

BNF, EBNF, and Parse Trees

BNF

We have seen previously seen the notation for formal grammars. However, that notation does not work well in practical grammars where elements of the grammar need to be defined not by single characters by character strings. Moreover, the “ \rightarrow ” is not easy to type in. A common alternative notation is BNF (Backus-Naur Form) as in

$$\begin{aligned} \langle \text{if stmt} \rangle &::= && \text{if } \langle \text{logic expr} \rangle \text{ then } \langle \text{stmt} \rangle \mid \\ &&& \text{if } \langle \text{logic expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \end{aligned}$$

In that notation “ $::=$ ” is used instead of “ \rightarrow ” and the variables (i.e. non-terminals) are denoted by angle brackets ($\langle \dots \rangle$). The terminal symbols are denoted without any special markings (but generally bolded in my notes). As we have seen previously, “ $|$ ” is used to separate rules for the same variable.

We also saw that derivations describe how a sentence can be generated by the grammar.

Let us consider the following grammar:

$$\begin{aligned} \langle \text{assign} \rangle &::= && \langle \text{id} \rangle = \langle \text{expr} \rangle \\ \langle \text{id} \rangle &::= && \mathbf{A} \mid \mathbf{B} \mid \mathbf{C} \\ \langle \text{expr} \rangle &::= && \langle \text{id} \rangle + \langle \text{expr} \rangle \mid \langle \text{id} \rangle * \langle \text{expr} \rangle \mid (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle \end{aligned}$$

The leftmost derivation for the following sentence:

$\mathbf{A = B * (A + C)}$ is

$$\begin{aligned} \langle \text{assign} \rangle &\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle \Rightarrow \mathbf{A} = \langle \text{expr} \rangle \Rightarrow \mathbf{A} = \langle \text{id} \rangle * \langle \text{expr} \rangle \Rightarrow \\ &\mathbf{A = B * } \langle \text{expr} \rangle \Rightarrow \mathbf{A = B * (} \langle \text{expr} \rangle \text{)} \Rightarrow \mathbf{A = B * (} \langle \text{id} \rangle + \langle \text{expr} \rangle \text{)} \Rightarrow \\ &\mathbf{A = B * (A + } \langle \text{expr} \rangle \text{)} \Rightarrow \mathbf{A = B * (A + } \langle \text{id} \rangle \text{)} \Rightarrow \mathbf{A = B * (A + C)} \end{aligned}$$

while the rightmost derivation is:

$$\begin{aligned} \langle \text{assign} \rangle &\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle \Rightarrow \langle \text{id} \rangle = \langle \text{id} \rangle * \langle \text{expr} \rangle \Rightarrow \\ &\langle \text{id} \rangle = \langle \text{id} \rangle * (\langle \text{expr} \rangle) \Rightarrow \langle \text{id} \rangle = \langle \text{id} \rangle * (\langle \text{id} \rangle + \langle \text{expr} \rangle) \Rightarrow \langle \text{id} \rangle = \\ &\langle \text{id} \rangle * (\langle \text{id} \rangle + \langle \text{id} \rangle) \Rightarrow \langle \text{id} \rangle = \langle \text{id} \rangle * (\langle \text{id} \rangle + \mathbf{C}) \Rightarrow \\ &\langle \text{id} \rangle = \langle \text{id} \rangle * (\mathbf{A + C}) \Rightarrow \langle \text{id} \rangle = \mathbf{B * (A + C)} \Rightarrow \mathbf{A = B * (A + C)} \end{aligned}$$

Notice the arrow we use at each step of the derivation. It is “ \Rightarrow ” which is different from the arrow used in productions.

