

Introduction to theory of languages

Patryk Kiepas

MINES ParisTech
AGH University Of Science and Technology
THALES

February 25, 2017

Course plan

- ① Saturday, 25th of February 2017 – lecture
 - Languages
 - Grammars
- ② Saturday, 4th of March 2017 – lecture
 - Parsing
 - ANTLR
- ③ Saturday, 11th of March 2017 – exercises
 - Grammars and languages
 - ANTLR
- ④ Saturday, 25th of March 2017 – exercises & exam (???)
 - ANTLR

Additional informations

Any questions?

Ask by mail: kiepas@agh.edu.pl

Course web-page

<http://home.agh.edu.pl/~kiepas> → **Teaching** → **Introduction to theory of languages (2017)**

This lecture

- Linguistics
- Languages
- Grammar
- Hierarchy of grammars
- Notations
- Automaton

Introduction

Linguistics

Scientific study of languages. Involves analysis of language:

- *form* – language evolution and task
- *context* – environment of language usage
- *semantics* – the meaning of the language

Some important aspects

- Phonetics
- Articulation
- Perception
- Acoustic features
- Morphology
- Syntax

1 Natural languages

- *Ordinary* – evolves naturally in humans without planning
- *Controlled* – a restricted subset of natural language in order to reduce or eliminate ambiguity and complexity

2 Artificial languages

- *Constructed* (planned *a priori* or *a posteriori*)
 - Engineered languages – experiments in *logic*, *philosophy*, *linguistics*
 - Auxiliary languages – international communication (e.g. Esperanto, Ido, Interlingua)
 - Artistic languages – aesthetic pleasure or humorous effect (e.g. Klingon)
- **Formal**
 - Computer programming languages (e.g. Java, Haskell, C, C++, Ruby)
 - Files and formats descriptions (e.g. YAML, JSON, XML)

Description of natural languages

A really small bit of history

- In the late 1950's Noam Chomsky tried to describe natural languages
- Important paper: "*Three models for the description of language*", Noam Chomsky (1956).
- In a result of his research two disciplines originated:
 - 1 **Theory of formal grammars**
 - 2 *Generative (transformational) grammars*



Figure 1: Professor of Linguistics (Emeritus) at MIT, Cambridge

Description of natural languages

What we know now?

- Description of natural languages is **hard**
- Description of any natural languages might be **impossible**

Why this is important?

- Better understanding of language creation processes
- More insights into functioning of our brain
- **Natural language processing (NLP)**
 - Translations (e.g. Google Translator)
 - Synthesis (e.g. speech generation)
 - Perceiving (e.g. robots, voice-control)

Description of formal languages

Result

Description of natural languages help us describe an artificial (formal) ones

Programming languages

- Protocol for communication with the computer
- Performing operations and computations
- Interpretation and execution
- Compilation
- Static code analysis

Data formats

- Structured data
- Interchangeable model for communication and data transmission

Alphabet

Alphabet

A set Σ of available symbols, the simplest elements in the language

Examples

- binary alphabet $\{0, 1\}$
- decimal numbers $\{0, 1, 2, 3, \dots, 9\}$
- Latin alphabet $\{a, b, c, d, \dots, z\}$
- Cyrillic

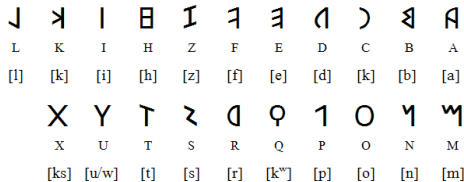


Figure 2: Ancient Latin alphabet

Word (I)

Word

Word w is a sequence of N symbols $w = x_1x_2...x_N$ (e.g. 010110, ABCDAAE)

Length

Length of word w is a number of symbols it consists of $|w| = N$ (e.g. $|010110| = 6$, $|ABCDAAE| = 7$)

Empty word

Special word ϵ with length $|\epsilon| = 0$

Word (II)

Words examples

- $w = 010110$ word over alphabet $\Sigma = \{0, 1\}$
- $w = abc13dj3$ word over alphabet $\Sigma = \{a, b, \dots, z, 0, 1, \dots, 9\}$
- $w = ACGTCCGGTA$ word over alphabet $\Sigma = \{A, C, G, T\}$

