

BNF



In 1958 John Backus proposed formalism for describing languages.

Peter Naur modified it. It became BNF used first for describing ALGOL.

BNF is metalanguage.

BNF uses abstract syntax structures.

Uses Rules or Productions.

Ex:

$\langle \text{assign} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$

Abstractions are called non-terminals.

Prof.R.Gururaj

CSF301

FPL

BITS Pilani, Hyderabad Campus

Token \rightarrow category of lexeme

LHS \rightarrow RHS \rightarrow either
non terminal
terminal

$\left\{ \begin{array}{l} \langle E \rangle \rightarrow \langle E \rangle + \langle E \rangle \\ \langle E \rangle \rightarrow \langle E \rangle - \langle E \rangle \\ \langle E \rangle \rightarrow \langle \text{Term} \rangle \\ \langle \text{Term} \rangle \rightarrow \text{id} \end{array} \right.$ brackets = things are non terminal
+ - id are terminal

lets eliminate $\langle \rangle$ for brevity

Sequence of replacements \rightarrow Derivation

$\left\{ \begin{array}{l} E \rightarrow E + E \\ \rightarrow E - E + E \\ \rightarrow \text{term} - E + E \end{array} \right. \rightarrow \left\{ \begin{array}{l} \text{id} - E + E \\ \text{id} - \text{term} + E \\ \text{id} - \text{id} + \text{id} \end{array} \right.$

Sentential form \rightarrow Language

Sentence \in Language

Start from the starting non terminal

↳ All intermediate steps are sentential form
↳ Every sentence is a sentential form
(RIP spelling)

By applying grammar rules, we get a derivation
These grammar rules are generators.

$\langle \text{program} \rangle \rightarrow \text{begin} \langle \text{stmt_list} \rangle \text{end}$

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle \langle \text{expression} \rangle$

A grammar for a small lang

- ① $\langle \text{program} \rangle \rightarrow \text{begin} \langle \text{stmt_list} \rangle \text{end}$
- ② $\langle \text{stmt_list} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmt_list} \rangle$
- ③ $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expression} \rangle$
- ④ $\langle \text{var} \rangle \rightarrow A \mid B \mid C$
- ⑤ $\langle \text{expression} \rangle \rightarrow \langle \text{var} \rangle + \langle \text{var} \rangle \mid \langle \text{var} \rangle - \langle \text{var} \rangle \mid \langle \text{var} \rangle$

```

<program> => begin <stmt_list> end
           => begin <stmt> ; <stmt_list> end
           => begin <var> = <expression> ; <stmt_list> end
           => begin A = <expression> ; <stmt_list> end
           => begin A = <var> + <var> ; <stmt_list> end
           => begin A = B + <var> ; <stmt_list> end
           => begin A = B + C ; <stmt_list> end
           => begin A = B + C ; <stmt> end
           => begin A = B + C ; <var> = <expression> end
           => begin A = B + C ; B = <expression> end
           => begin A = B + C ; B = <var> end
           => begin A = B + C ; B = C end

```

→ Last ~~a~~ sentential form is a sentence, and has only terminals

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A | B | C$

$\langle \text{expr} \rangle \rightarrow \langle \text{id} \rangle + \langle \text{expr} \rangle |$
 $\langle \text{id} \rangle * \langle \text{expr} \rangle |$
 $\langle \langle \text{expr} \rangle \rangle |$
 $\langle \text{id} \rangle$

→ whitespace is not allowed with a token

→ 001

$\langle \text{int} \rangle = \langle \text{digit-list} \rangle | \langle \text{digit} \rangle$

$\langle \text{digit-list} \rangle = \langle \text{digit} \rangle | \langle \text{digit} \rangle \langle \text{digit-list} \rangle$

$\langle \text{digit} \rangle = 0 | 1 | 2 | 3 | 4 | 5 \dots 9$

(assume 000 is valid)

Include sign

$\langle \text{int} \rangle \rightarrow \langle \text{sign} \rangle \langle \text{digit_list} \rangle$
 $\langle \text{digit_list} \rangle \rightarrow \langle \text{digit} \rangle \mid \langle \text{digit} \rangle \langle \text{digit_list} \rangle$
 $\langle \text{digit} \rangle \Rightarrow 0 \mid 1 \mid 2 \dots 9$
 $\langle \text{sign} \rangle \rightarrow + \mid -$

Exclude -01 -0036 1060 etc
[$+4036$ is valid]

(1w)

