01/15/2025 ch2. scanner.

1.



source program → [ front end [ ch2. scaner ] IR → Opt IR → backend ] → target program. (compiler)
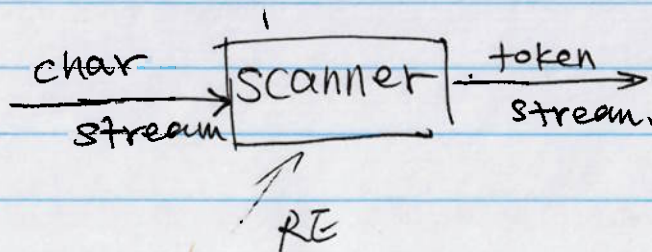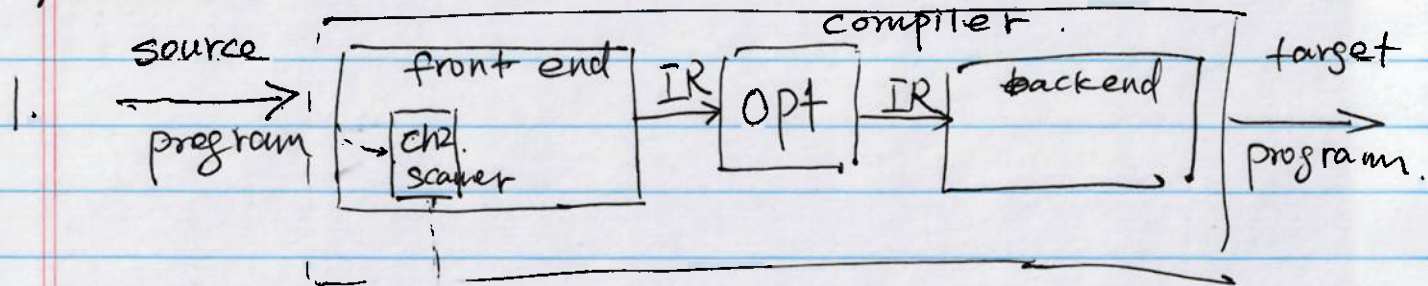
char stream → [scanner] token stream.
RE

2. tokens.

token : smallest element of a program which are identified by the compile that are meaningful.

different PLS. have different type of tokens. some. are. similar.

B. 2.1 Token. examples.

3†. Java :

⊔⊔⊔⊔ int  a = b + c ;

keywords :   int. if else.
identifiers :   names of variable. class. method.

literals:    constant.

　　　　 1　　　 3.14.　　 true　　 '1'　 "1"

　　　　 null

operator:　　　 +　 =　　 ++　 --　　 *=

separator.

　　　 - separate / punctuate other tokens.

　　　　　 ,　　 ;　　 .

　　　 - delimiters.

　　　　　 (　 )　 { }　 [ ]

12　　　　　 1 2　　　　　 "1 ⊔ 2"

X >= Y　　　　　 X >>> 3　 unsigned right
                                      shift.
f123　　 f ⊔ 123

Note: about ~~tokens~~ in Java.

①. Java always tries to construct the longest ~~token~~

②. white space :　 space. tab. newlines.
　　 - used to separate tokens.
　　 unless. it is inside stry literals.

③ comments.
　 // . single line
　 /*　 multiple line.
　 **
　 */

/** Java Doc

~~#*/~~

scanner recognize comments ∧ but exclude them
from further processing .

whitespace

↓    ↓    ↓    ↓    ↓ ↓ ↓↓ ↓ ~~⟷~~
public static void main (Stry[] args ~~☐☐~~ )

→ {    /* this */

     System. out. println ( "1 + 2 =" + (1+2) );
       ↑   ↑   ↑ ↑   ↑    ↑     ↑     ↑ ↑↑↑↑↑↑↑

→ }

2.2. tokens in python .

↓ ↓ ↓ ↓ ↓
if i > j : ↓   ↓
~~|4|~~ ~~|4|~~ ~~|4|~~   print (i)
            ↑   ↑
→ else:←
|4| |4| |4|   print (j)
print ("hello") ↑ ↑↑ ↑

# comment

                        |    else:     print (j) .

3. from characters to tokens.

① Recognize a single word. "new".

c ← nextChar( )

if c == 'n'

    c ← next char( )

    if c == 'e'

        c ← nextChar( )

        if c == 'w'

            return true

        else      retur false

    else

else.

② recognize several words.
    new.    not

③ recognize more complex words.

0        001        00        1,000

1_000

Tools from

we need formal language theory. automata

**4.**

Regular language
RE / FSA

context free language

context. sensitive language.

PDA

FSA { DFA
      NDFA

TM.

specify. tokens/:words

RL / RE /. FSA

$$RE \xrightarrow[\text{construction}]{\text{2.4.2} \atop \text{Thompson}} NFA \xrightarrow[\text{construction}]{\text{2.4.3} \atop \text{subset}} DFA$$

DFA minimization.
2.4.4
2.6.3

RE
2.3

2.4.1

2.4.

01/22/2025    RE

Regular expression ( 2.3 ).    lexical structure

1. definition.

Notation for specifying  microsyntax  such as the

spelling of    a    positive integer.

