

Ch3. Syntax and Semantics

Syntax – specifies how programs in the language are built up;

Semantics – specifies what the program means;

Same syntax can have different semantics in different world (language).

ex) syntax: DD/DD/DDDD

string (program): 01/02/2017

possible semantics: Jan. 02, 2017 (USA), Feb. 01, 2017 (India?)

Lexical syntax – spelling of words;

notation: regular expression, ex) a^* , ab^+ , $(a|b)^+$, . . .

sample words: a, aa, aaa, ab, abb, abab, . . .

grammar (CFG)

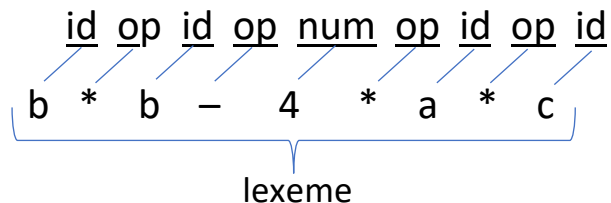
notation: BNF (Backus-Naur Form)

EBNF (extended Backus-Naur Form)

syntax chart (graphical)

ex) Lexical syntax

$b * b - 4 * a * c$ → Lexical analyzer → token stream (terminals)



Expression notation

Expressions such as $a+b*c$ have been in use for centuries and were a starting point for the design of programming languages.

- **prefix** notation – easy to decode

→ *evaluate right-most op first*

$$+ 3 * 5 7 = 38$$

$$* + 20 30 60 = 3000$$

ex) `read(x); max(x, y); min(x, y, z);` //operation first, then data

- **postfix** notation – evaluated by stack

← *evaluate left-most op first*

$$3 5 7 * + = 38$$

$$20 30 + 60 * = 3000$$

+	*
*	60
7	+
5	30
3	20

- **infix** notation – user friendly (easy to read)

decoding (evaluation) needs rules for precedence & associativity;

() breaks the rules;

only infix can use ()

*higher
prec*

^

*, /

+, -

left-assoc:

+, -, *, /

ex) $8/4/2 = 1$

right-assoc:

^

ex) $2^3^4 = 2^{81}$

- **Abstract syntax tree (AST)** – syntactic structure of an expression

independent from the three notations (pre/post/infix);

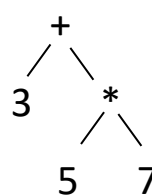
identifies the meaningful components of each construct in the language;

ex) infix: $3 + 5 * 7 (=38)$

prefix: $+ 3 * 5 7$

postfix: $3 5 7 * +$

AST:

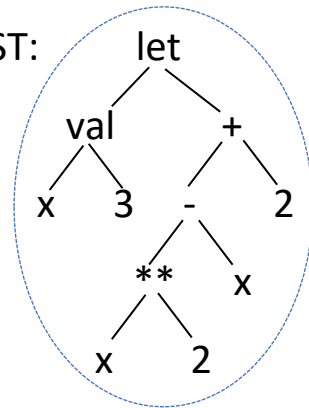


lower precedence up;
higher precedence
down;

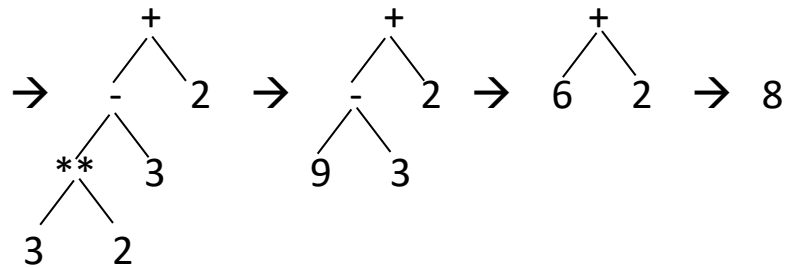
ex) ML (a functional language) let construct

[let val x=3 in
 $x**2 - x + 2$

→ AST:



