

Chapter 3: Describing Syntax and Semantics

- Introduction
- Formal methods of describing syntax (BNF)

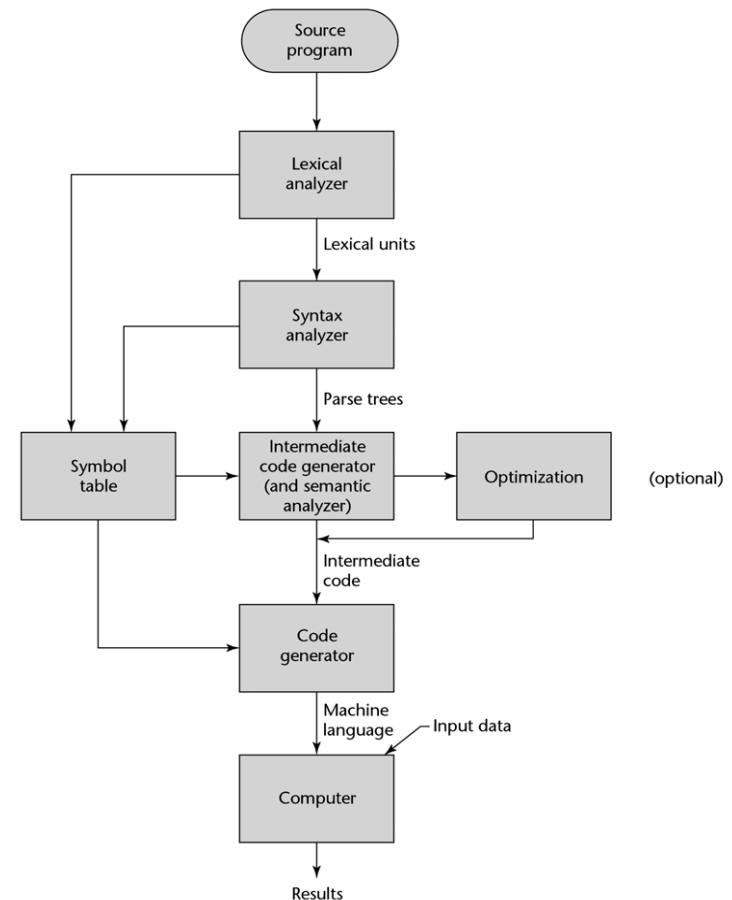
We can analyze **syntax** of a computer program on two levels:

1. Lexical level
2. Syntactic level

- Lexical analyzer collect characters into tokens.
- Syntactic analyzer determine syntax structure and determine whether the given programs are syntactically correct.

Figure 1.3

The compilation process



Derivation

- A derivation is a repeated application of rules, starting with the start symbol and ending with a sentence (all terminal symbols)

```

<program> => <stmts>
=> <stmt>
=> <var> = <expr>
=> a = <expr>
=> a = <term> + <term>
=> a = <var> + <term>
=> a = b + <term>
=> a = b + const

```

```

<program> → <stmts>
<stmts> → <stmt> | <stmt> ; <stmts>
<stmt> → <var> = <expr>
<var> → a | b | c | d
<expr> → <term> + <term> | <term> -
                                     <term>
<term> → <var> | const

```

Parse Tree: A hierarchical representation of a derivation

For example $x + y + z$ is parsed:

$\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle \langle \text{binary} \rangle \langle \text{exp} \rangle$

$\langle \text{exp} \rangle \rightarrow \langle \text{identifier} \rangle$

 | $\langle \text{literal} \rangle$

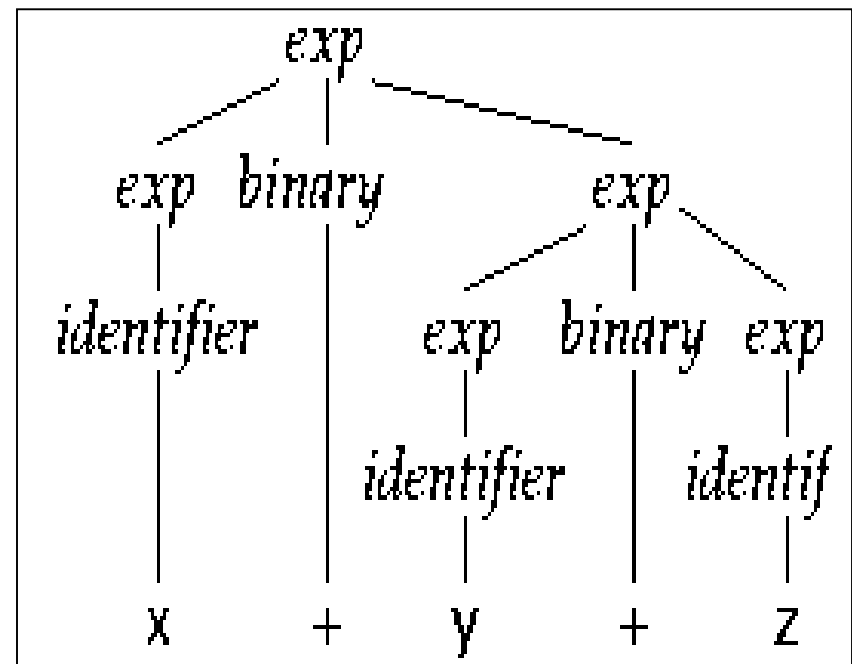
 | $\langle \text{unary} \rangle \langle \text{exp} \rangle$