RE  over  an  alpabet $\Sigma$.  each RE  is  associated

with  a  set  of  strings  over  $\Sigma$.

~~~~~~~~~~~~~~~~~~~~~~

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

\* $\varepsilon$  is  a  regular expression.    empty string.

\* if $x \in \Sigma$, then  x  is a RE

e.g.    a,            b.          $L(a) = \{a\}$

$L(b) = \{b\}$.

\* if X and y  are  REs.  then

- concatination    $xy$  is a RE.

eg.  ab,    aa,    aab    $L(ab)$

$= \{L(a)L(b)$

$= \{a\}\{b\}$.

$L(xy) = L(x)\,L(y)$

$= \{pq : p \in L(x), q \in L(y)\} = \{ab\}$

$L(ba) = \{ba\}$.

- union / alternation.

$x|y$    is a RE

$L(x|y) = $ ~~$\emptyset$~~ $L(x) \cup L(y)$

$\qquad = \{ p : p \in L(x) \text{ or } p \in L(y) \}$

$L(a \mid ab) = \{ a, ab, \underline{\underline{~~aab~~}} \}$

$L(a) = \{ a \}$

$L(ab) = \{ ab \}$.

- klenee closure. $\cancel{X}^{*}$

$\quad L(x^{*}) = L(x)^{0} \cup L(x)^{1} \cup L(x)^{2} \cup \cdots$

$\quad L(x)^{0} = \{ \varepsilon \}$      ''''

$\quad L(x)^{i} = L(x)^{i-1} L(x)$

$L(a^{*}) = \{ \varepsilon, a, aa, aaa, aaaa, \cdots \}$

$L(ab)^{*} = \{ \varepsilon, ab, abab, ababab, \cdots \}$ ?.

$L(a|b)^{*} = \{ \varepsilon, a, b, ab, aa, bb, ba, \cdots \cdots$

$L(a|b)^{2} = L(a|b) L(a|b)$

$\qquad = \{ a, b \} \{ a, b \}$.

$\qquad = \{ \underline{a}\,\underline{a}, \underline{b}\,\underline{b}, \underline{a}\,\underline{b}, \underline{b}\,\underline{a} \}$

precedence    ( )     highest

$*$

concatination

alteration.     lowest .

$$L(ab^*) = \{a, ab, abb, abbb, \cdots \}$$

$$= L(a) L(b^*)$$
$$= \{a\} \{\varepsilon, b, bb, bbb, \cdots \}.$$

$$= \{ a\varepsilon, \ ab, \ abb, \ \ldots \ldots \}$$

$$x\varepsilon = \varepsilon x = x$$

2. closure properties of REs./RLs.

   REs/RLs are closed under
           concatination
           alternation
           kleene clousure$^*$

Note:   $x^+$ :    positive clousure

$$L(x^+) = L(x) \cup L(x^2) \cdots$$

$$ab^+ \equiv abb^*$$

$\overline{x^i}$ : finite closure.

$L(x^0) \cup L(x^1) \cup L(x)^2 \cdot L(x^i)$

not recommend.     two possible meaning:

$$L(x^i)$$

$1^{2}$      $\{ \mathcal{E}, 1, 11 \}$

$\{ 11 \}$

3. more examples.

① RE $\rightarrow$ RL.                    Languages.

$(0|1)(0|1)$                    $\{ 00, 01, 10, 11 \}$

$0(0|1)^*$                    $\{ 0, 00, 01, 000, 001 \cdots \}$

~~10~~  ~~110~~

$1\mathcal{E}$                    $\{ 1 \}$

$(1|01)^*$          $\{ \mathcal{E}, 1, 01, 11, 111, 101, 011$

every "0" has to be     $\underline{1}$    $\underline{01}$
followed by a '1'    $\underline{1\ 1}$    $\underline{01}\ 1$              $0101$
                   $\underline{1\ 01}$    $\underline{01\ 01}$
                   $\underline{111}$    .

$(1|01)^* (\varepsilon|0)$

$\varepsilon, 1, 0, 10, 01, 11, \cancel{00}$

$\underline{0}1\underline{0}$  $\underline{0}1\underline{1}$.   $\underline{1}\underline{0}1$  $\underline{1}\underline{1}\underline{0}$  $\underline{1}\underline{1}\underline{1}$

$\varepsilon$  or any binary strg w/o  two

consecutive "0's.

$219989$  $(0|1|2|3|4|5|6|7|8|9)$ (     ) (     ) (   )

$[0\cdots9]$                    $[0\cdots9]$ $[0\cdots9]$ $[0\cdots9]$

② . RL  $\Longrightarrow$   RE .

all binary strgs ending with "1"   $(0|1)^* 1$

$\{$ public, static, if $\}$.   $\cancel{\{}$ public$|$ static$|$ if$\cancel{\}}$.

unsigned integers w/o leading 0-s.

$[0\cdots9]^*$       $\cancel{\varepsilon}$

$[0\cdots9][0\cdots9]^*$     $\cancel{00}$ $\cancel{01}$

$[0\cdots9][2\cdots9]^*$     $\cancel{11}$ $\cancel{02}$

$[1\cdots9][0\cdots9]^*$     $\cancel{0}$

$0 | ([1\cdots9][0\cdots9]^*)$