# Language Engineering | Lecture 8

## Parsing

Typically, languages are defined in terms of a grammar that describes the valid construction of sentences in that language.

The designers of Algol 58 invented BNF to describe the language (Backus-Naur Form)

eg.

$$\langle expr \rangle ::= \langle term \rangle \ | \ \langle expr \rangle \ " + " \ \langle term \rangle$$

things in "<" ">" brackets are non-terminals

things in quotes are terminals

this says that the LHS can be produced by the rules on the RHS

this gives us the choice between productions

$$\langle term \rangle ::= \langle factor \rangle \ | \ \langle term \rangle \ "*" \ \langle factor \rangle$$
$$\langle factor \rangle ::= \langle constant \rangle \ | \ \langle variable \rangle$$
$$\ | \ "(" \ \langle expression \rangle \ ")"$$

$\langle variable \rangle ::= "x" \mid "y" \mid "z".$

$\langle constant \rangle ::= \langle digit \rangle \mid \langle digit \rangle \langle constant \rangle$

$\langle digit \rangle ::= "0" \mid "1" \mid "2" \mid ... \mid "9".$

BNF consists of

- ~~nat~~ erminals          $\langle expr \rangle$

- ~~non~~terminals          "3"

- alternations          p | q

This was extended to Extended BNF or EBNF, which has a few more constructs:

- optionals :

     $\langle term \rangle ::= [ "-" ] \langle factor \rangle$

                ↑

         optional "-"

In BNF we could write this as:

     $\langle term \rangle ::= \langle factor \rangle$

         $\mid "-" \langle factor \rangle$

- repetition:

$$\langle args \rangle ::= \langle arg \rangle \ \{ \ "," \ \langle arg \rangle \ \}$$

the {ad} brackets tell us that
this can be repeated 0 or more times.

- grouping:

$$\langle expr \rangle ::= \langle term \rangle \ ( \ "+" \ | \ "-" ) \ \langle expr \rangle$$

the ( and ) allow us to have
scope for some BNF rules.

"~~#~~    " + 3 * 9"

"1 + 9 × 3"

$\langle digit \rangle$ " + 9 × 3"

Parsing is traditionally broken up
into:
    * lexical analysis
    * parsing.

* lexical analysis breaks up an input
  stream into tokens.
  The tokens correspond to ~~non~~ terminals

  this stage usually deals with
  whitespace.
  Not all whitespace is necessarily
  stripped eg. newlines /tabs
  might turn into tokens.

  tools that did this phase
  are called lexers popular tools
  are lex and flex.

* A parser traditionally consumes
  tokens from the lexical analysis,
  and maps rules onto datatypes.
  Tools that do this are
  deeply tied into the language
  the compiler is written in.
  Traditionally, the tools bison,
  yacc, or (antlr) were used
  for this (java)

We will be using parser combinators instead. This builds up parsers from smaller one.

We will be looking at parsec.