

September, 18th

---

---

---

---



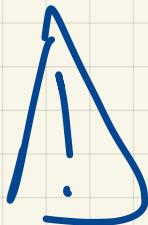
Gilles GEERAERTS

gilles.geeraerts@ulb.be

Office: Place du Campus

No building 8th floor  
N side, last office

TA : Nathalie Sessoles  
Leimond Brice



Timetables

## Uni-Nurste- Vink nelle web page

- lecture notes (PDF)
- Slides
- Project

## Evaluation

Project → group 2 people  
→ 3 parts 8 / 20  
No second choice!

Exam

12 / 20

## "Teaching style"

Lecture notes  $\approx$  text book  
complete  
mathematically  
precise -

Lectures  $\rightarrow$  higher level  
intuitions.

# What is a language?

Vocabulary  
Syntax

Way to communicate

"

Colours green ideas sleep furiously"

Syntax ✓

Semantics X

what we can  
write

Syntax

Meaning  
Semantics

"We can write a program that  
compiles but does not do what  
we want"

Syntax ✓

→ Semantics ✗

# Formal language

Definition: an **Alphabets** is **finite** set of symbols

$\{a, b, c, \dots, z\}$  : alphabet

$\{0, 1\}$  : alphabet

IN : not an alphabet

Definition: a **word** is a **finite** sequence of symbols from the alphabet

Alphabets:  $\{0, 1\}$

00110 → word  
110 → word

Empty word:  $\epsilon$

Definition A language is a set of words.

Example: Alphabet:  $\{0, 1\}$

Language = set of all binary numbers which are even.

$\{0, 10, 100, 110, \dots\}$

# What about programming languages?

alphabet

a ..	z	{ }
A ..	z	( )
Q ..	g	+,-,%
l	- - -	

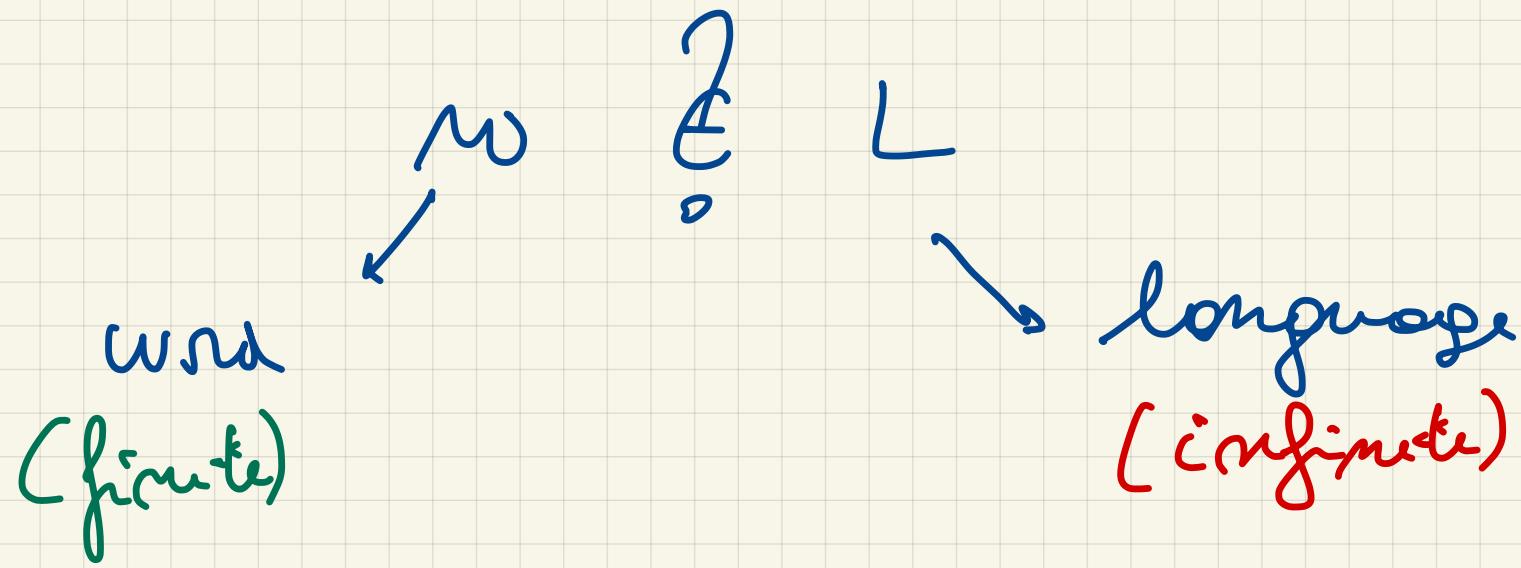
word

program.

