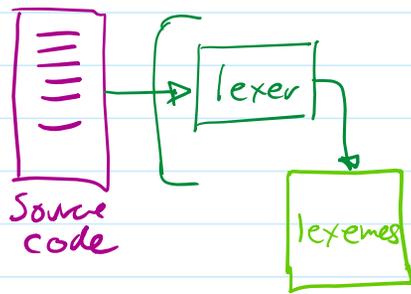


# 4 Syntax and Regular Expressions

Wednesday, February 4, 2026 3:00 PM



- Borrow from Linguistics & Mathematics

- Orthography : How to write words?

apple manzana pomme. 苹果

A complete list of words: Dictionary



in a programming language: Keywords.

- Syntax :- How to write sentences?

"has jump apple dog green"

Syntax has rules.

also in programming languages.

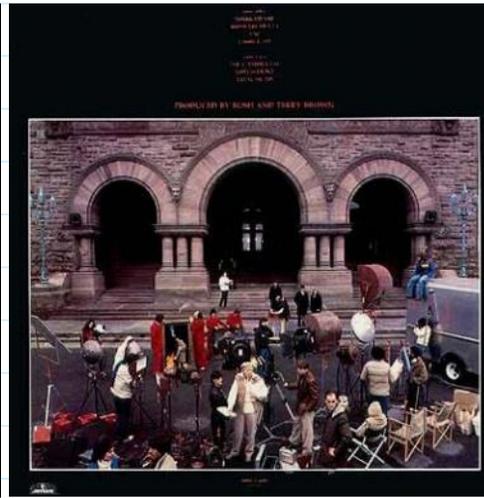
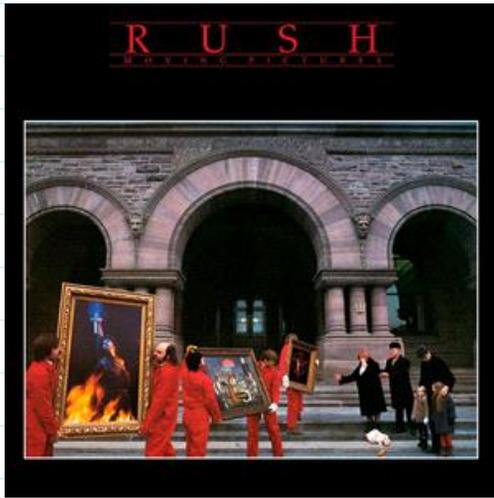
- Semantics :- What do sentences mean?

"colorless green ideas sleep furiously"  
Adj Adj noun verb Adverb

"the chickens are ready to eat"

 or 

"Moving Pictures"



In programming languages we should avoid

- non-sensical sentences
- ambiguous sentences.

- Pragmatics : How meaning changes over time?  
 "the container has to handle an unreliable cloud"

• Programming Language Orthography

What are the words of a programming language?

Symbols :- ASCII

- Keywords.
- Operators
- Literals :- constant values
- identifiers :- named entities.
  - Variable
  - constants
  - functions
  - classes.
  - modules
  - namespaces
  - templates.

How many identifiers are possible?  
 infinite !!

- We need rules!

We borrow from mathematics

"Mathematical Languages" S. Kleene  
 E. Post

DEF: An alphabet is a finite set of symbols

e.g.  $\Sigma = \{a, b\}$        $\Sigma = \{0, 1\}$



Eg.  $\Sigma = \{a, b\}$

- $(a|b)b = \{ab, bb\}$
- $a = \{a\}$
- $b = \{b\}$
- $(a|b) = \{a, b\}$
- $(ab)^* = \{\epsilon, ab, abab, ababab, \dots\}$   
↑  
empty string

Note: an infinite set is clearly specified by finite means

- $(aa|b)^* = \{\epsilon, b, aa, bb, aab, baa, bbb, aaaa, aabb, \dots\}$
- $aa = \{aa\}$
- $b = \{b\}$
- $(aa|b) = \{aa, b\}$

- Shorthands:
  - $r? \equiv (r|\epsilon)$  "to r or not to r"
  - $r+ \equiv r(r^*)$  "one or more repetitions of r"

e.g.  $a(b?)(a^+) = \{aa, aba, aaa, abaa, aaaa, \dots\}$

- $a = \{a\}$
- $b? = \{b, \epsilon\}$
- $a^+ = \{a, aa, aaa, aaaa, \dots\}$

- Regular Expression and Programming Languages

$\Sigma = \text{ASCII keyboard symbols}$

$L_{\text{C++}}$  = The set of all C++ programs      RegEx not enough.