Closures

- Σ^* – set of all words over Σ
- Σ^+ – set of all nonempty words $\Sigma^+ = \Sigma^* \setminus \{\epsilon\}$

Closures examples

- if $\Sigma = \{a\}$ then $\Sigma^* = \{\epsilon, a, aa, aaa, aaaa, aaaaa, aaaaaa, \dots\}$
- if $\Sigma = \{a, b\}$ then $\Sigma^+ = \{a, b, aa, bb, ab, ba, aaa, bbb, \dots\}$
- if $\Sigma = \{a, b, \dots, z\}$ then $\Sigma^+ = \{cat, dog, a, aa, aaa, \dots\}$

- Alphabet Σ
- Word s , $|s| > 0$
- Empty word ϵ , $|\epsilon| = 0$
- Σ^* – set of all words over Σ
- Σ^+ – set of all nonempty words $\Sigma^+ = \Sigma^* \setminus \{\epsilon\}$

Grammar

- Description of a language
- A recipe for composing elements into sentence
- Describes syntax of a language

Each grammar describe some language.

Components

- 1 Terminals
- 2 Non-terminals
- 3 Production rules
- 4 Start symbol (non-terminal)

Grammar example

Examples

Digits separated by plus or minus signs

$$\textit{list} \rightarrow \textit{list} + \textit{digit}$$
$$\textit{list} \rightarrow \textit{list} - \textit{digit}$$
$$\textit{list} \rightarrow \textit{digit}$$
$$\textit{list} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

Chomsky's hierarchy

Grammar	Language	Production rules
Type-0 : Unrestricted	Recursively enumerable	$\alpha \rightarrow \beta$
Type-1 : Context-sensitive		$\alpha A \beta \rightarrow \alpha \gamma \beta$
Type-2 : Context-free		$\alpha \rightarrow \gamma$
Type-3 : Regular		$A \rightarrow a$ and $A \rightarrow aB$

Limiting condition

First condition

Production rules can't decrease length of word. $\forall (\alpha \rightarrow \beta) \in P : |\alpha| \leq |\beta|$

Second condition

Third condition

Type-0: Unrestricted grammar

Type-1: Context-sensitive grammar

Type-2: Context-free grammar

Type-3: Regular grammar

Backus-Naur form (BNF)

Backus-Naur form (BNF)

Notation technique for *context-free grammars*. Frequently used to describe syntax of *programming languages*, *document formats* etc.

Syntax

$\langle \text{term} \rangle ::= _ _ \text{expression} _ _$

- $\langle \text{term} \rangle$ is a *nonterminal*
- $_ _ \text{expression} _ _$ is a sequence of one or more terminal and/or nonterminal symbols separated by vertical line |
- Terminal symbols: a, b, c, A, 0, 1, 2 etc.
- Nonterminal symbols: $\langle \text{digit} \rangle$, $\langle \text{postal-code} \rangle$ etc.

Backus-Naur form (BNF)

Meta-symbols

- $::=$ – meaning "is defined as"
- $|$ – meaning "or"
- $\langle \rangle$ angle brackets used to surround category names.
- $'''$ –
- $\langle EOL \rangle$ –

Examples

$\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$\langle \text{postal-code} \rangle ::= \langle \text{digit} \rangle \langle \text{digit} \rangle \langle \text{digit} \rangle \langle \text{digit} \rangle \langle \text{digit} \rangle$

BNF example : Palindrome

Palindrome grammar

```
<letter>      ::= a | b | c | ... | y | z
<palindrome>  ::= <letter> |
<palindrome>  ::= a <palindrome> a | b <palindrome> b |
                  c <palindrome> c | d <palindrome> d |
                  e <palindrome> e | ...
                  | z <palindrome> z
```

Results

```
a
bb
bab
pop
hannah
```

BNF example : Postal address

Postal address grammar

```
<postal-address> ::= <name-part> <street-address> <zip-part>  
<name-part> ::= <first-name> <last-name> <EOL>  
<street-address> ::= <number> <street-name> <apt-num> <EOL>  
<zip-part> ::= <postal-code> <town-name> <EOL>  
<apt-num> ::= <number> | ""
```

Extended Backus-Naur form (EBNF)