 | $\langle \text{exp} \rangle \langle \text{binary} \rangle \langle \text{exp} \rangle$

$\langle \text{binary} \rangle \rightarrow '<' \mid '>' \mid '+' \mid '-'$

$\langle \text{unary} \rangle \rightarrow '-'$

$\langle \text{identifier} \rangle \rightarrow x|y|z$



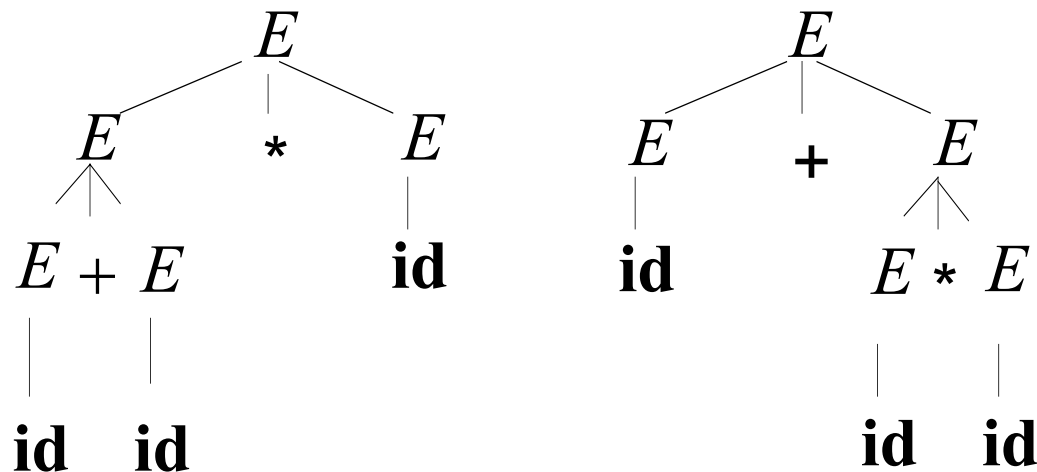
An Ambiguous Expression Grammar

A grammar that can have more than one parse tree generating a particular string is **ambiguous**.

