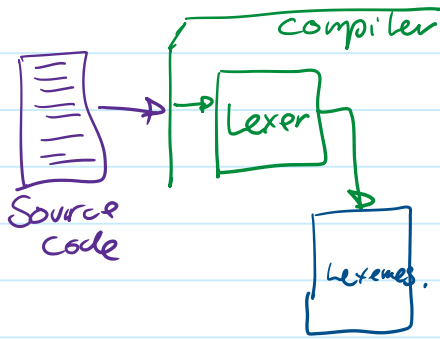


# 4 Syntax and Regular Expressions

Wednesday, September 4, 2024 12:00 PM



from Linguistics:

- Orthography: - How to write words?

apple      manzana      pomme      𐀀𐀁𐀂

this gives a Dictionary.

In a programming Language : keywords.

for    while    if    do    else.

- Syntax: - How to write sentences?

"has jump apple dog green"

syntax has rules.

in programming languages. also

- Semantics: - what do sentences mean?

"colorless green ideas sleep furiously"

Adj

Adj

noun

verb

adverb.

N. Chomsky

"The chickens are ready to eat"

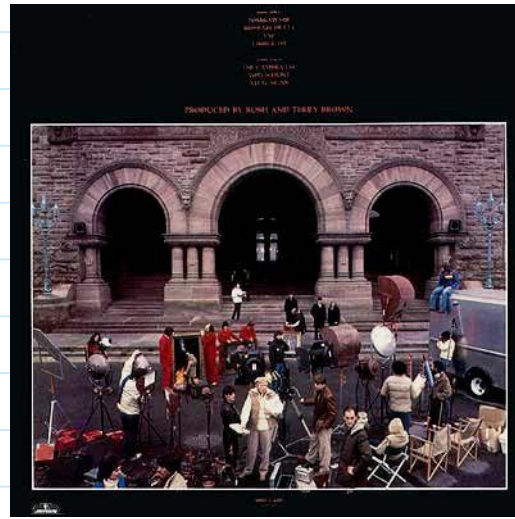
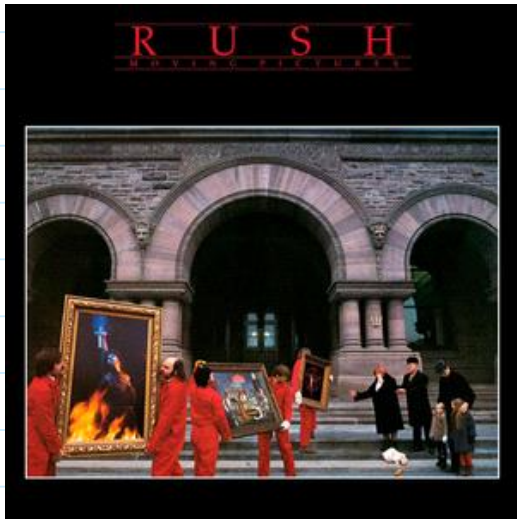


or



"Morina Pictures"

# "Moving Pictures"



## - Pragmatics.-

How meaning changes over time.

"the program has to handle an unreliable cloud"

## - Syntax: for a programming language.

### - How to write words?

Ascii

- keywords

- Identifiers, - names for things  
variable, functions, classes, <sup>templates</sup>  
namespaces .....

- Literals, -

constant values: strings  
integer floats

To answer this question:

"Mathematical Languages" S. Kleene  
E. post.

DEF: • Alphabet - a finite set of symbols.

eg  $\Sigma = \{a, b\}$

- Language: a set of sequences (strings) (words) made of symbols in the alphabet  
e.g.  $L = \{a, aa, aba, baa, bbb\}$

## 1.- Specification Problem:

How to precisely describe a language?

$$L = \{a, ab, aba, abb, \dots\}$$

??

## 2.- Recognition Problem

given a word  $w$  and a Language  $L$

answer:  $w \in L?$

- Answer to #1 - specification: Regular expressions

built from:  $\Sigma$  plus  $|$   $*$   $($   $)$

DEF Regular Expressions:

- a word  $w$  is a regular expression
- if  $r_1, r_2$  are regular expressions
  - $r_1 r_2$
  - $r_1 | r_2$
  - $r_1^*$
 } are also regular expressions.

A regular expression specifies a language.