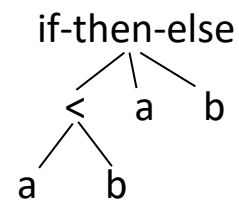
→ reduction steps (evaluation)



ex)	<u>infix</u>	<u>prefix</u>	<u>postfix</u>	<u>AST</u>
	$3 + 5 * 7 = 38$	$+ 3 * 5 7 = 38$	$3 5 7 * + = 38$	
	$(3 + 5) * 7 = 56$ <i>only infix can use ()</i>	$* + 3 5 7 = 56$	$3 5 + 7 * = 56$	
		read(x); max(x, y); min(x, y, z);		

ex) mix-fix expression:

if $a < b$ then a else b



Practice with: $b^2 - 4ac$

Context-Free Grammar (CFG)

is used to specify syntax of a language; infinite set of strings;

notations: BNF, EBNF, syntax chart

Formal definition:

$G = (V, T, P, S)$, where

- V: finite set of non-terminals
 - T: finite set of terminals
 - P: finite set of productions; l.h.s. (V) \rightarrow r.h.s. (V + T)
 - S: start symbol (from V)
- derivation symbol: ::=
- simplified notation: \rightarrow

ex) $\langle \text{sentence} \rangle \rightarrow \langle \text{noun-phrase} \rangle \langle \text{verb-phrase} \rangle$

$\langle \text{noun-phrase} \rangle \rightarrow \langle \text{adjective} \rangle \langle \text{noun-phrase} \rangle$

$\langle \text{noun-phrase} \rangle \rightarrow \langle \text{noun} \rangle$

$\langle \text{noun} \rangle \rightarrow \underline{\text{boy}} \mid \underline{\text{girl}}$

$\langle \text{adjective} \rangle \rightarrow \underline{\text{little}}$

$\langle \text{verb-phrase} \rangle \rightarrow \dots$

terminals

Q: Is “little boy” a part of a legal string in this language? \rightarrow Yes

Is “boy little” a part of a legal string in this language? \rightarrow No

Is “little little boy” a part of a legal string in this language? \rightarrow Yes

- Original motivation of CFG: description of natural language;

ex) syntax of integer numbers (1, 123, 23, ...)

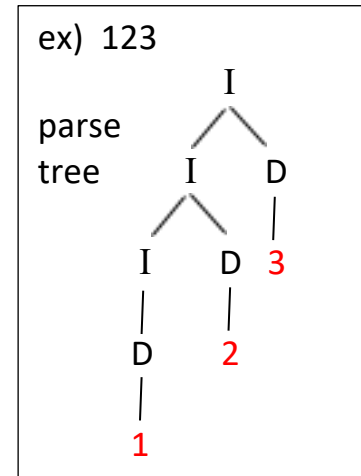
$G = (V, T, P, S)$, where

$V = \{I, D\}$

$T = \{0, 1, 2, 3, \dots, 9\}$

$S = \{I\}$

$P = \left. \begin{array}{l} I \rightarrow D \\ I \rightarrow ID \\ D \rightarrow 0 \\ D \rightarrow 1 \\ \dots \\ D \rightarrow 9 \end{array} \right\} \equiv \left[\begin{array}{l} I \rightarrow D \mid ID \\ D \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9 \end{array} \right]$
 //total 10 productions



ex) syntax of real numbers (123.12, 123, .123, ...)