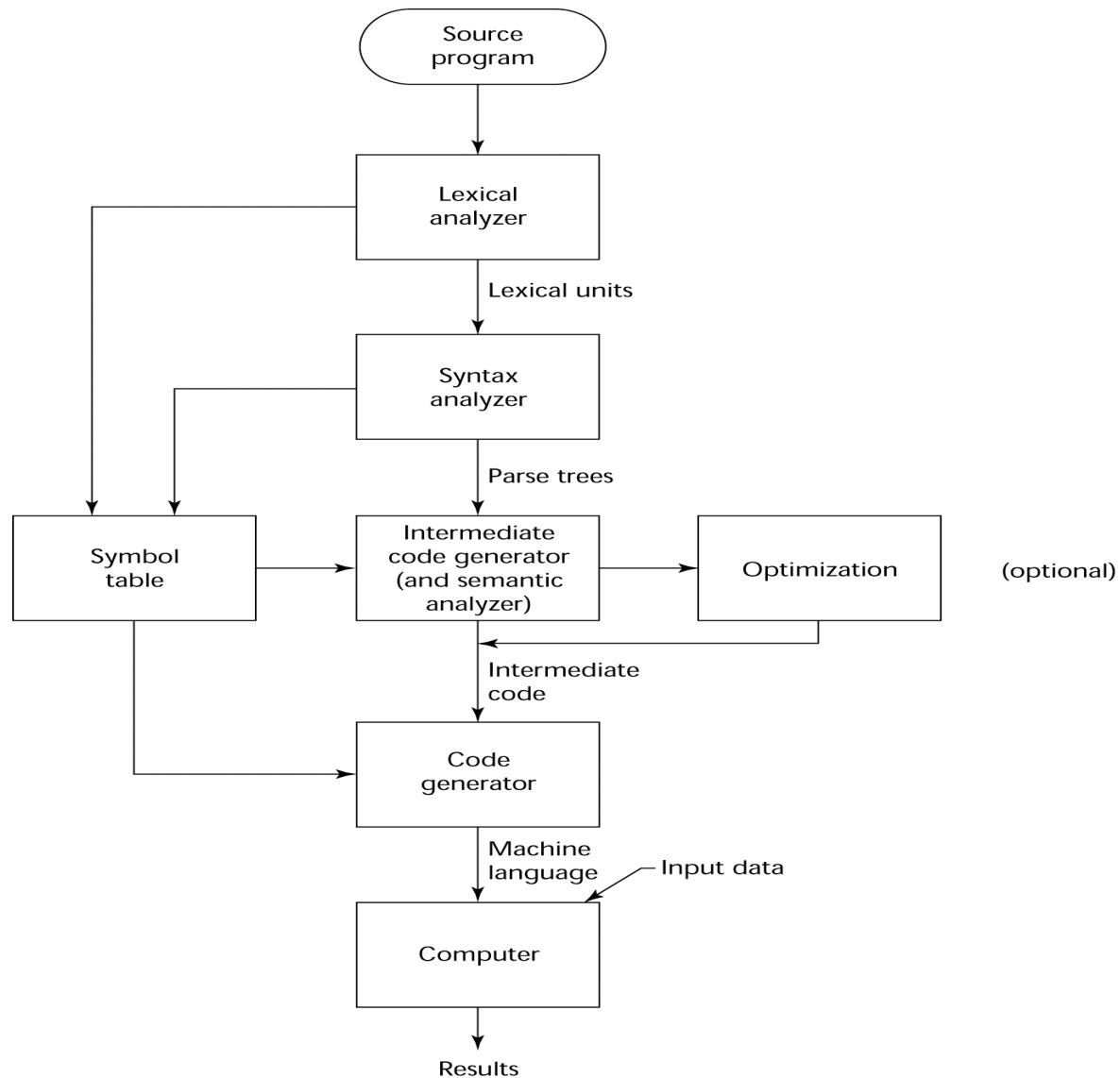
$\langle E \rangle \rightarrow \langle E \rangle + \langle E \rangle \mid \langle E \rangle - \langle E \rangle \mid \langle E \rangle * \langle E \rangle \mid \langle E \rangle /$

$\langle E \rangle \mid \text{id}$

There are two parse trees for the expression $\text{id} + \text{id} * \text{id}$



Ambiguity is a problem: compiler chooses the code to be generated for a statement by examining its parse tree. If more than one parse tree, the meaning of the structure cannot be unique.