language

set of all  
syntactically  
correct programs

# The membership Problem



1. The set  $L_{Cid}$  of all non-empty words on  $\Sigma_C$  (see example 1.7) that do not begin with a digit, is a language. It contains all valid C identifiers (variable names, function names, etc) and all C keywords (**for**, **while**, etc).
2. The set  $L_{odd}$  of all non-empty words on  $\{0, 1\}$  that end with a 1 is a language. It contains all the binary encodings of odd numbers.
3. Similarly to the previous example, the set  $L_0$  of all words on  $\Sigma = \{(, )\}$  which are well-parenthesised, i.e., s.t. each closing parenthesis matches a previously open and still pending parenthesis, and each open parenthesis is eventually closed. For example  $(00) \in L_0$ , but neither  $)()$  nor  $(()$  do.

This language is also known as the *Dyck language*, named after the German mathematician Walter von DYCK (1856–† 1934). It is mainly of theoretical interest: we will rely on it several times later to discuss the kind of formalism we need to recognise languages of expressions that contain parenthesis, such as the language  $L_{alg}$  defined in the next item:

$$\textcircled{1} \quad \Sigma_c = \{a, b, \dots, z, A, B, \dots, Z, 0, \dots, 9, -\}$$

all the words in  $\Sigma_c$  that do not start with a digit.

→ easy to check

just look at the first symbol

\textcircled{2}

0110101      easy to  
                  check.  
                  3

③

$$\Sigma = \{ , ) \}$$

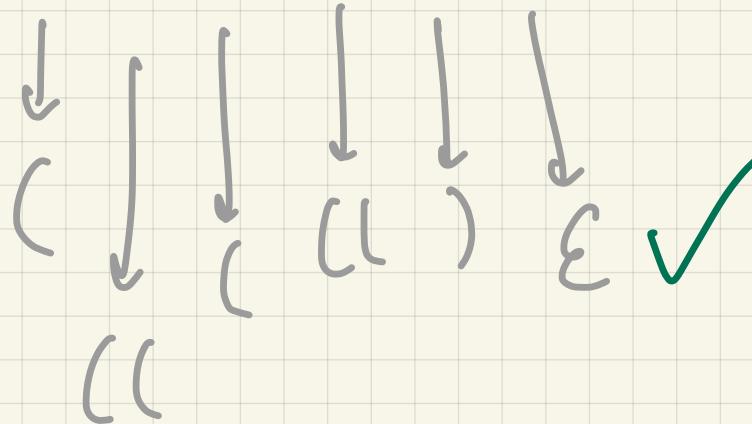
( ) ✓

( ( ) ( ) ) ✓

) ( X

(( )) X

(( ) ( ))