$\langle \text{int} \rangle \rightarrow$
 $\langle \text{digit_list} \rangle$

$\langle \text{sign} \rangle \langle \text{digit_list} \rangle$

Sentence list

$E \rightarrow E + E$ $E \rightarrow \text{term}$
 $E \rightarrow E - E$ $\text{term} \rightarrow \text{id}$

NTs: $\{ E, \text{term} \}$
Terminals: $\{ +, -, \text{id} \}$
Start NT: E

- Your derivation should always start from starting non terminal
- Your starting non terminal contains only 1 lexeme

Then choose a non terminal and replace it with a suitable non terminal / terminal
repeat till you reach all terminals

if you start from the left-most non terminal
and where & repeat

You get a left-most Derivation (LMD)
they RMD for right-most Derivation (RMD)

LMD

$$E \rightarrow \underline{E} + E$$

$$E \rightarrow \underline{E} + E + E$$

$$E \rightarrow \underline{\text{term}} + E + E$$

$$E \rightarrow \text{id} + \underline{E} + E$$

$$E \rightarrow \text{id} + \text{term} + E$$

⋮

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

RMD

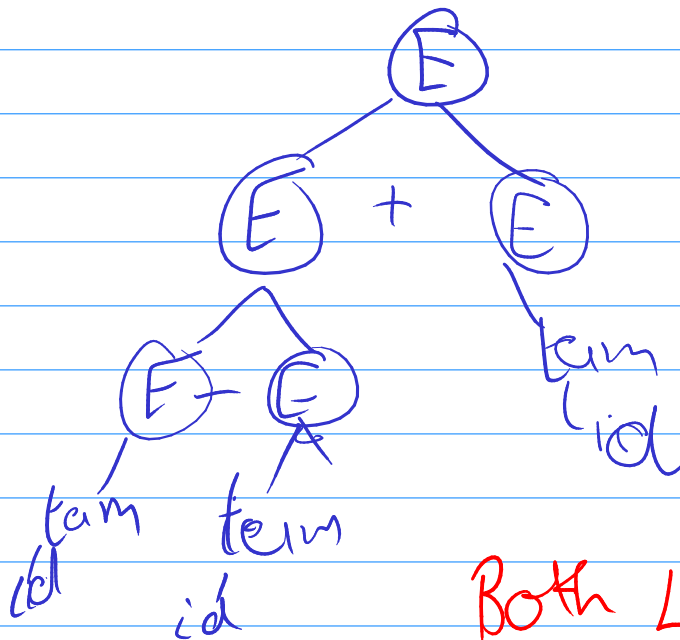
$$E \rightarrow E + \underline{E}$$

$$E \rightarrow E + E - \underline{E}$$

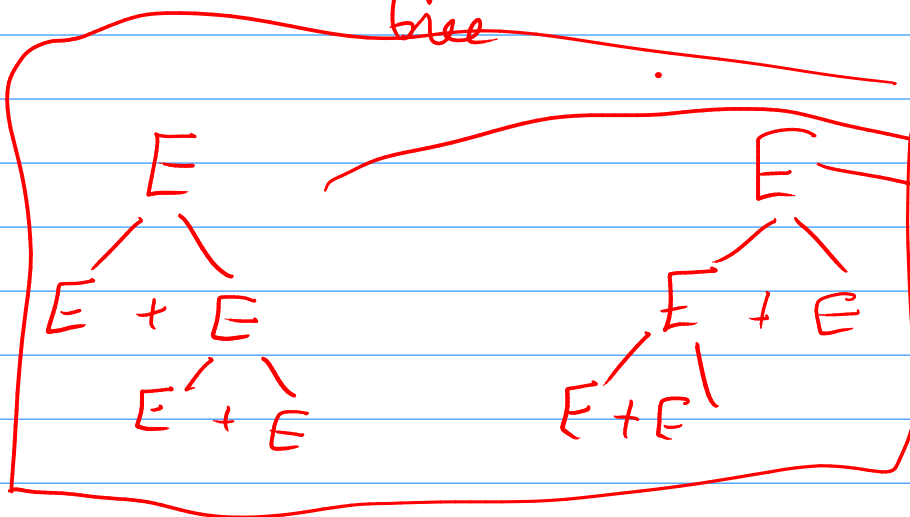
$$E \rightarrow E + E - \underline{\text{term}}$$

⋮

$$E \rightarrow \text{id} + \text{id} - \text{id}$$



Both LMD & RMD are represented the same way in a parse tree



However sometimes different parse tree for some expression \Rightarrow ambiguous grammar