB: this is a different grammar from the one in Slide 2)

$Goal \rightarrow Expr$

$Expr \rightarrow Term + Expr$   
 $\quad | Term - Expr$   
 $\quad | Term$

$Term \rightarrow Factor * Term$

$\quad | Factor / Term$   
 $\quad | Factor$

$Factor \rightarrow number$   
 $\quad | id$

We have a problem with the different rules for  $Expr$  as well as those for  $Term$ . In both cases, the first symbol of the right-hand side is the same ( $Term$  and  $Factor$ , respectively). E.g.:

$FIRST(Term) = FIRST(Term) \cap FIRST(Term) = \{number, id\}.$

$FIRST(Factor) = FIRST(Factor) \cap FIRST(Factor) = \{number, id\}.$

## Applying left factoring:

$Expr \rightarrow Term Expr'$

$Expr' \rightarrow + Expr \mid - Expr \mid \varepsilon$

$Term \rightarrow Factor Term'$

$Term' \rightarrow * Term \mid / Term \mid \varepsilon$

$FIRST(+) = \{+\}; FIRST(-) = \{-\}; FIRST(\varepsilon) = \{\varepsilon\};$   
 $FIRST(-) \cap FIRST(+) \cap FIRST(\varepsilon) = \emptyset$

$FIRST(*) = \{*\}; FIRST(/) = \{/ \}; FIRST(\varepsilon) = \{\varepsilon\};$   
 $FIRST(*) \cap FIRST(/) \cap FIRST(\varepsilon) = \emptyset$

# Example (cont.)

1.  $Goal \rightarrow Expr$
2.  $Expr \rightarrow Term Expr'$
3.  $Expr' \rightarrow + Expr$
4.       |  $- Expr$
5.       |  $\varepsilon$
6.  $Term \rightarrow Factor Term'$
7.  $Term' \rightarrow * Term$
8.       |  $/ Term$
9.       |  $\varepsilon$
10.  $Factor \rightarrow number$
11.       |  $id$

Rule	Sentential Form	Input
-	<i>Goal</i>	x - 2*y
1	<i>Expr</i>	x - 2*y
2	<i>Term Expr'</i>	x - 2*y
6	<i>Factor Term' Expr'</i>	x - 2*y
11	<i>id Term' Expr'</i>	x - 2*y
Match	<i>id Term' Expr'</i>	x   - 2*y
9	<i>id <math>\varepsilon</math> Expr'</i>	x   - 2*y
4	<i>id - Expr</i>	x   - 2*y
Match	<i>id - Expr</i>	x -   2*y
2	<i>id - Term Expr'</i>	x -   2*y
6	<i>id - Factor Term' Expr'</i>	x -   2*y
10	<i>id - num Term' Expr'</i>	x -   2*y
Match	<i>id - num Term' Expr'</i>	x - 2   *y
7	<i>id - num * Term Expr'</i>	x - 2   *y
Match	<i>id - num * Term Expr'</i>	x - 2*   y
6	<i>id - num * Factor Term' Expr'</i>	x - 2*   y
11	<i>id - num * id Term' Expr'</i>	x - 2*   y
Match	<i>id - num * id Term' Expr'</i>	x - 2*y
9	<i>id - num * id Expr'</i>	x - 2*y
5	<i>id - num * id</i>	x - 2*y

The next symbol determines each choice correctly. No backtracking needed.

# Conclusion

- Top-down parsing:
  - recursive with backtracking (not often used in practice)
  - recursive predictive
- Nonrecursive Predictive Parsing is possible too: maintain a stack explicitly rather than implicitly via recursion and determine the production to be applied using a table (Aho, pp.186-190).
- Given a Context Free Grammar that doesn't meet the LL(1) condition, it is undecidable whether or not an equivalent LL(1) grammar exists.
- Next time: Bottom-Up Parsing
- Reading: Aho2, Sections 4.3.3, 4.3.4, 4.4; Aho1, pp. 176-178, 181-185; Grune pp.117-133; Hunter pp. 72-93; Cooper, Section 3.3.