

# Lecture 13

## Top-Down Parser

A top-down parser starts with the root of the parse tree. The root node is labeled with the goal (start) symbol of the grammar. The top-down parsing algorithm proceeds as follows:

1. Construct the root node of the parse tree
2. Repeat until the fringe of the parse tree matches input string
  - a. At a node labeled  $A$ , select a production with  $A$  on its lhs
  - b. for each symbol on its rhs, construct the appropriate child
  - c. When a terminal symbol is added to the fringe and it does not match the fringe, backtrack

The key is picking right production in step a. That choice should be guided by the input string. Let's try parsing using this algorithm using the expression grammar.

$$x - 2 * y$$

$P$	Sentential Form	input
-	Goal	- <u>x</u> - <u>2</u> * y
1	expr	- <u>x</u> - <u>2</u> * y
2	expr + term	- <u>x</u> - <u>2</u> * y
4	term + term	- <u>x</u> - <u>2</u> * y
7	factor + term	- <u>x</u> - <u>2</u> * y
9	<id,x> + term	- <u>x</u> - <u>2</u> * y
9	<id,x> + term	x - <u>-</u> <u>2</u> * y

This worked well except that “-” does not match “+”. The parser made the wrong choice of production to use at step 2. The parser must backtrack and use a different production.

$P$	Sentential Form	input
-	Goal	- <u>x</u> - <u>2</u> * y
1	expr	- <u>x</u> - <u>2</u> * y
3	expr - term	- <u>x</u> - <u>2</u> * y
4	term - term	- <u>x</u> - <u>2</u> * y
7	factor - term	- <u>x</u> - <u>2</u> * y
9	<id,x> - term	- <u>x</u> - <u>2</u> * y
9	<id,x> - term	x - <u>-</u> <u>2</u> * y

This time the “-” and “-” matched. We can advance past “-” to look at “2”. Now, we need to expand “term”

<i>P</i>	<i>Sentential Form</i>	<i>input</i>
-	$\langle \text{id}, \underline{x} \rangle - \text{term}$	$\underline{x} - - \underline{2} * \underline{y}$
7	$\langle \text{id}, \underline{x} \rangle - \text{factor}$	$\underline{x} - - \underline{2} * \underline{y}$
9	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, 2 \rangle$	$\underline{x} - - \underline{2} * \underline{y}$
-	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle$	$\underline{x} - \underline{2} * \underline{y}$

The 2's match but the expansion terminated too soon because there is still unconsumed input and there are no non-terminals to expand in the sentential form  $\Rightarrow$  Need to backtrack.

<i>P</i>	<i>Sentential Form</i>	<i>input</i>
-	$\langle \text{id}, \underline{x} \rangle - \text{term}$	$\underline{x} - - \underline{2} * \underline{y}$
5	$\langle \text{id}, \underline{x} \rangle - \text{term} * \text{factor}$	$\underline{x} - - \underline{2} * \underline{y}$
7	$\langle \text{id}, \underline{x} \rangle - \text{factor} * \text{factor}$	$\underline{x} - - \underline{2} * \underline{y}$
8	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, 2 \rangle * \text{factor}$	$\underline{x} - - \underline{2} * \underline{y}$
-	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, 2 \rangle * \text{factor}$	$\underline{x} - \underline{2} - * \underline{y}$
-	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, 2 \rangle * \text{factor}$	$\underline{x} - \underline{2} * - \underline{y}$
9	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, 2 \rangle * \langle \text{id}, \underline{y} \rangle$	$\underline{x} - \underline{2} * - \underline{y}$
-	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, 2 \rangle * \langle \text{id}, \underline{y} \rangle$	$\underline{x} - \underline{2} * \underline{y} -$

This time the parser met with success. All of the input matched.

## Left Recursion

Consider another possible parse:

<i>P</i>	<i>Sentential Form</i>	<i>input</i>
-	<i>Goal</i>	$\underline{x} - \underline{2} * \underline{y}$
1	<i>expr</i>	$\underline{x} - \underline{2} * \underline{y}$
2	<i>expr + term</i>	$\underline{x} - \underline{2} * \underline{y}$
2	<i>expr + term + term</i>	$\underline{x} - \underline{2} * \underline{y}$
2	<i>expr + term + term + term</i>	$\underline{x} - \underline{2} * \underline{y}$
2	<i>expr + term + term + term + ....</i>	$\underline{x} - \underline{2} * \underline{y}$

Parser is using productions but no input is being consumed.

Top-down parsers cannot handle left-recursive grammars. Formally, a grammar is left recursive if  $\exists A \in NT$  such that  $\exists$  a derivation  $A \Rightarrow^* A a$ , for some string  $a \in (NT \cup T)^*$ .

Our expression grammar is left recursive. This can lead to non-termination in a top-down parser. Non-termination is bad in any part of a compiler! For a top-down parser, any recursion must be a *right recursion*. We would like to convert left recursion to right. To remove left recursion, we can transform the grammar. Consider a grammar fragment:

$$A \rightarrow A a \mid b$$

where neither  $a$  nor  $b$  starts with  $A$ . We can rewrite this as:

$$A \rightarrow b A'$$

$$A' \rightarrow a A' \mid e$$

where  $A'$  is a new non-terminal. This grammar accepts the same language but uses only right recursion. The expression grammar we have been using contains two cases of left-recursion. Applying the transformation yields

$$\begin{array}{l} \text{expr} \rightarrow \text{term expr}' \\ \text{expr}' \rightarrow + \text{term expr}' \\ \quad \quad \quad | - \text{term expr}' \\ \\ \text{term} \rightarrow \text{factor term}' \\ \text{term}' \rightarrow * \text{factor term}' \\ \quad \quad \quad | / \text{factor term}' \\ \quad \quad \quad | \mathbf{e} \end{array}$$

These fragments use only right recursion. They retain the original left associativity. A top-down parser will terminate using them.

## Predictive Parsing

If a top down parser picks the wrong production, it may need to backtrack.

Alternative is to *look-ahead* in input and use context to pick the production to use correctly. How much look-ahead is needed? In general, an arbitrarily large amount of look-ahead symbols are required.. Fortunately, large classes of CFGs can be parsed with limited lookahead. Most programming languages constructs fall in those subclasses

The basic idea in predictive parsing is: given  $A \rightarrow a \mid b$ , the parser should be able to choose between  $a$  and  $b$ . To accomplish this, the parser needs FIRST and FOLLOW sets.

**Definition:** FIRST sets: for some rhs  $a \in G$ , define  $\text{FIRST}(a)$  as the set of tokens that appear as the first symbol in some string that derives from  $a$ . That is,  $x \in \text{FIRST}(a)$  iff  $a \xrightarrow{P} xg$ , for some  $g$ .