Note: RegEx do not have the power to specify such set.

Why?

e.g.  $\{a^n b^n \mid n \in \mathbb{Z}^+\} = \{ab, aabb, aaabbb, aaaa bbbb, \dots\}$

consider matching brackets

$\{ ( ) [ ] \langle \rangle \}$

$\{ ( [ ] ) \}$

consider matching brackets

{ ( ) [ ] } < >

{ ( [ ] ) }

• Famous stack overflow question:

what RegEx for validating HTML?

<b> ..... </b>  
<h3> ..... </h3>

closed.

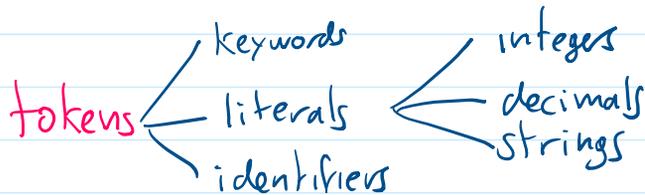
• RegEx can be used to validate

- Phone Number
- Email addresses
- URL
- Serial Numbers
- Credit Card Numbers.

• RegEx is good enough to specify **tokens**. !!

**token**:- a category of atomic strings

**lexeme**:- an instance of a token.



e.g. while foo 3.7

• Describing Tokens using RegEx

$\Sigma$  = ASCII keyboard characters.

• RegEx for keywords

C++ (if | while | for | switch | else | ... | void)

• RegEx for operators

(+ | \* | - | / | = | ...)

• Shorthands

[0-9]  $\equiv$  (0 | 1 | 2 | 3 | ... | 9)

[A-Z]  $\equiv$  (A | B | C | D | ... | Z)

[a-z]  $\equiv$  (a | b | c | d | ... | z)

$$\begin{aligned}
 [A-Z] &\equiv (A|B|C|D|\dots|Z) \\
 [a-z] &\equiv (a|b|c|d|\dots|z) \\
 [A-Za-z] &\equiv ([A-Z]|[a-z])
 \end{aligned}$$

• Regex for integer Literals

$$(+|-)?[0-9]+ = \left\{ \begin{array}{l} 15 \quad 27 \quad 54 \quad 333 \\ 1257 \quad 01 \quad 000 \quad +7 \quad +85 \\ -5 \quad -27 \quad -000 \quad +0 \end{array} \right\}$$

• Regex for decimal literals

$$(+|-)?[0-9]+\.[0-9]+ \left\{ \begin{array}{l} 17.0 \quad 0.03 \quad 6.3 \\ +0.0 \quad -11.3 \\ -0.0 \quad -27.0 \end{array} \right\}$$

• RegEx for identifiers

BASIC  $[A-Z]^+$  A ABC APPLE

Fortran  $[A-Z]([A-Z]|[0-9])^*$  A P2 K7F APPLE3

Pascal  $([A-Za-z])([A-Za-z]|[0-9]|_)^*$   
           Var1       point\_2D       init\_x       ThisIs Pascal Case.

Python  $([A-Za-z]_)([A-Za-z]|[0-9]|_)^*$        thisIs Camel Case  
 --main--       --init--       this\_is\_snake\_case  
 --plus--       var\_1       ----       this-is-kebab-case

• RegEx for Scientific Numbers

e.g 3.7e-5   -8.9e7   123e3   +12.7e-8   -5e+7

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

[1-9][0-9]\*  
2.5e107

WHENEVER I LEARN A NEW SKILL I CONCOCT ELABORATE FANTASY SCENARIOS WHERE IT LETS ME SAVE THE DAY.

OH NO! THE KILLER MUST HAVE FOLLOWED HER ON VACATION!



BUT TO FIND THEM WE'D HAVE TO SEARCH THROUGH 200 MB OF EMAILS LOOKING FOR SOMETHING FORMATTED LIKE AN ADDRESS!

IT'S HOPELESS!



EVERYBODY STAND BACK.



I KNOW REGULAR EXPRESSIONS.



PERL!

TAP TAP