The notation \Rightarrow^* denotes multiple steps. As in

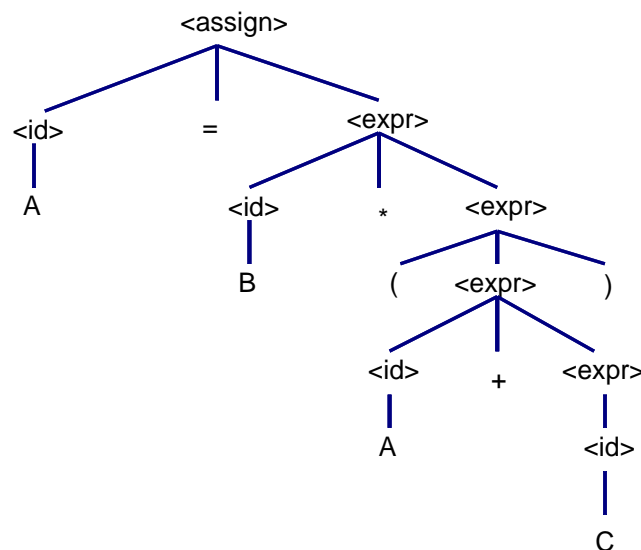
$$\langle \text{assign} \rangle \Rightarrow^* \langle \text{id} \rangle = \langle \text{id} \rangle * (\langle \text{id} \rangle + \mathbf{C})$$

Think of it as $\langle \text{assign} \rangle \Rightarrow \dots \Rightarrow \langle \text{id} \rangle = \langle \text{id} \rangle * (\langle \text{id} \rangle + \mathbf{C})$

We will now examine another way to describe the structure of a sentence.

Parse Tree

Another, somewhat more convenient, way to display the same information is via a parse tree. Parse trees have the advantage that, since they are two dimensional, the order in which the rules are applied is not reflected in them. The parse tree for the sentence above is:



Some things to notice about parse trees:

- The start symbol is always at the root of the tree
- Every leaf node is a terminal
- Every interior node is a non-terminal
- The sentence appears in a left-to-right traversal of the leaves.
- Each parse tree corresponds to a rightmost derivation and to a left most derivation (i.e. there is a one-to-one correspondence between them)

When designing a grammar one must be careful not to create ambiguity. This is when two or more parse trees (or leftmost derivations, or rightmost derivations) could exist. For example:

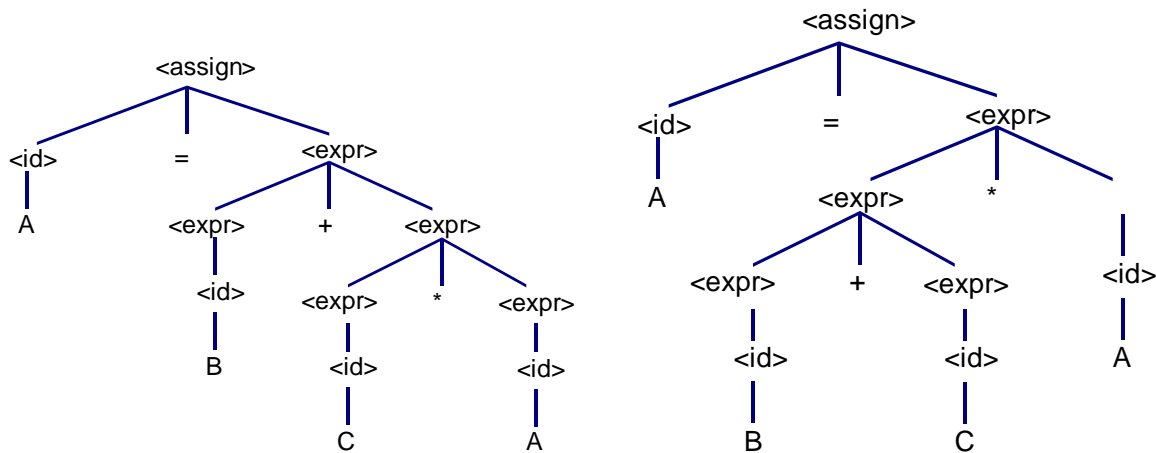
Consider the following grammar:

```

<assign>      ::= <id> = <expr>
<id>          ::= A | B | C
<expr>        ::= <expr> + <expr> | <expr> * <expr> | ( <expr> )
                | <id>