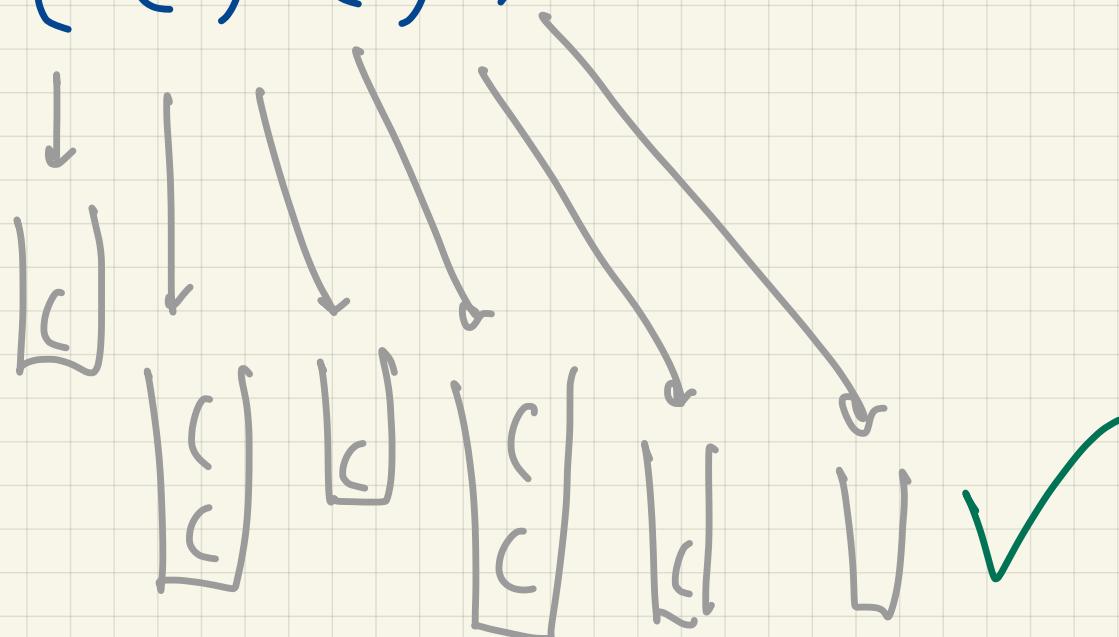
Queue  
versus -

( : add it to the queue

) : remove from the queue (if possible)  
otherwise  
reject

In the end: empty queue.

( ( ) ( ) )

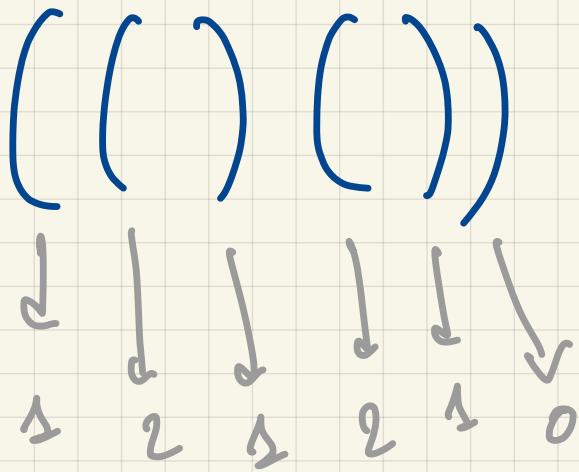


Stack  
version

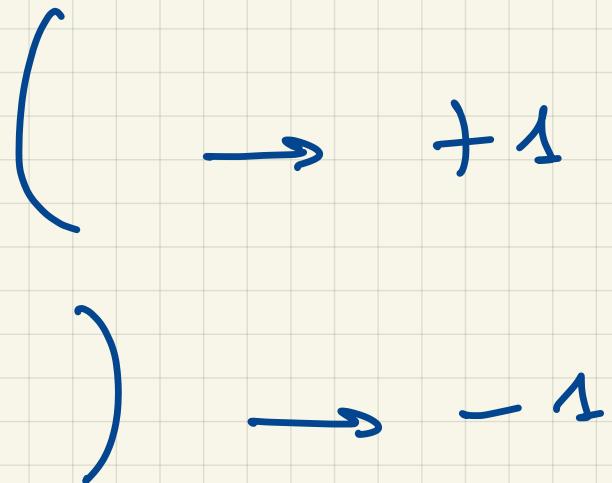
( → Push

) → Pop

Empty at the end.