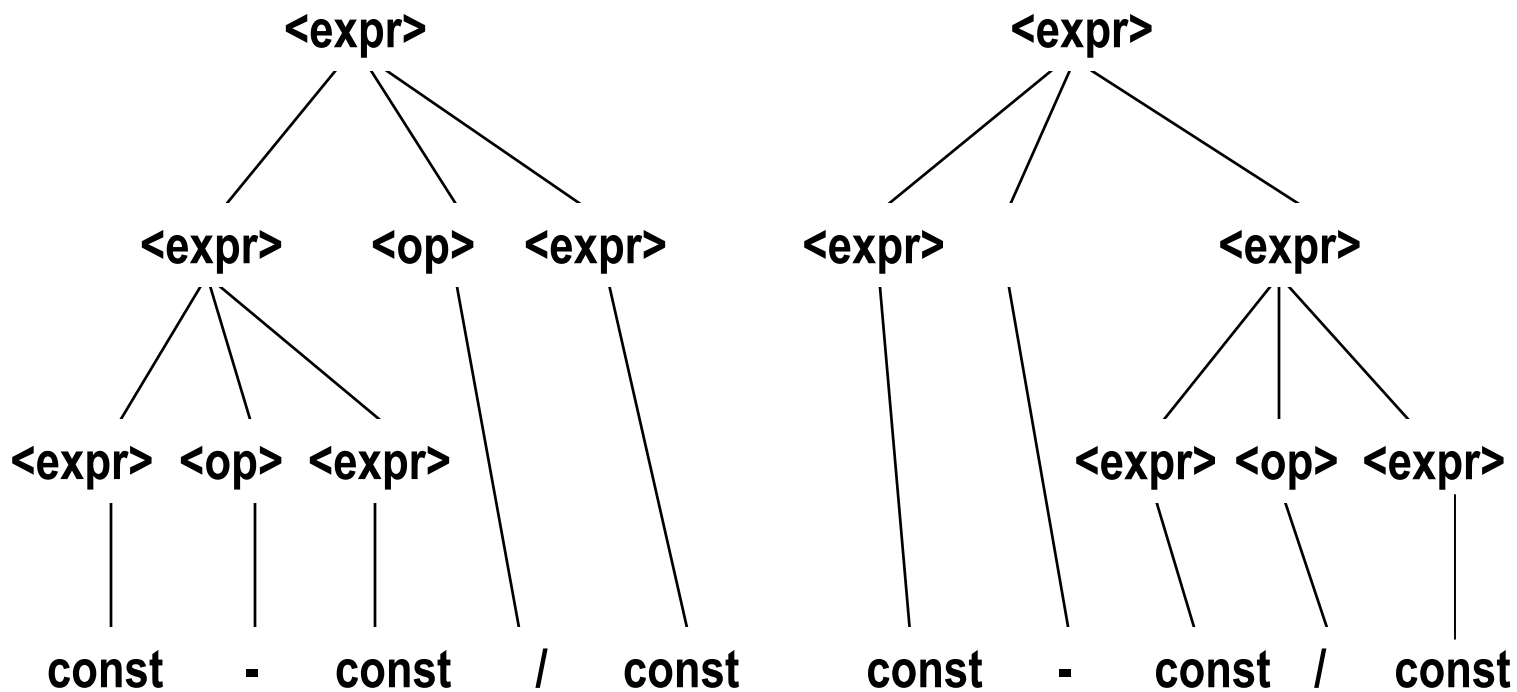


Another example:

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle | \text{const}$

$\langle \text{op} \rangle \rightarrow / | -$

Two parse trees for expression $\text{const} - \text{const} / \text{const}$

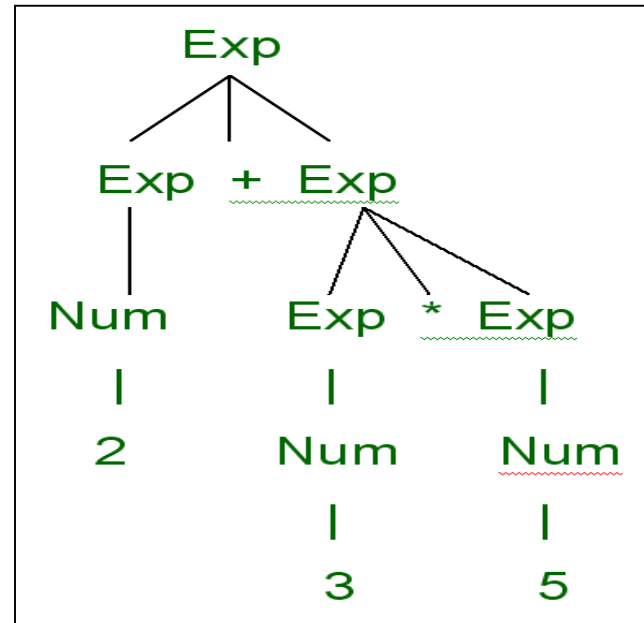
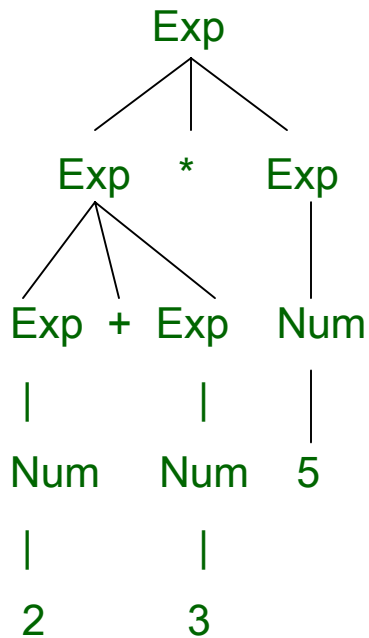


Exercise:

$\langle \text{Exp} \rangle \rightarrow \langle \text{Num} \rangle \mid \langle \text{Exp} \rangle + \langle \text{Exp} \rangle \mid \langle \text{Exp} \rangle * \langle \text{Exp} \rangle$

$\langle \text{Num} \rangle \rightarrow 2 \mid 3 \mid 5$

Two parse trees for string $2 + 3 * 5$



In order to avoid ambiguity, it is essential that the grammar generate only one possible structure for each string in the language.

The ambiguity can be eliminated by imposing the precedence of one operator over the other.

Imposing the precedence of one operator over the other.

We say that **op1** has precedence over **op2** if an expression of the form

$$e_1 \text{ op1 } e_2 \text{ op2 } e_3$$

is interpreted only as

$$(e_1 \text{ op1 } e_2) \text{ op2 } e_3$$

In other words, **op1** binds tighter than **op2**.

From the point of view of derivation trees, the fact that $e_1 \text{ op1 } e_2 \text{ op2 } e_3$ is interpreted as $(e_1 \text{ op1 } e_2) \text{ op2 } e_3$ means that the introduction of **op1** must be done at a level strictly lower than **op2**. In order to modify the grammar so that it generates only this kind of tree, a possible solution is to introduce a new syntactic category producing expressions of the form $e_1 \text{ op } e_2$, and to force an order to op1 and op2.