```

and the sentence **A = B + C * A** which has two distinct parse trees (from what we saw earlier it would also have two leftmost derivations and two rightmost derivations):



The left tree corresponds to the interpretation of the string as $A = B + (C * A)$, while the right tree corresponds to $A = (B + C) * A$. As you can see the results will be quite different. Assuming $A = 4$, $B = 3$, and $C = 6$. The first interpretation will set A to 27 and the second to 36. *Which one is the one used by most languages?*

We can rewrite the grammar to reflect the precedence of operators $+$ and $*$ and eliminate this ambiguity

```

<assign>      ::= <id> = <expr>
<id>          ::= A | B | C
<expr>        ::= <expr> + <expr> | <factor>
<factor>      ::= <factor> * <factor> | ( <expr> ) | <id>

```

The string $A = B + C * A$ has now a single parse tree. However, the string $A = B + C + A$ has two parse trees as the above grammar does not reflect the associativity of the $+$ operator. This does not cause serious problems with addition but it would with subtraction. For example:

$5 - 3 - 2$ is equal to 0 if “-” is left associative [i.e. $(5 - 3) - 2$] but is equal to 4 if it is right associative $[5 - (3 - 2)]$.

Note: addition and subtraction are left associative as are division and mod.
power is right associative

The following grammar eliminates this ambiguity:

```

<assign> ::= <id> = <expr>
<id>      ::= A | B | C
<expr>    ::= <expr> + <term> | <term>
<term>    ::= <term> * <factor> | <factor>
<factor>  ::= ( <expr> ) | <id>

```

This grammar reflects left associativity (associated with left recursive rule). Right associativity is associated with right recursion.

Extended BNF (EBNF)

Three extensions are commonly included (highlighted in red):

- optional part of an RHS
 $\langle \text{selection} \rangle ::= \text{if } (\langle \text{expression} \rangle) \langle \text{statement} \rangle \text{ [else } \langle \text{statement} \rangle]$
- indefinite repeat
 $\langle \text{ident list} \rangle ::= \langle \text{identifier} \rangle \{ , \langle \text{identifier} \rangle \}$
- multiple choice option
 $\langle \text{term} \rangle ::= \langle \text{term} \rangle (* | / | \%) \langle \text{factor} \rangle$

Note that by using the braces you may lose some information:

$\langle \text{expr} \rangle ::= \langle \text{term} \rangle \{ + \langle \text{term} \rangle \}$

does not imply left associativity the way the BNF rule

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle$

does. This can be addressed by enforcing the correct associativity during the semantic analysis phase.

	BNF	EBNF
1	$\langle a \rangle ::= A \mid AB$	$\langle a \rangle ::= A [B]$
2	$\langle snum \rangle ::= - \langle num \rangle \mid \langle num \rangle$	$\langle snum \rangle ::= [-] \langle num \rangle$
3	$\langle a \rangle ::= \langle a \rangle A \mid A$	$\langle a \rangle ::= A \{ A \}$
4	$\langle a \rangle ::= \langle a \rangle A \mid \lambda$	$\langle a \rangle ::= \{ A \}$
5	$\langle block \rangle ::= \text{begin } \langle stmts \rangle \text{ end}$ $\langle stmts \rangle ::= \langle cmd \rangle \mid \langle cmd \rangle ; \langle stmts \rangle$	$\langle block \rangle ::= \text{begin } \langle cmd \rangle \{ ; \langle cmd \rangle \}$ end
6	$\langle snum \rangle ::= \langle num \rangle \mid + \langle num \rangle \mid - \langle num \rangle$	$\langle snum \rangle ::= [(+ \mid -)] \langle num \rangle$

EBNF has the effect of

- removing most “or’s”(|) which reduces the number of rules
- removing most recursion which is replaced by { } loops
- removing occurrences of the null string (λ)

It is also used to build syntax diagrams.