Counter



$\neq 0$  at all  
times

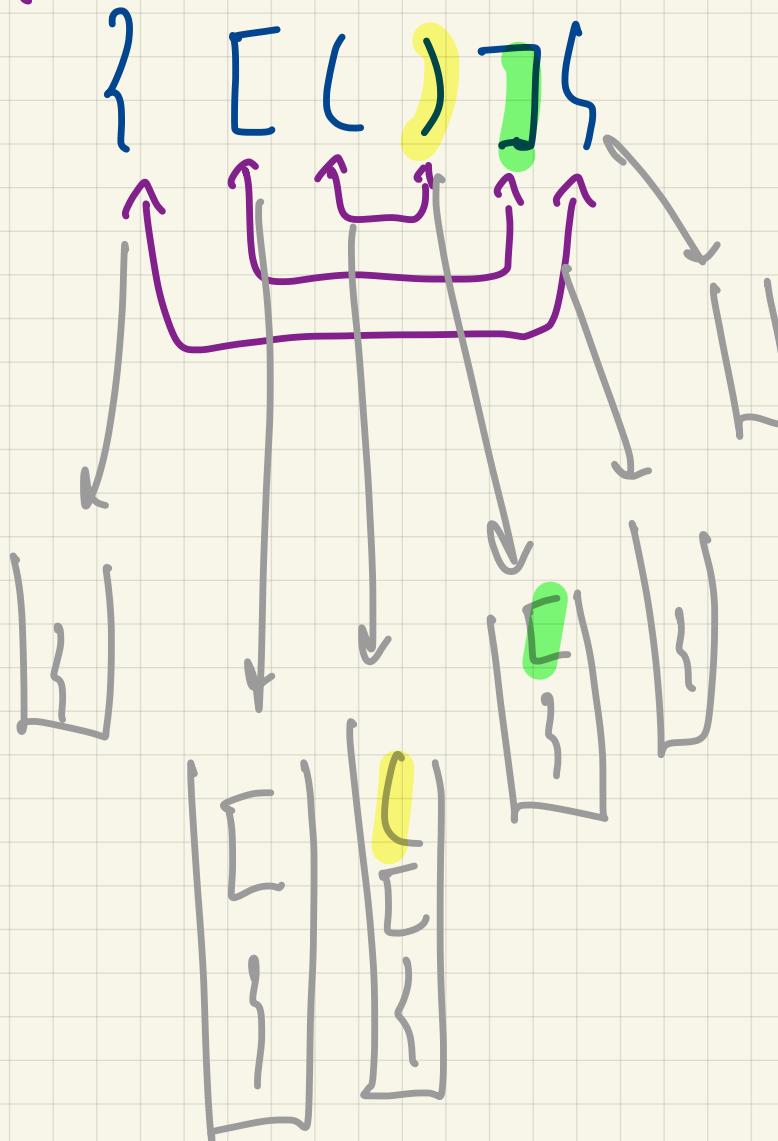
$= 0$  in the end



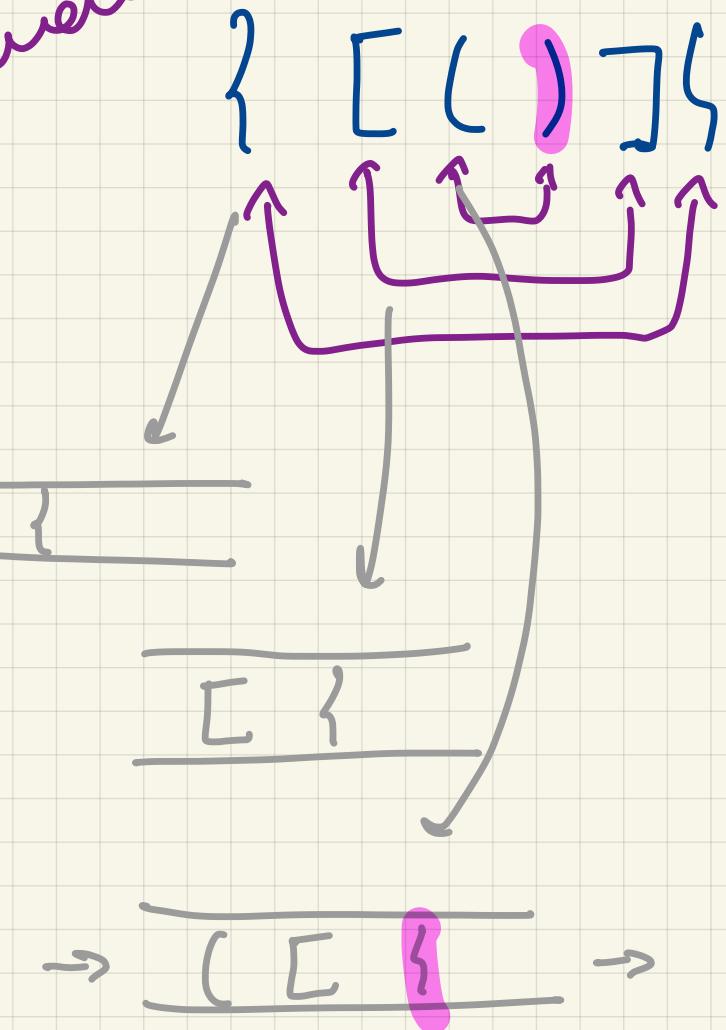
We need unbounded memory.

What if we lose ≠ brackets?