- $w$              $\{w\}$
  - $r_1 r_2$         $\{w_1 w_2 \mid w_1 \in r_1 \wedge w_2 \in r_2\}$
  - $r_1 | r_2$       $\{w \mid w \in r_1 \vee w \in r_2\}$
  - $r^*$             $\{w \mid w \text{ is zero or more repetitions of a word in } r\}$
- ↑  
Kleene Star.

E.g:  $(a | b) b = \{ab, bb\}$

$a : \{a\}$

E.g.  $(a|b)^*$   
 $a : \{a\}$   
 $b : \{b\}$   
 $(a|b) : \{a, b\}$

Eg:  $(ab)^* = \{\epsilon, ab, abab, ababab, \dots\}$   
 $ab : \{ab\}$   
↑ empty string.

★ an infinite language is specified by finite means.

Eg:  $(a|b)^* = \{\epsilon, a, b, aa, ab, aabba, \dots\}$   
 $a : \{a\}$   
 $b : \{b\}$   
 $a|b : \{a, b\}$

Eg:  $(aa|b^*)a = \{a, ba, aaa, bba, bbba, bbbba, \dots\}$   
 $aa : \{aa\}$   
 $b^* : \{\epsilon, b, bb, bbb, bbbb, \dots\}$

• Shorthands

$r? \equiv (r|\epsilon)$  "to r, or not to r"

$r+ \equiv r(r^*)$  "one or more r's"

⇒ Regular expressions & Programming Languages.

$\Sigma = \text{ASCII keyboard characters.}$

Note: you cannot use regular expressions to describe a programming language.

$D_{C++}$  = the set of all valid C++ programs.

Why?

consider simple example:

$\{a^n b^n \mid n \in \mathbb{Z}^+\} = \{ab, aabb, aaabbb, aaaabbbb, \dots\}$

consider matching brackets.

$\{( ) [ ] \}$

$\{([ ])\}$

Famous stack overflow Question:

what Regular expression to validate HTML?

$\langle a \rangle \dots \langle /a \rangle$   
 $\langle h3 \rangle \dots \langle /h3 \rangle$

closed

- RegEx are good enough to describe "tokens"  
token.- a category of atomic strings.  
lexeme.- an instance of a token

e.g. token  $\begin{cases} \rightarrow \text{identifiers} \\ \rightarrow \text{integers} \\ \rightarrow \text{decimals.} \end{cases}$

e.g. lexemes foo 5 0.1

while | ( | flag | == | 62.3 | else | if | ) | 14 | 7 | != | x | + | var\_1 |  
↓            ↓            ↓            ↓            ↓            ↓            ↓            ↓            ↓            ↓            ↓  
keyword    identifier   op    float literal   keyword   keyword   int Lit   or   ident   op   ident.

- Describing tokens using RegEx  
 $\Sigma$  = keyboard symbols.

- Shorthands

$[0-9] \equiv (0|1|2|3|4|5|6|7|8|9)$   
 $[A-Z] \equiv (A|B|C|D|...|Y|Z)$   
 $[a-z] \equiv (a|b|c|d|...|y|z)$   
 $[A-Za-z] \equiv ([A-Z]|[a-z])$

- Regex for keywords

e.g C++ (if | while | for | do | else | .....)

- Regex for operators

(+ | \* | - | / | % | = | <= | >= | .....)

- Regex for integer literals.

$(+|-)?([0-9])^+$  =  $\left( \begin{array}{cccc} 1 & 2 & 0 & 00\ 001 \\ 27 & -5 & +10 & +00 \\ 235 & -12 & -001 & -00 \end{array} \right)$

- RegEx for identifiers.

Fortran  $([A-Z])([A-Z]|[0-9])^*$

A A2 K7F APPLE3

Pascal  $([A-Za-z])([A-Za-z]|[0-9]|'.'|'-'|^*)^*$

var1 point\_2D initX PascalCase

Python  $([A-Za-z]|'.'|'-'|^*)([A-Za-z]|[0-9]|'.'|^*)^*$

--init-- --main-- this\_is\_snake  
var\_1 -----

- RegEx for scientific notation.

... 2.7 -5 0.0 7 1.2 3 1.12 7e-0

- RegEx for scientific notation.

eg.  $3.7e^{-5}$      $-8.9e^7$      $123e^3$      $+12.7e^{-8}$

$(+|-)? [0-9]^+ (\ '.' [0-9]^+ )? 'e' (+|-)? [0-9]^+$

- RegEx can also be used to describe other string formats

- phone numbers

- SSN

- URLs

- Serial numbers

- e-mail addresses

- Licence plates

- etc..

→ • → EOF