$\langle \text{Exp} \rangle \rightarrow \text{Num} \mid \langle \text{Exp} \rangle + \langle \text{Exp} \rangle \mid \langle \text{Exp} \rangle * \langle \text{Exp} \rangle$

$\text{Num} \rightarrow 2 \mid 3 \mid 5$

We can eliminate the ambiguities from the grammar by introducing a new syntactic category **Term** producing expressions of the form $\langle \text{Exp} \rangle * \langle \text{Exp} \rangle$

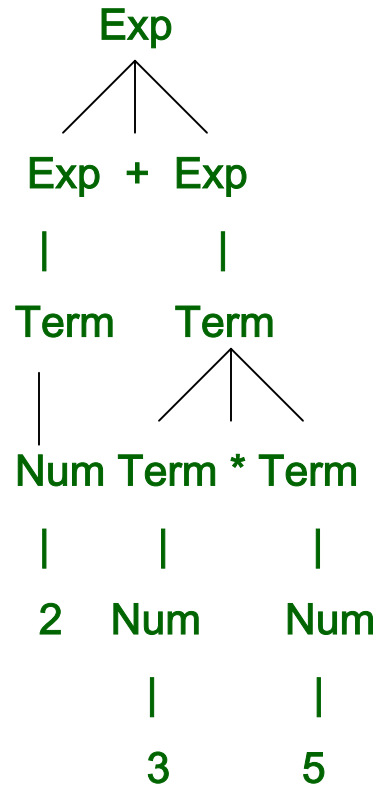
$\langle \text{Exp} \rangle \rightarrow \langle \text{Exp} \rangle + \langle \text{Exp} \rangle \mid \langle \text{Term} \rangle$

$\langle \text{Term} \rangle \rightarrow \langle \text{Term} \rangle * \langle \text{Term} \rangle \mid \text{Num}$

$\text{Num} \rightarrow 2 \mid 3 \mid 5$

This modification corresponds to assigning $*$ a higher priority than $+$.

In the new grammar there is only one tree which can generate it:



$\langle \text{Exp} \rangle \rightarrow \langle \text{Exp} \rangle + \langle \text{Exp} \rangle \mid \langle \text{Term} \rangle$
 $\langle \text{Term} \rangle \rightarrow \langle \text{Term} \rangle * \langle \text{Term} \rangle \mid \text{Num}$
 $\text{Num} \rightarrow 2 \mid 3 \mid 5$

Associativity

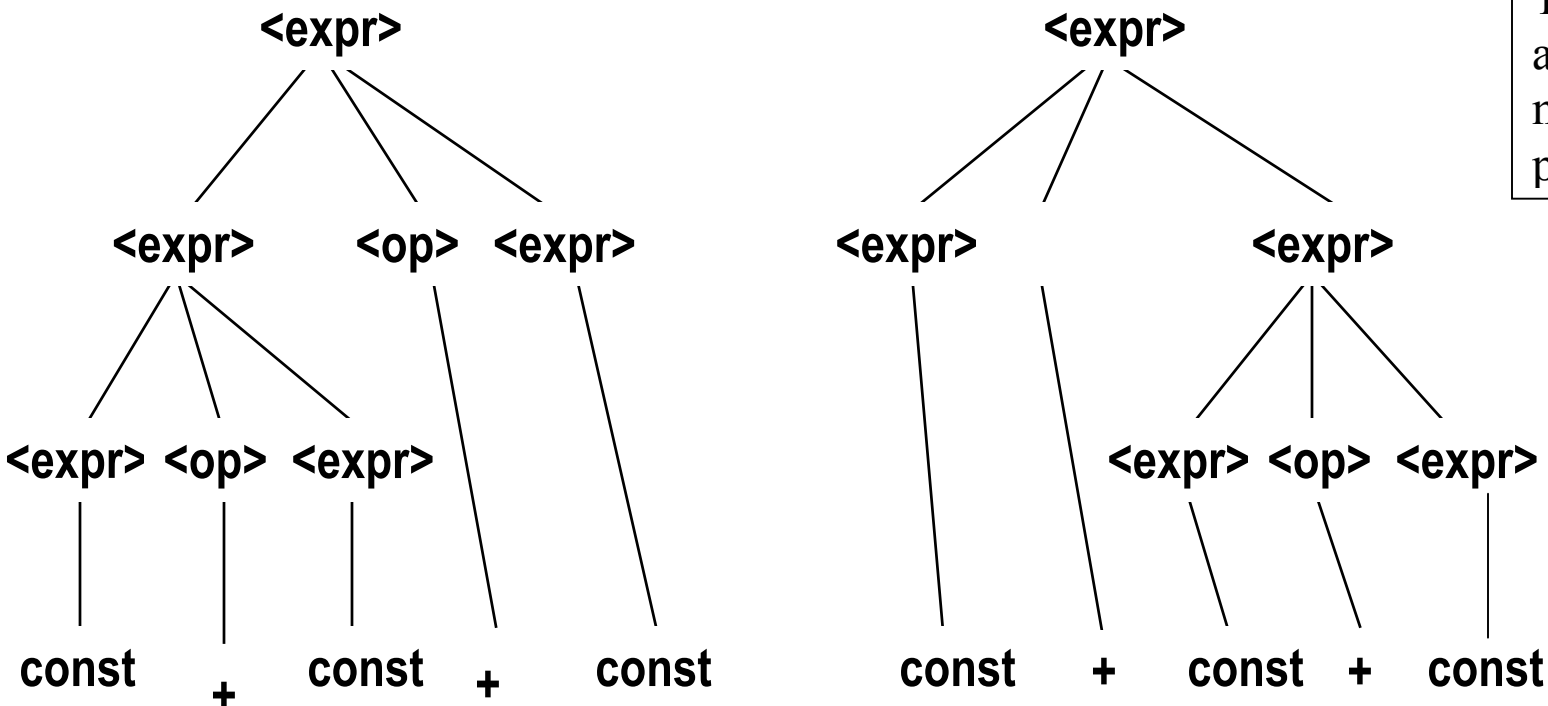
Previous grammar is unambiguous regarding the precedence of $*$ and $+$, but still has ambiguities of another kind.