Stack



Queue



LIFO

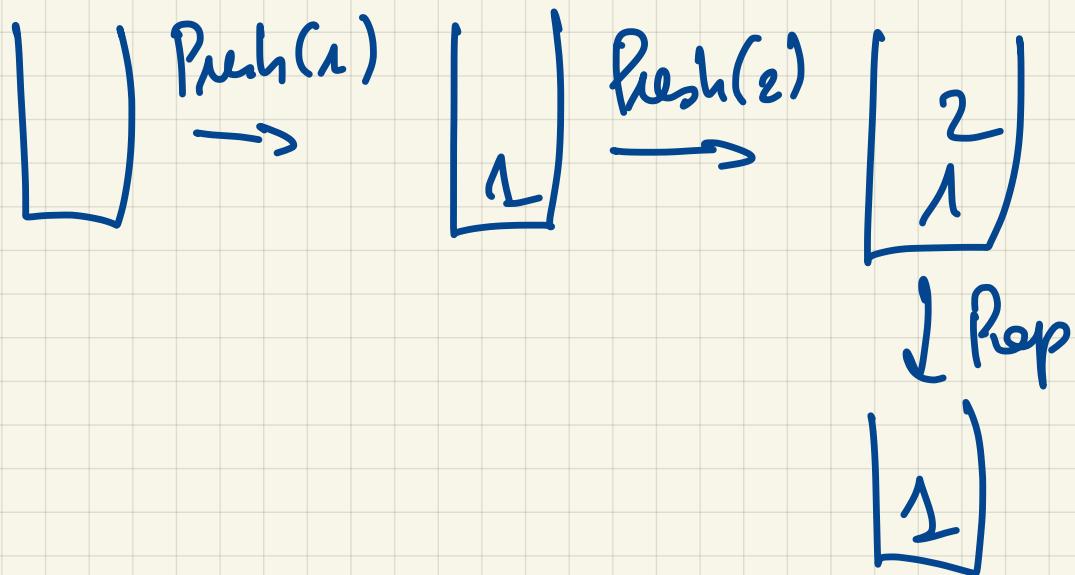
Stack

Push

Pop

Manipulate only  
the "Top"

One end of the  
sequence



FIFO

Queue

En queue

E(L) → 1

Dequeue

2 1  
DC ↗

Manipulate ≠  
ends of the  
sequence

More examples ! Test !

4. The set  $L_{alg}$  of all algebraic expressions that use only the  $x$  variable, the  $+$  and  $*$  operators and parenthesis, and which are well-parenthesised, is a language on the alphabet  $\Sigma = \{(,), x, +, *\}$ . For instance  $((x+x)*x)+x$  belongs to this language, while  $)(x + x$  does not, although it is a word on  $\Sigma$ .
5. The set  $L_C$  of all syntactically correct C programs is a language.
6. The set  $L_{Cterm}$  of syntactically C programs that terminate whatever the input given by the user is a language.

(

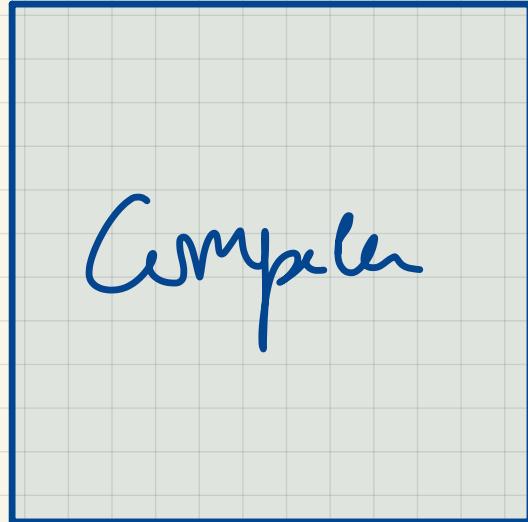
"beg free"

