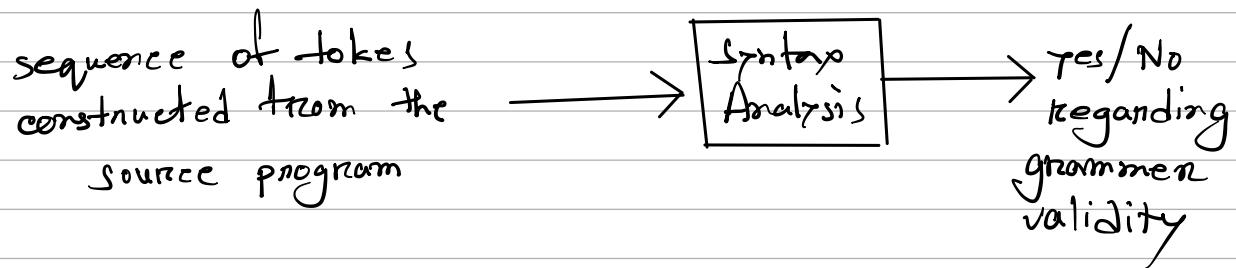


## Lecture 3: Syntax Analysis

\* During implementation syntax & semantic analysis are combined together.

Syntax analysis  $\Rightarrow$  Grammar Checking



Noam Chomsky

$\hookrightarrow$  invented the hierarchy of language.

\* Language is a set of strings.

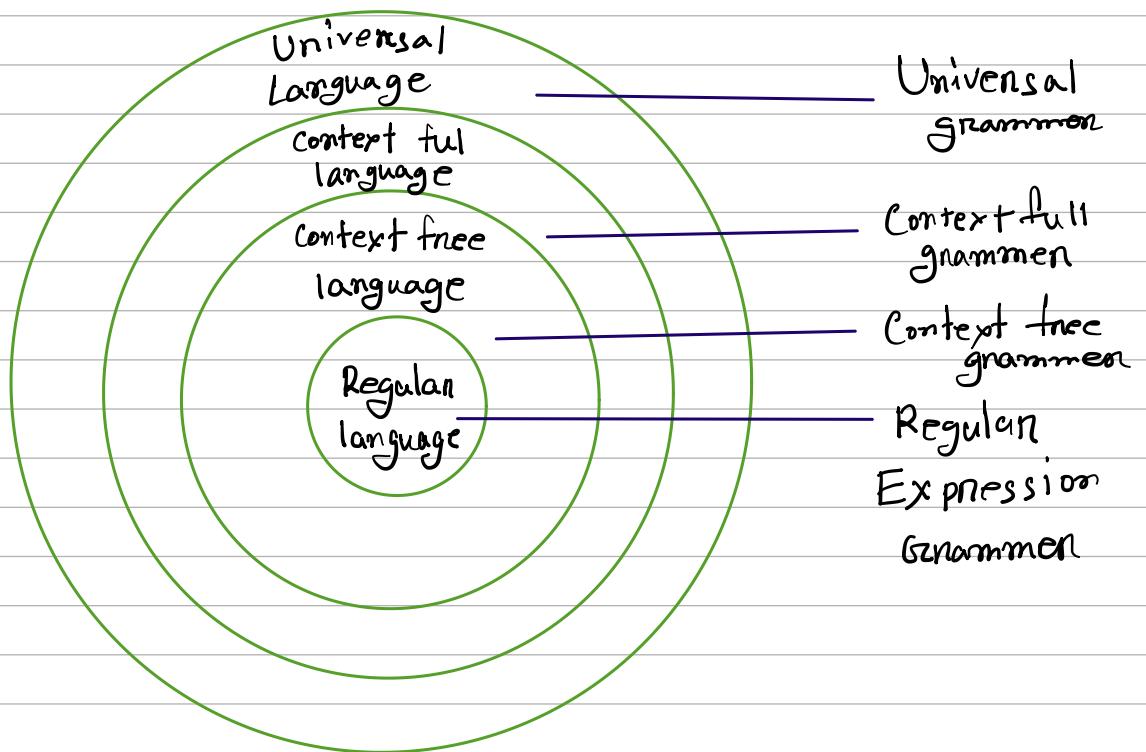
—————  
Grammar of the language permits

Ate rice he  $\rightarrow$  not permitted  $\times$

He ate rice  $\rightarrow$  permitted  $\checkmark$

$\therefore$  Grammar is the tool that defines what is acceptable.