It allows two different derivation trees for the string $2 + 3 + 5$, one corresponding to the structure $(2 + 3) + 5$ and one corresponding to the structure $2 + (3 + 5)$.

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \quad | \quad \text{const}$

$\langle \text{op} \rangle \rightarrow +$

Two parse trees for expression $\text{const} + \text{const} + \text{const}$



This kind of ambiguity does not cause problems

However, an operator might be not associative. For instance the case for the - and ^ exponentiation operators: $(5 - 3) - 2$ and $5 - (3 - 2)$ have different values, as well as $(5 ^ 3) ^ 2$ and $5 ^ (3 ^ 2)$.

In order to eliminate this kind of ambiguity, we must establish whether the operator is **left-associative** or **right-associative**.

- Left-associative: $e_1 \text{ op } e_2 \text{ op } e_3$ is interpreted as $(e_1 \text{ op } e_2) \text{ op } e_3$ (op associates to the left).
- Right-associative: $e_1 \text{ op } (e_2 \text{ op } e_3)$ (op associates to the right).

We can impose left-associativity (resp. right-associativity) by using a left-recursive (resp. right-recursive) production for op.

$$\langle \text{Exp} \rangle \rightarrow \langle \text{Exp} \rangle + \langle \text{Exp} \rangle \mid \langle \text{Term} \rangle$$
$$\langle \text{Term} \rangle \rightarrow \langle \text{Term} \rangle * \langle \text{Term} \rangle \mid \langle \text{Num} \rangle \quad \text{is changed again to}$$
$$\langle \text{Exp} \rangle \rightarrow \langle \text{Exp} \rangle + \langle \text{Term} \rangle \mid \langle \text{Term} \rangle$$
$$\langle \text{Term} \rangle \rightarrow \langle \text{Term} \rangle * \text{Num} \mid \text{Num}$$

This grammar is now unambiguous.

Another Example

$$\langle \text{Exp} \rangle \rightarrow \text{Num} \mid \langle \text{Exp} \rangle - \langle \text{Exp} \rangle$$

This grammar is ambiguous since it allows both the interpretations $(5 - 3) - 2$ and $5 - (3 - 2)$.

If we want to impose the left-associativity:

$$\langle \text{Exp} \rangle \rightarrow \text{Num} \mid \langle \text{Exp} \rangle - \text{Num}$$

Consider the following grammar:

$$\langle \text{Exp} \rangle \rightarrow \text{Num} \mid \langle \text{Exp} \rangle ^ \langle \text{Exp} \rangle$$

This grammar is ambiguous since it allows both the interpretations $(5 ^ 3) ^ 2$ and $5 ^ (3 ^ 2)$.