# Compilers

input  
source  
language →

C, C ++

java



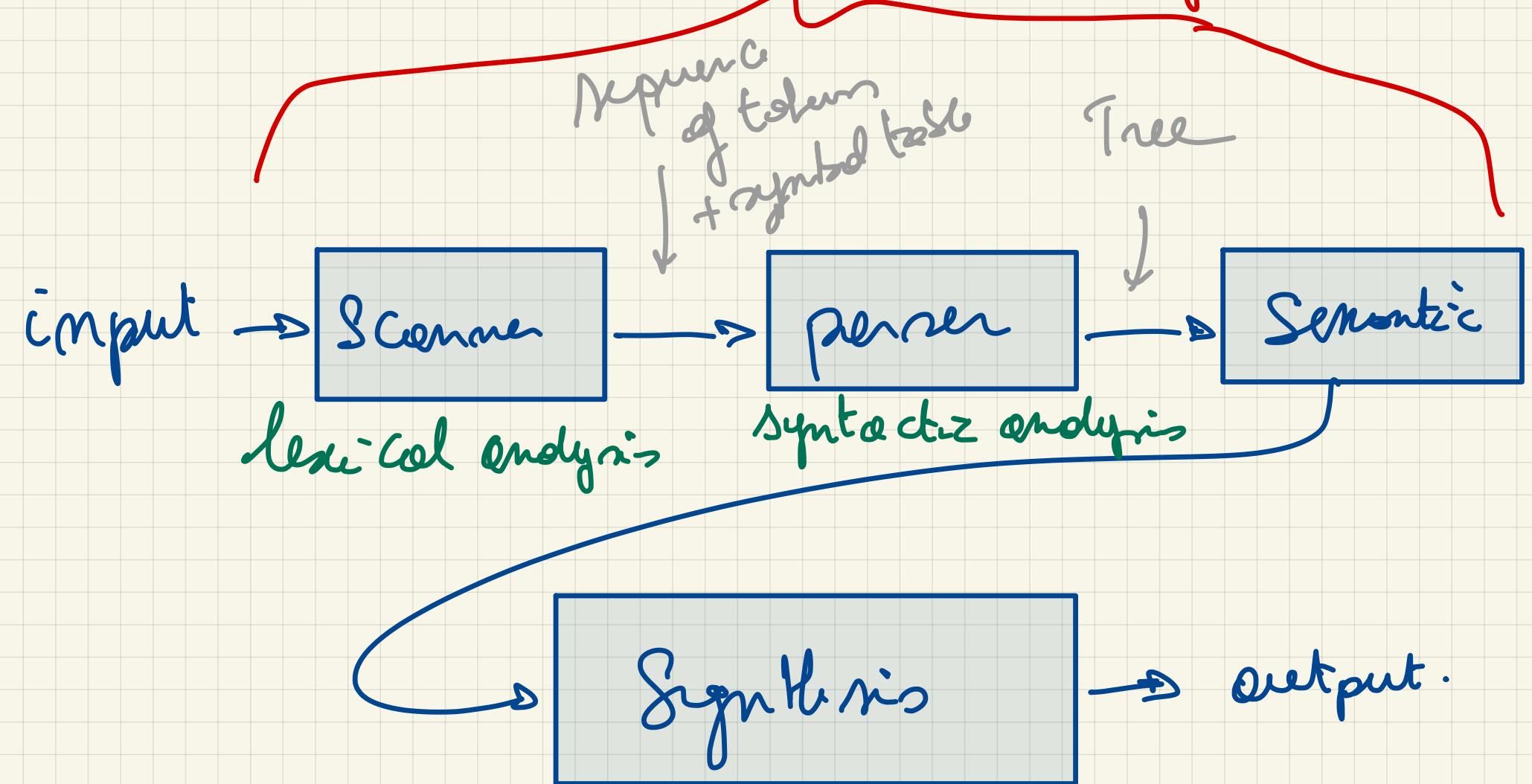
output  
target  
language

Machine  
Code

↓  
Check the syntax.  
translate

# Compiler

Analysis



# ① Scanning.

```
int i = 5 ;
```

```
int f ( int j ) {  
    int i = j ;  
    return i + 1 ;  
}
```

```
int main () {  
    printf ( "Hello_World_!" ) ;  
    printf ( "%d_%d" , i , f ( i + 1 ) ) ;  
    return 0 ;  
}
```

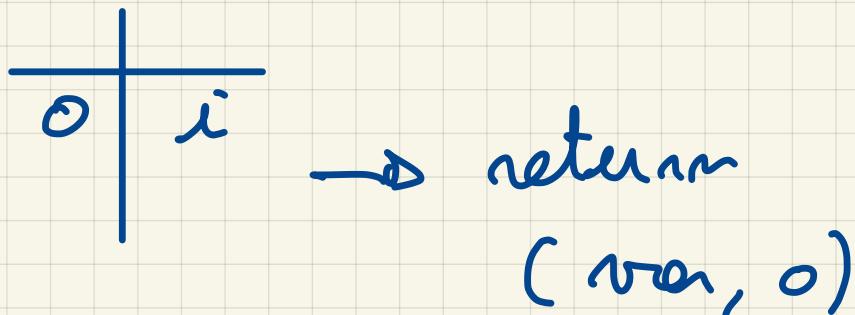
Split the input into tokens.