Defined by



→ Human / Natural language is contextfull.

Selim stucked Kanim, he was crying.

who?  
requires context → tough to get context

\* All programming languages are contextfree. (CFGs)

- the purpose of compiler is to produce executable that exactly matches the programmer's idea. There cannot be any vagueness. To translate exactly we don't need context.
- ) It is hard to construct contextfull compiler.

reason  
to  
use  
context free  
grammars

## Components of CFG:

- ①  $\Sigma \leftarrow$  alphabet of terminal symbols.
- ② A set of non-terminals
- ③ Production rules that convert a non-terminal into a string of terminals and/or non-terminals.
- ④ A starting non-terminal.

Production rules for arithmetic expressions:

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E * E$$

$$E \rightarrow E / E$$

!

$$E \rightarrow \text{number}$$

$$E \rightarrow \text{identification} \rightarrow$$

$$E \rightarrow (E) \qquad \begin{matrix} \text{variable} \\ \text{names} \end{matrix}$$

$\Rightarrow$  on the left side we have a single non-terminal than it is context-free. otherwise contextfull.

$$\Sigma = \{ +, -, *, /, \text{number}, \text{identification}, (, ) \}$$

$\hookrightarrow$  set of token types  
from lexical analysis

id + number \* id

$$E \rightarrow E + E$$

$$\rightarrow E + E * E$$

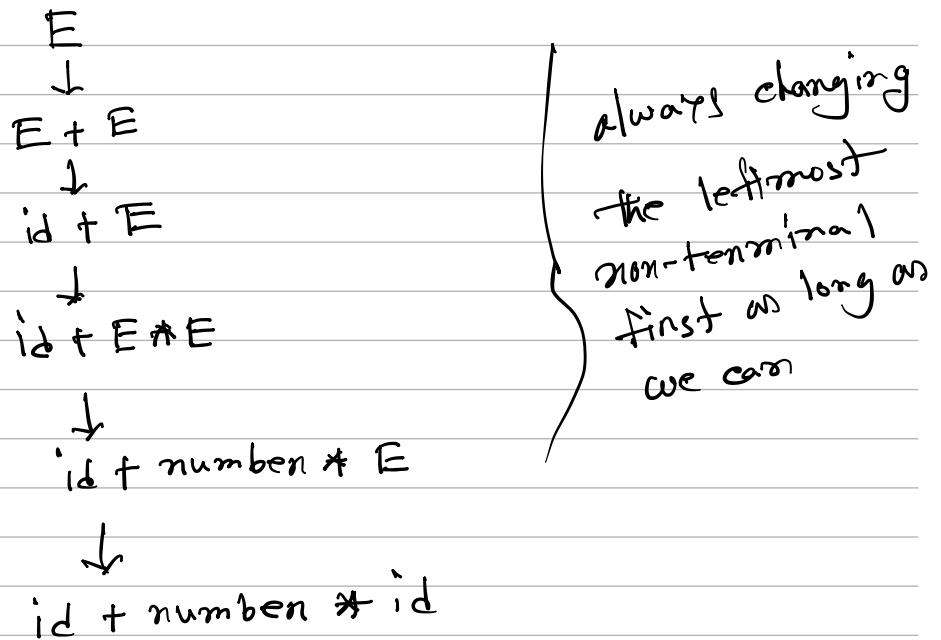
$$\rightarrow id + E * E$$

$$\rightarrow id + \text{number} * E$$

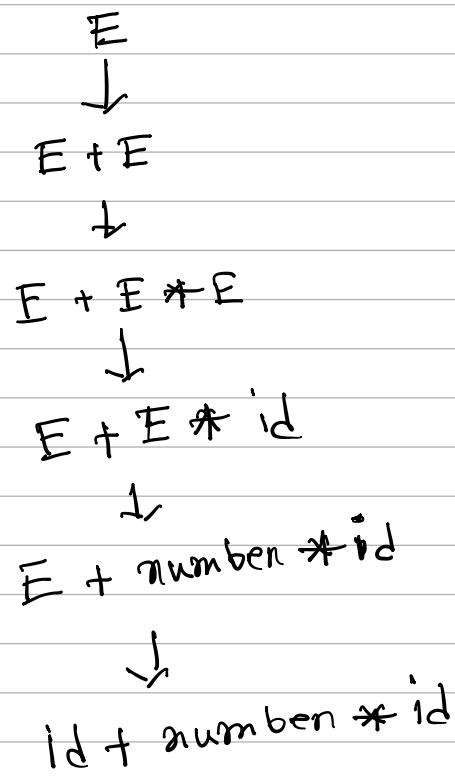
$$\rightarrow id + \text{number} * id$$

derivation

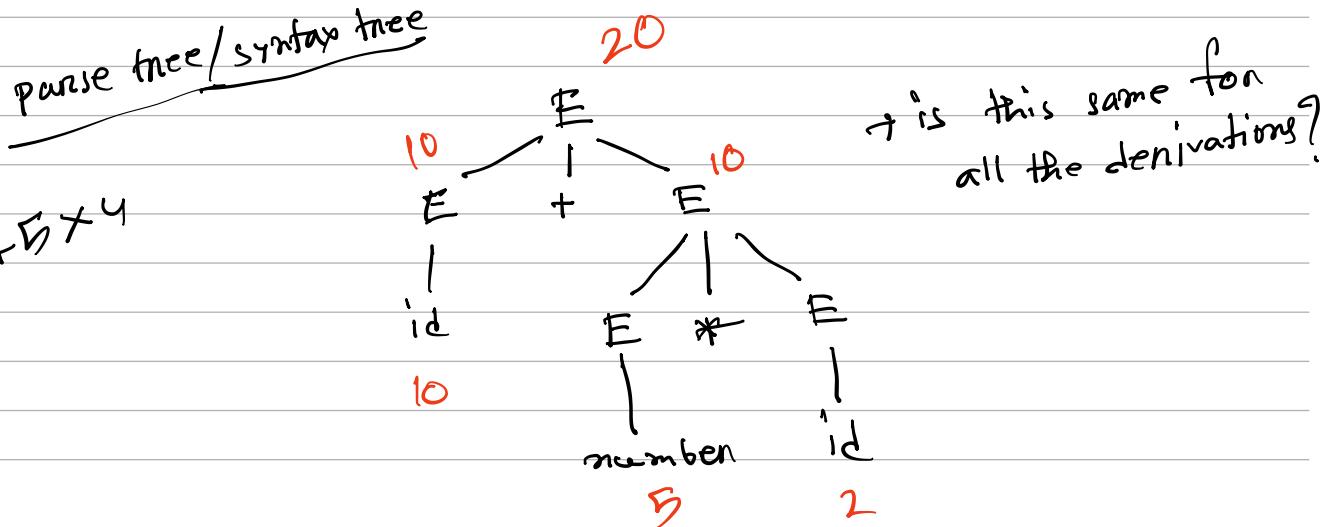
leftmost derivation:



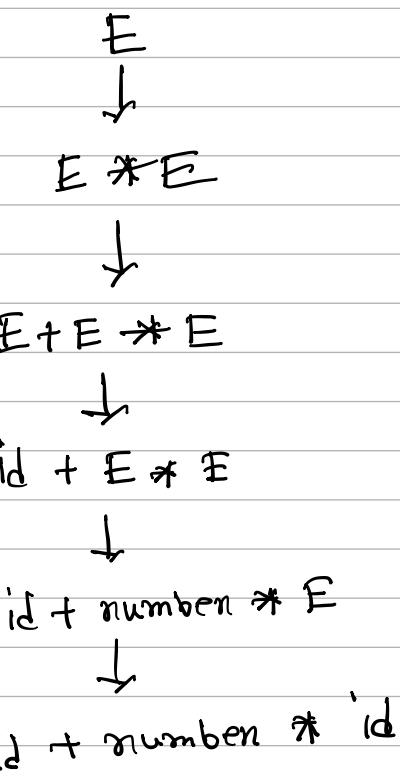
rightmost!



