TinyJ Assignment 1 email*: Do the installation before next Monday.
* With assigned reading and exercises on syntax and method overriding in c++.

## Syntax of Programming Languages.

The syntax of programming language is usually specified in terms of

① Units called <u>tokens</u>, such as: <u>while</u>, <u>for</u>, ~~IDENTIFY~~ IDENTIFIER (x, while, for)

UNSIGNED-int-Literal. (72. 47193) ; ( >= - (java/c++ tokens).

Some tokens (e.g. IDENTIFIER and UNSIGNED-INT-LITERAL in Java/c++)

have more than one instance.

Different instances of the same token are semantically different (ie. have different

meanings) but syntax specification written in BNF or some similar notation such as

EBNF will <u>not</u> distinguish between them.

Note: It is also possible for the <u>same instance</u> of token to have more than one spelling,

        Apple     73.0
        ApplE     73.00

A token (like <u>while</u> or <u>for</u> in java/c++) that looks like an identifier but <u>cannot</u>

be used as identifiers are called <u>reserved words</u>. [Some languages, such as

Fortran and Lisp, has no reserved words, in such languages words that have special

[(defun nil (if (+ if 1)) is a legal Lisp function definition]  meaning (like IF in Lisp)
                                                                  are called

<u>keywords</u>    Some author refer to ~~reserved~~ reversed words as keywords,

② Certain choozen language constructs (e.g. <while stmt>, <stmt> or

<expression>).    In a syntax speciafication of an entire language ( rather just

some part of a language) ~~one~~ of these constructs represents acceptable input for

a complier or interpreter — e.g. it might be called <program>.

The purpose of a syntax specification is to answer the following 2 questions:

<u>Q1</u>: For each kind of token exactly which sequences of characters constitute a valid
      instance of that token.

**Q2** For each of the language constructs of ②, exactly which sequences of tokens constitute a syntactically valid instance of that language construct?

[Reminder] A syntactically valid $\left(\begin{smallmatrix} x \\ y \end{smallmatrix} = 4*(y-w)\right)$ instance of a construct need _not_ be a valid instance — it may contain errors such as ~~understand~~ underclared identifier type ~~mismath~~ mismatch or devided by zero.

Q1 and Q2 are usually answered using context-free grammar notation/BNF or some variant of this (such as EBNF)

Context-Free grammer Notation is a relatively concise (and yet completely price) notation for ~~specifiy~~ specifying a (possible infinite) set of finite sequences of symbols. The symbols are called the Terminal of the grammer.

In ~~instance~~ answering Q1 (for tokens that have many instances, like IDENTIFIER or Unsigned-float-Litera) we can use a grammer whose ~~term~~ terminals are character; [Regular expression notation is another way to do this]

In answing Q2, we can ~~use~~ use a grammer whose terminal are tokens It is also possible to use a single grammar, whose ~~ston~~ terminals are character, to answer ~~both~~ both Q1 and Q2.

Context-free grammer notation was intended by Chamsky, who invented a hierarchy of 4 types of grammer for specifying the syntax of English and other similar natural languages. Context-free grammars are Type 2 in Chamsky's Hierarchy.

the idea of a context-free grammer was independently reinvented by Backus a few years later. Backus proposed the use of his grammar notation for specifying the syntax of Algol 60.

the Algo 60 Report (which specified the language Algol 60 in a way that was widely admired by computer scientists) adopted Backus's proposal.

Naur was the editor of the Algol 60 report. Backus's notation is now call BNF

(Backus Naur Form)

$\langle E \rangle ::= \langle E \rangle + \langle T \rangle$  (BNF)

$E \rightarrow E + T$.  (not BNF)

Like many authors today⒜ we will use BNF to mean some ~~p cont~~ commonly used

notation for writing a context-free grammar.

Ⓧ But unlike Sethi:

cf Fig 2.6 on ——— P42.

and Fig 2.10 on P46.

Sethi would **not** consider Fig 2.6 to be BNF.
But we will consider it to be BNF.

## Ex of a Grammar.

UNSIGNED-FLOAT-LITERAL tokens

can be specified in BNF as follows : ( assuming they have the form: ~~one or more~~

~~digits. on~~  <u>one or more digits</u> . <u>one or more digits</u> .        (Fig 2.3 on P36.)
                     ip (integer part)       f (fraction) ← fraction.

① ufpl ::= ip.f

② ip ::= d        ip ::= d / ip d.   ("d" = "digit").

③ ip ::= ip d ⑤
④ f ::= d / d f
   d ::= ⓪0 / ⓪1 / 2 / 3 / 4 / 5 / 6 / 7 / 8 / 9 ⑮. ── undefined constant symbols

The ~~∅~~ 11 symbols  0, 1 —, 9. are the <u>terminals</u> of this grammar.

The 4 symbols ufpl, ip, f, d are the <u>nonterminals</u> of this grammar.
                                                    ↓
                        variables, each of which denotes of finite sequences of terminals.

There are 15 rules called productions ~~each~~ each of the form

$$N ::= \alpha \underline{\quad}$$

a single <u>nonterminal</u>          a sequence of zero or move terminals and/or nonterminals.

→ indicated, the
starting nonterminal
~~on the left~~ is the
nonterminal
on the left
side of
the
first
produc-
-tion.

The <u>non</u>terminal ufpl in the starting nonterminal

This nonterminal denotes the set of sequences that the grammar is intended to define
── this set is called the <u>language</u> of / generated by the grammar. Unless otherwise