The scanner can build the symbol table  
= list of all identifiers

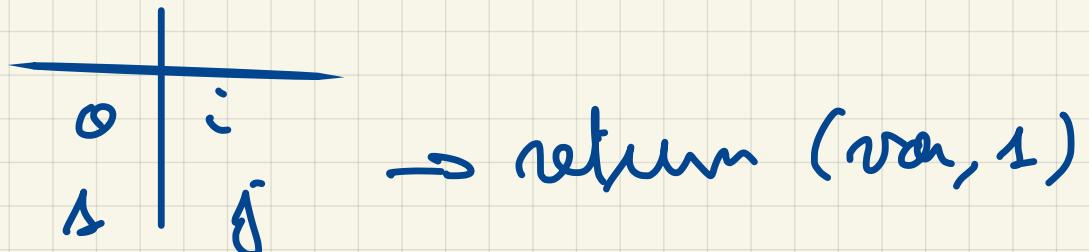
**Example 1.13.** Let us consider the simple code excerpt:

```
1 int i = 5;  
2 int j = 3 ;  
3 i = 9 ;
```

after line 1



after line 2



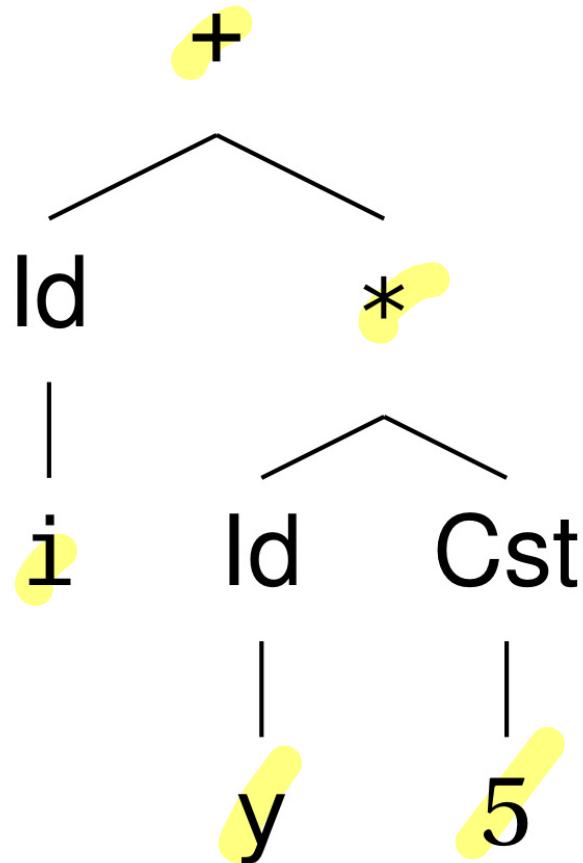
after line 3

// → return (var, 0)

## ② Parsing

Analyzing the global nature of the program.

d be:



$i + y * 5$

### ③ Semantics analysis.

#### → Scoping / Types

```
1 struct S {int i;} ;
2 int main() {
3     int i, j ;
4     struct S s ;
5     struct S * p = &s ;
6
7     i = 3 ;
8     j = i + 4 ;           int int
9     j = p ;             ✓
10    j = s ;            -i → type mismatch
11    j = s ;            X No conversion
12 }
```

Figure 1.2. Three syntactically correct as

# → Control flow

```
1 #include <iostream>
2
3 int main () {
4     for(int i=1;i<10;++i) {
5         infor: std::cout << i ;
6         std::cout << std::endl ;
7     }
8     goto infor ;
9 }
10 }
```

i exists  
in loop  
inside  
the  
control  
of the for

# Operations on words

$$w = w_1 \dots w_m$$

$$v = v_1 v_2 \dots v_r$$

- \* Concatenation of  $w$  and  $v$

$$w \cdot v = w_1 \dots w_m v_1 v_2 \dots v_r$$

$$abc \cdot de = abcd e$$

- \* Repeated concatenation

$$w^n = \underbrace{w \cdot w \cdot w \cdots w}_{n \text{ times}}$$

$$(ab)^3 = ab ab ab$$

## Operations on languages

$$L_1 \cdot L_2 = \{ w_1 \cdot w_2 \mid w_1 \in L_1, w_2 \in L_2 \}$$

$$\{ab, cd\} \cdot \{e, f\}$$

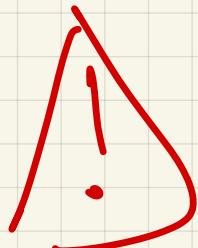
$$= \{ abe, abf, cde, cdf \}$$

?