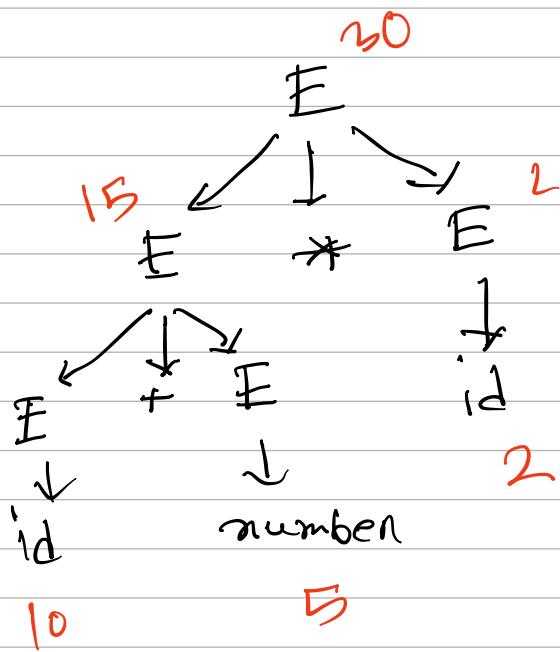
parses tree / syntax tree



*another leftmost derivation*



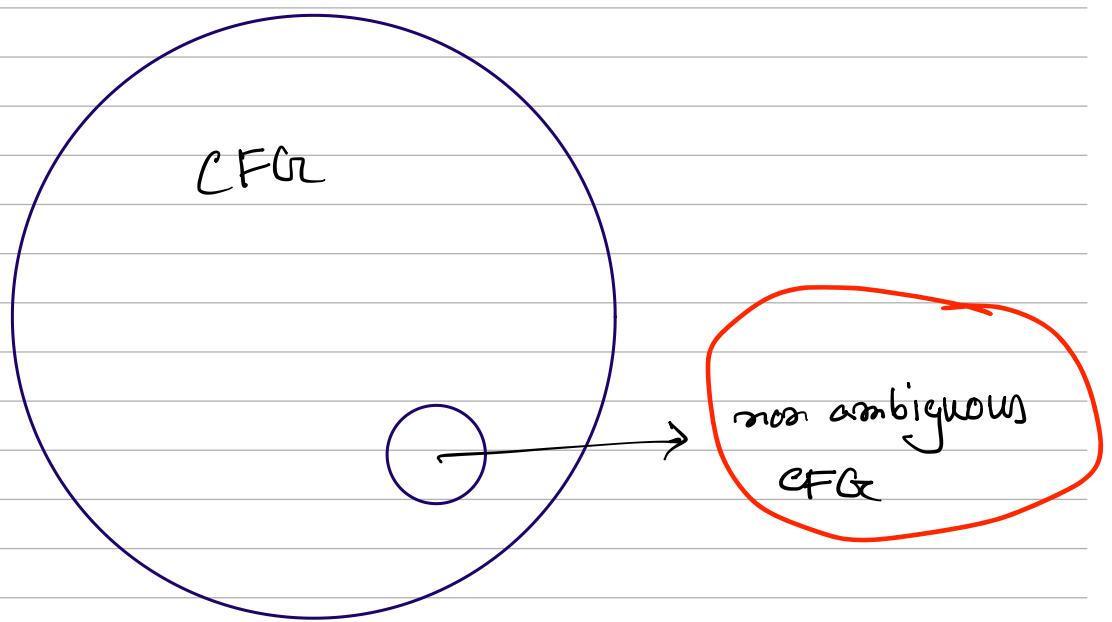
Parse tree for  
this:



For the same values, we get two different outputs for two different parse trees. This type of grammar is called ambiguous.

We cannot use ambiguous grammars  
for our programming language.

∴ All CFGs cannot be used for  
programming languages.



{ why we have to use semicolon in C ?  
" " " " " indents in Python ?

all these are to maintain the non ambiguity  
of grammars .

context free → Parser generator.  
grammars