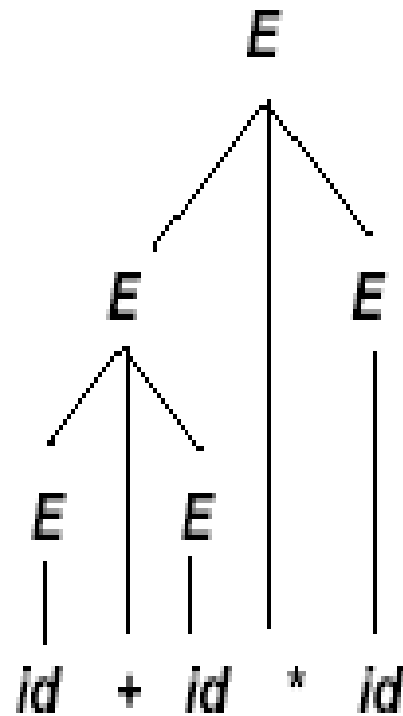
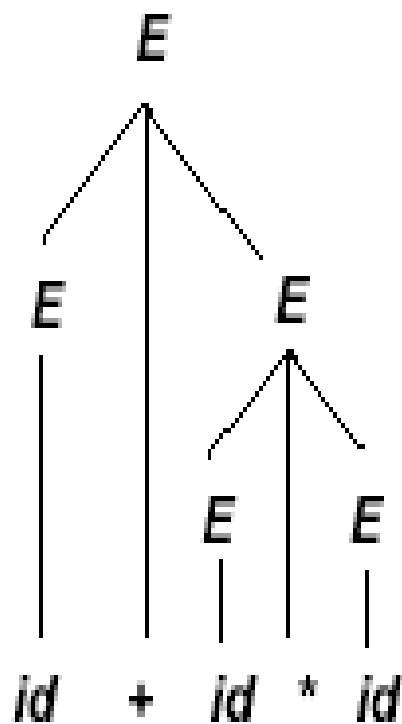
If we want to impose the right-associativity:

$$\langle \text{Exp} \rangle \rightarrow \text{Num} \mid \text{Num} ^ \langle \text{Exp} \rangle$$

Generally, we can eliminate ambiguity by revising the grammar.

Grammar: $\langle E \rangle \rightarrow \langle E \rangle + \langle E \rangle \mid \langle E \rangle * \langle E \rangle \mid (E) \mid id$

Two parse trees for expression $id + id * id$



It is possible to write a grammar for arithmetic expressions that

1. is unambiguous.
2. enforces the precedence of * and / over + and -.
3. enforces left associativity.

Removal of Ambiguity:

Grammar: $\langle E \rangle \rightarrow \langle E \rangle + \langle E \rangle \mid \langle E \rangle * \langle E \rangle \mid (E) \mid \text{id}$

1. Enforce higher precedence for *

$$\begin{aligned}\langle E \rangle &\rightarrow \langle E \rangle + \langle E \rangle \mid \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle * \langle T \rangle \mid \text{id} \mid (E)\end{aligned}$$

2. Eliminate right-recursion for $\langle E \rangle \rightarrow \langle E \rangle + \langle E \rangle$ and $\langle T \rangle \rightarrow \langle T \rangle * \langle T \rangle$.

$$\begin{aligned}\langle E \rangle &\rightarrow \langle E \rangle + \langle T \rangle \mid \langle T \rangle \\ \langle T \rangle &\rightarrow \langle T \rangle * \text{id} \mid \langle T \rangle * (E) \mid \text{id} \mid (E)\end{aligned}$$

