

# Top-Down Parsing

# Parsing:

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

# 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: Parse $x-2*y$

## Example:

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

| 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    |

# Example: Parse $x-2*y$

Example:

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

| <b>Rule</b> | <b>Sentential Form</b>                 | <b>Input</b> |
|-------------|--|--------------|
| -           | <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</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

## Example:

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

Applying the transformation to the Grammar of the Example 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!

# 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  $A \rightarrow a / b$  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.  $Expr \rightarrow Term Expr'$

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

$Term \rightarrow Factor Term'$

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

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

$FIRST(Expr') = \{+, -, \epsilon\}$ ,  $FIRST(Term') = \{*, /, \epsilon\}$ ,  $FIRST(Factor) = \{number, id\}$

# 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!

# Left Factoring

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

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

## Algorithm:

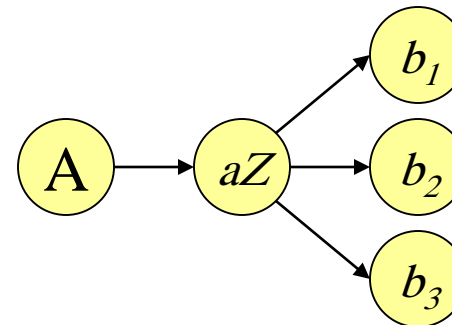
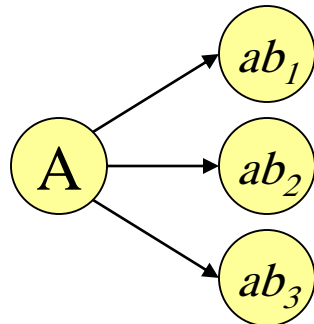
1. For each non-terminal  $A$ , find the longest prefix, say  $a$ , common to two or more of its alternatives

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

*Goal*  $\rightarrow$  *Expr*

*Expr*  $\rightarrow$  *Term* + *Expr*  
/ *Term* - *Expr*  
/ *Term*

*Term*  $\rightarrow$  *Factor* \* *Term*

/ *Factor* / *Term*  
/ *Factor*  
*Factor*  $\rightarrow$  *number*  
/ *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* / - *Expr* /  $\epsilon$

*Term*  $\rightarrow$  *Factor Term'*  
*Term'*  $\rightarrow$  \* *Term* / / *Term* /  $\epsilon$

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

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



# 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.      $/ id$

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

| 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 ε 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    |



# 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