$G = (V, T, P, S)$, where

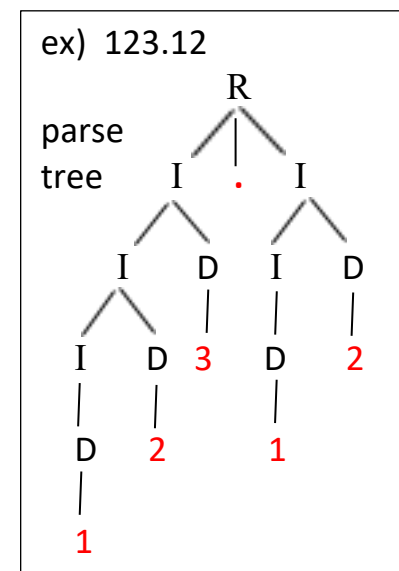
$V = \{R, I, D\}$

$T = \{0, 1, 2, 3, \dots, 9, \cdot\}$ decimal point

$S = \{R\}$

$P = \left. \begin{array}{l} R \rightarrow I.I \\ R \rightarrow I \\ R \rightarrow .I \\ \left. \begin{array}{l} I \rightarrow D \\ I \rightarrow ID \\ D \rightarrow 0 \\ D \rightarrow 1 \\ \dots \\ D \rightarrow 9 \end{array} \right\} \end{array} \right\} \equiv \left[\begin{array}{l} R \rightarrow I.I \mid I \mid .I \\ I \rightarrow D \mid ID \\ D \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9 \end{array} \right]$
 //total 15 productions

same as integer



ex) Language: simplified infix arithmetic expression with
no (), no power operation, single-digit int num.

$G = (V, T, P, S)$, where

$V = \{ E, T, F, \text{Num} \}$

$T = \{ +, -, *, /, 0, 1, 2, 3, \dots, 9 \}$

$S = \{ E \}$

$P = E \rightarrow E + T \mid E - T \mid T$

$T \rightarrow T * F \mid T / F \mid F$

$F \rightarrow \text{Num}$

$\text{Num} \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

E – for expression
T – for term
F – for factor

Parse tree

depicts concrete syntax;

shows how a string is built from a given syntax (grammar);

A string (e.g., C++ program) is in a language (e.g., C++) iff it is generated by a parse tree, i.e., \exists derivation.

The construction of a parse tree is called parsing.

Parse tree structure:

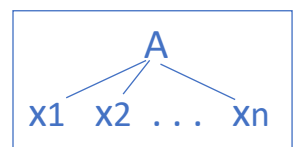
root – start symbol from grammar;

each leaf – labeled by a terminal symbol;

internal nodes – non-terminals;

if a non-terminal A has children labeled x_1, x_2, \dots, x_n ,

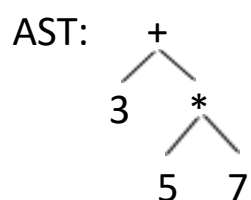
then $A \rightarrow x_1 x_2 \dots x_n$ is a production.



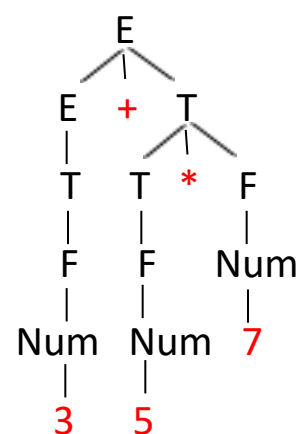
- A grammar for a language is designed to reflect AST;
the productions are chosen, s.t. parse trees are as close as possible to AST.

ex) string: 3 + 5 * 7

parse tree:



CFG: $E \rightarrow E + T \mid E - T \mid T$
 $T \rightarrow T * F \mid T / F \mid F$
 $F \rightarrow \text{Num}$
 $\text{Num} \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$



AST – built from string only;

parse tree – built from G and string;

Lower subtree is evaluated first;

Derivation – snap shots of leaf nodes during building a parse tree

- right-most derivation: derive the right-most non-T first (in r.h.s. of a prod.);
- left-most derivation: derive the left-most non-T first (in r.h.s. of a prod.);

ex) G: above example, string: 3+5*7

left-most: $E \rightarrow E + T$
 $\rightarrow T + T$
 $\rightarrow F + T$
 $\rightarrow \text{Num} + T$
 $\rightarrow 3 + T$
 \dots
 $\rightarrow 3 + 5 * 7$

right-most: $E \rightarrow E + T$
 $\rightarrow E + T * F$
 $\rightarrow E + T * \text{Num}$
 $\rightarrow E + T * 7$
 $\rightarrow E + F * 7$
 \dots
 $\rightarrow 3 + 5 * 7$