Example:

```
position := initial + rate * 60
```

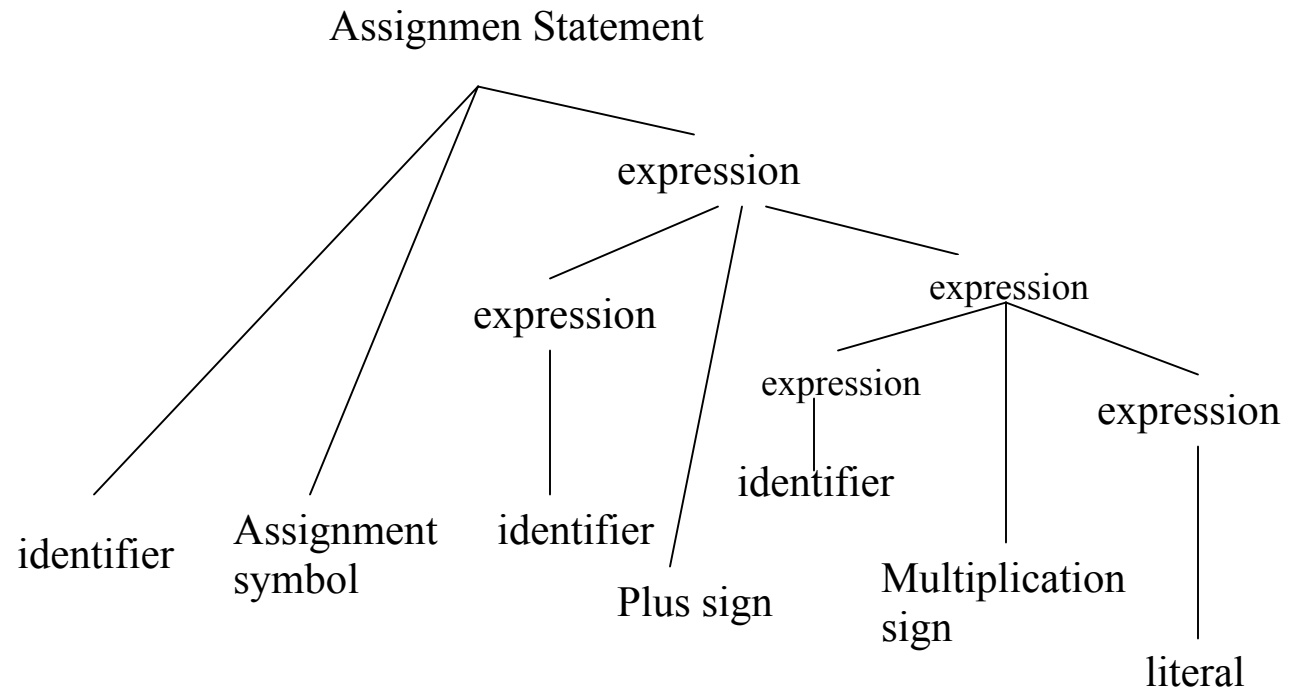
Lexical Analysis

Group these character as follows tokens:

	Token
• position	identifier
• :=	assignment symbol
• initial	identifier
• +	plus sign
• rate	identifier
• *	multiplication sign
• 60	literal

Next, we want to determine that this is a structurally correct statement. This is the main province of syntax analysis or parsing. The result of parsing is a *parse tree*

<Assignment statement> → <identifier><assignment symbol> <expression>
 <expression> → <expression> <plus sign> <expression> | <expression> <minus sign>
 <expression> | <expression> <multiplication sign> <expression>
 <expression> → <identifier> | <literal>



```
position := initial + rate * 60
```

Exercise 1. A grammar of binary numbers

Write a grammar of the language whose elements are all and only those unsigned binary numbers that contain at least three consecutive digits 1. The language includes, for example, 111, 00001111111010 and 1111110, but not 0011000101011 or 1010101010.

Answer 1. A grammar of binary numbers

```
<string> -> <term> | <mix> <term> | <term> <mix> | <mix> <term> <mix>
<mix> -> <bit> <mix> | <bit>
<bit> -> 0 | 1
<term> -> 1 1 1
```

or

```
<string> -> <term> | <bit> <string> | <string> <bit>
<bit> -> 0 | 1
<term> -> 1 1 1
```

Exercise 2 Parse trees

Consider the following grammar with three terminal symbols: a b g.

$$\begin{aligned}\langle S \rangle &\rightarrow a \langle A \rangle \mid \langle B \rangle g \\ \langle A \rangle &\rightarrow \langle C \rangle g \mid b \langle A \rangle \mid a \langle A \rangle \\ \langle B \rangle &\rightarrow \langle C \rangle b \mid \langle C \rangle \langle B \rangle \\ \langle C \rangle &\rightarrow a\end{aligned}$$

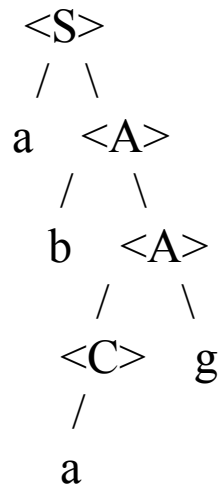
The start symbol is $\langle S \rangle$. Consider the string:
a b a g

Show the leftmost derivation of a b a g,

Next, draw a parse tree for the string a b a g.

Answer 2. Parse trees

$\langle S \rangle \Rightarrow a \langle A \rangle \Rightarrow a b \langle A \rangle \Rightarrow a b \langle C \rangle g \Rightarrow a b a g$



Exercise:

(1) Describe, in English, the language defined by the following grammar:

$$\langle S \rangle \rightarrow \langle A \rangle \langle B \rangle \langle C \rangle$$
$$\langle A \rangle \rightarrow a \langle A \rangle \mid a$$
$$\langle B \rangle \rightarrow b \langle B \rangle \mid b$$
$$\langle C \rangle \rightarrow c \langle C \rangle \mid c$$

(2) Write a grammar for the language consisting of strings that have n copies of the letter a followed by the same number of copies of the letter b , where $n > 0$. For example, the strings ab , $aaaabbbb$ are in the language but a , abb are not.

(3) Convert the following BNF to EBNF

$$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$$
$$\langle \text{id} \rangle \rightarrow A \mid B \mid C$$
$$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle$$
$$\mid \langle \text{expr} \rangle * \langle \text{expr} \rangle$$
$$\mid (\langle \text{expr} \rangle)$$
$$\mid \langle \text{id} \rangle$$

Answer:

(1) One or more a's followed by one or more b's followed by one or more c's.

(2)

$$\langle S \rangle \rightarrow a \langle S \rangle b \mid a b$$

(3)

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