.

$$L \cdot \{\epsilon\} = L = \{\epsilon\} \cdot L$$

$$L \cdot \emptyset = \emptyset$$



$\epsilon$   
empty word

$\{\epsilon\}$

↓  
non empty  
language not  
contains the empty word

$\emptyset$   
↓

empty language

$$L^m = \underbrace{L \cdot L \cdot L \dots L}_{m \text{ times}}$$

Kleen closure -

$$L^* = \{ w_1 w_2 \dots w_m \mid m \geq 0, \forall i: w_i \in L \}$$

$$\{\alpha\}^* = \{ \epsilon, \alpha, \alpha\alpha, \alpha\alpha\alpha, \dots \}$$

$\{0, 1\}^*$  = all binary words (including  $\epsilon$ )

the set of all representations of natural numbers

(in a given base, say 10)

$$\{0, \dots, 9\}^* \setminus \{\epsilon\}$$

# Regular languages

= family of languages

2 kinds of tools

Specification

Regular expressions

Machine ≈ code

Automata  
(finite)

**Definition 2.1** (Regular languages). Let us fix an alphabet  $\Sigma$ . Then, a language  $L$  is regular iff:

- 1. either  $L = \emptyset$ ;
- 2. or  $L = \{\epsilon\}$ ;
- 3. or  $L = \{a\}$  for some  $a \in \Sigma$ ;
- 4. or  $L = L_1 \cup L_2$ ;
- 5. or  $L = L_1 \cdot L_2$ ;
- 6. or  $L = L_1^*$

where  $L_1$  and  $L_2$  are regular languages on  $\Sigma$ .



*Bone*  
*Individually*

$$\Sigma = \{a, b, c, d, \epsilon\}$$

$$\{a\} \quad \{b\} \quad \{c\}$$

$$\{abc\}$$

$$\{abc, def\}$$

$$L = \{abc, def\}$$

$$\{d\} \quad \{c\} \quad \{f\}; \text{ R.L.}$$

$$\{def\}$$

(3)  
(5)  
(4)

Consequence: any finite language is regular.

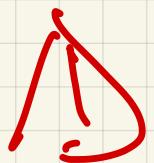
$$\begin{aligned} L &= \{ w_1 \dots w_m \} \\ &= \{ w_1 \} \cup \{ w_2 \} \cup \dots \cup \{ w_m \} \end{aligned}$$



Can be obtained by concatenation

$w_i$  is regular

$\Sigma^*$  = regular language which  
is infinite



$L \subseteq \Sigma^* \not\Rightarrow L$  is regular.

Remark

$L_{(,)} = \text{long resp. of well-parenthesized words}$

is not regular!

$$\{(, )\}^* \neq L_{(,)}$$

$$)( \in \{(, )\}^*$$

New operation on languages !!

$$L^+ = L \cdot L^*$$

$$= \{ w_1 \dots w_n \mid \forall i \cdot w_i \in L, n \geq 1 \}$$

So if  $L$  is regular, then  $L^+$  is regular.

**Definition 2.3** (Regular expressions). Given a finite alphabet  $\Sigma$ , the following are regular expressions on  $\Sigma$ :

1. The constant  $\emptyset$ . It denotes the language  $L(\emptyset) = \emptyset$ .
  2. The constant  $\varepsilon$ . It denotes the language  $L(\varepsilon) = \{\varepsilon\}$ .
  3. All constants  $a \in \Sigma$ . Each constant  $a \in \Sigma$  denotes the language  $L(a) = \{a\}$ .
  4. All expressions of the form  $r_1 + r_2$ , where  $r_1$  and  $r_2$  are regular expressions on  $\Sigma$ . Each expression  $r_1 + r_2$  denotes the language  $L(r_1 + r_2) = L(r_1) \cup L(r_2)$ .
  5. All expressions of the form  $r_1 \cdot r_2$ , where  $r_1$  and  $r_2$  are regular expressions on  $\Sigma$ . Each expression  $r_1 \cdot r_2$  denotes the language  $L(r_1 \cdot r_2) = L(r_1) \cdot L(r_2)$ .
  6. All expressions of the form  $r^*$ , where  $r$  is a regular expression on  $\Sigma$ . Each expression  $r^*$  denotes the language  $L(r^*) = (L(r))^*$ .

In addition, parenthesis are allowed in regular expressions to group sub-expressions (with their usual semantics). 

३

## Example

$$l = \alpha + b + c + \dots + z + A + B + \dots + Z$$

$$d = \emptyset + 1 + 2 + \dots + 9$$

C identifiers

$$l \cdot (l+d)^*$$

Each regular expression represents a regular language

Each regular language can be represented by a regular expression