


Sept 17th



Introduction to languages theory and compiling

Gilles

Mandato

Notthuen

GEERAERTS

DALACHANDER

SASSOLAS

Université
Vintello.



Lecture notes
Podcasts.

Evaluation

→ Exam January / August
/12

→ Project: write a compiler!

1/8 no second
run in!

You must achieve at least $\frac{1}{3}$ in part.

Schedule → v.v.

Longmap Theory

What is a longmap?

Word

longmap = set of words.

Rules



Colorless green ideas sleep furiously

Noam Chomsky.

Grammar: \checkmark \rightarrow syntax

Poem: \times \rightarrow semantics

Compiler? → check the syntax.

Semantics?

= "What the program
should do"

Formal definition of language

Alphabet : finite set of symbols

$\{0, 1\}$ ✓

$\{a, b, \dots, z, A, B, \dots, Z\}$ ✓

ASCII set ✓

IN X

Word : finite sequence of symbols on a given alphabet

alphabet $\Sigma = \{0, 1\}$

⚠ $\epsilon \notin \Sigma$

0 110 101

Empty word / sequence ϵ

Language : set of words

$\Sigma = \{0, 1\}$

$L = \{0, 1, 00, 11, 01, 10, \dots\}$
 $L = \{\epsilon, 0, 00, \underbrace{000, 0000, \dots}_{\text{infinite}}\}$

1. The set L_{Cid} of all non-empty words on Σ_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).

$a b 1 2$ ✓
 $- c 3$ ✓
 $9 a 1$ ✗

$$\Sigma_C = \{a, b, \dots, z, A, B, \dots, Z, 0, \dots, 9, -\}$$

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 $(()) \in L_0$, 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:

$$1. \quad w \stackrel{?}{\in} L_{\text{cid}} \xrightarrow{\text{infinite}}$$

No memory needed!

$$w = ab12$$

$$\uparrow \\ \in \{0, \dots, 9\} \rightarrow \text{no} \rightarrow \checkmark$$

$$2. \quad w \stackrel{?}{\in} L_{\text{odd}}$$

$$w = 01011$$

$$\uparrow \quad \dots \quad \uparrow \\ \text{no} = 1 \quad \checkmark$$

L()

($\rightarrow +1$
) $\rightarrow -1$ $\neq 0$

(() ()) ✓

1 2 1 2 1 0

↓

= 0 ✓

(()) X

1 2 1

↓

$\neq 0$ X

()) X

1 0 -1

↓

< 0 X

We need memory: counter.

Here are some ...

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 Σ .
5. The set L_C of all syntactically correct C programs is a language.
6. The set L_{Cterm} of syntactically correct C programs that terminate whatever the input given by the user is a language.

5. The set L_C of all syntactically correct C programs is a language.

1 Word = 1 program

language = all programs whose syntax
is correct.

Checking the syntax of a C program
 $\Leftrightarrow \text{prog} \in L_C$

Operations on words

Concatenation of words

$$W = w_1 w_2 \dots w_m$$

$$V = v_1 v_2 \dots v_l$$

$$W \cdot V = w_1 \dots \dots w_m v_1 \dots v_l$$

↑
concatenation

$$\varepsilon \cdot W = W = W \cdot \varepsilon$$

$$W^n = \underbrace{W \cdot W \cdot \dots \cdot W}_{n \text{ times}}$$

Operations on language

Concatenation

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

$$L_1 = \{aa, b\}$$

$$L_2 = \{c, dd\}$$

$$L_1 \cdot L_2 = \{aac, aadd, bc, bdd\}$$

$$L^n = \{w_1 w_2 \dots w_m \mid w_1, w_2, \dots, w_m \in L\}$$

$$L = \{a, b\}$$

$$L^2 = \{aa, ab, ba, bb\}$$

Kleene closure

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

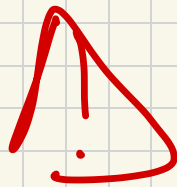
$$L = \{a\}$$

$$L^* = \{\epsilon, a, aa, aaa, \dots\}$$



$L = \phi \rightarrow$ empty language
it contains no word.

$L = \{\epsilon\} \rightarrow$ non-empty language!
it contains 1 word
which is empty.



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

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

does not exist

\emptyset

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

Regular languages

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

1. either $L = \emptyset$;
2. or $L = \{\varepsilon\}$;
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 Σ .



Base cases
Induction