Syntactic ambiguity

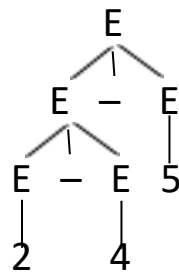
A grammar for a language is ambiguous if some string in the language has more than 1 parse trees.

ex) $G: E \rightarrow E - E$

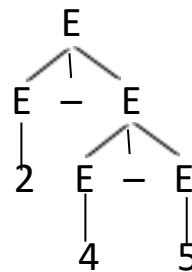
$| 0 | 1 | 2 | \dots | 9$

string: $2 - 4 - 5$

parse tree:



$(2 - 4) - 5 = -7$ //correct



$2 - (4 - 5) = 1$ //incorrect

\exists two parse trees for a given string \rightarrow grammar is ambiguous;

incorrect grammar

- for writing unambiguous grammar,

- keep precedence, i.e., lower precedence up;
- keep associativity, i.e.,

left-assoc. op – left-recursive grammar,
right-assoc. op – right recursive grammar;

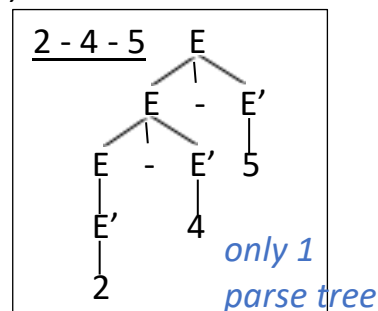
$E \rightarrow E + T \mid T$	----	$+$	
$T \rightarrow T * F \mid F$	----	$*$	
$F \rightarrow \text{Num} \wedge F$	----	\wedge	
....			

ex) since $(-)$ is left-associative operator, G should be left-recursive.

$E \rightarrow E - E'$

$| E'$

$E' \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$



Associativity – reflect associativity in the grammar (for unambiguous G);

left-associative op (+, -, *, /)

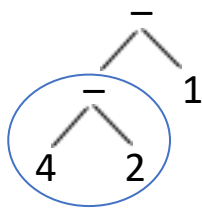
→ G should be left-recursive;

ex) +, -

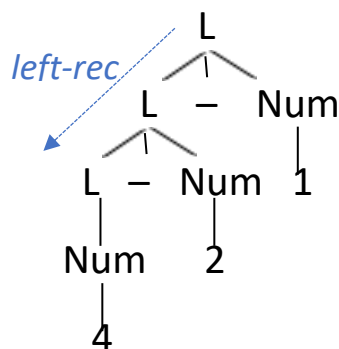
$$\left[\begin{array}{l} L \rightarrow L + \text{Num} \\ \quad | L - \text{Num} \\ \quad | \text{Num} \\ \text{Num} \rightarrow 1|2|\dots|9 \end{array} \right]$$

ex. string: 4 - 2 - 1

AST:



parse tree:



right-associative op (^{assignment}^, :=)

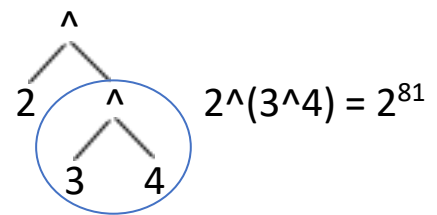
→ G should be right-recursive;

ex) ^ (power)

$$\left[\begin{array}{l} R \rightarrow \text{Num} ^ R \\ \quad | \text{Num} \\ \text{Num} \rightarrow 1|2|\dots|9 \end{array} \right]$$

ex. string: 2 ^ 3 ^ 4

AST:



parse tree:

