

# ID1020 – Algorithms and Data Structures

## Project 1 – VT14 P2

### 1 Organisation

The project is a programming task that is somewhat more involved than the lab. The results will be presented orally and a grade assigned based on the quality of the solution and understanding of the involved concepts.

The project itself is split into two parts, *Project 1 (P1)* and *Project 2 (P2)*, such that P1 forms the basis and P2 expands and improves on P1.

#### 1.1 Dates

**Submission** Wednesday, Dec. 18th, 23:59 in Bilda

**Presentation** Thursday, Dec. 19th, time-slots will be assigned.

#### 1.2 Goals

The project has the following goals:

- Work with the algorithms and data structures presented in the course
- Reason about usage patterns in a software system and leverage this to make implementation decisions
- Work on a real problem within the context of the course
- Get an idea of CFGs and their “real” applications

#### 1.3 Requirements

For the project you will need the following:

- Java
- Maven

cf. Lab 2 for details.

## 1.4 Time

This depends heavily on your experience in writing code for more involved projects.

We can only recommend to start early, to get a feeling on how much time you will need to invest.

Make sure to think through the problem first before you throw yourself head first into coding. A good design from the beginning can save you many hours of wading through convoluted code later. Especially since you will have to reuse parts of your solution for P2.

## 1.5 Notations & Definitions

We denote the set of all natural numbers with  $\mathbb{N} = \{1, 2, 3, 4, \dots\}$  and  $\mathbb{N}_0 = \{0\} \cup \mathbb{N}$ . Similarly  $\mathbb{R}$  denotes the set of all real numbers and  $\mathbb{R}^+$  the set of all positive real numbers. The set of common complexity classes is denoted as follows:

$$\mathcal{C} = \left\{ f : \mathbb{N} \rightarrow \mathbb{R} \mid f(n) \mapsto \begin{cases} 1 \\ \log n \\ n \\ n \log n \\ n^r \text{ for some } r \in \mathbb{R}^+ \\ r^n \text{ for some } r \in \mathbb{R}^+ \\ n! \end{cases} \right\}$$

Let  $\mathbb{R}^{\mathbb{N}}$  refer to the set of all functions  $f : \mathbb{N} \rightarrow \mathbb{R}$  then clearly  $\mathcal{C} \subseteq \mathbb{R}^{\mathbb{N}}$ .

## 2 Background

The goal of the project is to build a simple search engine for natural language text documents.

When completed the search engine will support the following operations:

- Index documents based on their content.
- Find and list documents which contain a single provided key word.
- Order search results by different properties (e.g. relevance, popularity).
- Use a *query language* to support more involved searches.

The documents themselves are provided in a specific format as part of the project framework and you don't have to worry about processing the raw files. However, it is important to notice that the given documents are english language texts of various types (e.g. short stories, newspaper articles, etc) and natural language texts differ in their treatment from more formal sources like programming languages, for example.

The specific documents we provide are from the so called Brown Corpus<sup>1</sup>. To give you some quick context, we provide a short summary of the field of natural language processing in section 2.1. Furthermore, the query language in the last task will require you to understand the concepts of parsing a limited syntax and to that end section 2.2 introduces the concept of context-free grammars and their relationship with parsing syntax.

## 2.1 Natural Language Processing

Natural language processing (NLP) is a field of computer science, artificial intelligence, and linguistics concerned with the interactions between computers and human (natural) languages. To a large extent the purpose of NLP is to allow computers to derive and subsequently leverage information from documents written for and by humans. This information could be anything from simply observing word usage pattern in certain kinds of documents, to understanding and classifying the content of said documents on a semantic basis. NLP has a number of sub-fields of which the following might be of interest to our search engine:

**Part-of-speech tagging** Given a sentence, determine the part of speech (POS) for each word. Many words, especially common ones, can serve as multiple parts of speech. For example, “book” can be a noun (“the book on the table”) or verb (“to book a flight”). Some languages have more such ambiguity than others. Languages with little inflectional morphology, such as English are particularly prone to such ambiguity. POS tagging is a necessary pre-requisite in most cases for any kind of syntactical analysis.

**Parsing** Determine the parse tree (grammatical analysis) of a given sentence. The grammar for natural languages is ambiguous and typical sentences have multiple possible analyses. In fact, perhaps surprisingly, for a typical sentence there may be thousands of potential parses (most of which will seem completely nonsensical to a human).

**Sentence breaking** (also known as sentence boundary disambiguation)

Given a chunk of text, find the sentence boundaries. Sentence boundaries are often marked by periods or other punctuation marks, but these same characters can serve other purposes (e.g. marking abbreviations).

**Word sense disambiguation** Many words have more than one meaning; we have to select the meaning which makes the most sense in context. For this problem, we are typically given a list of words and associated word senses, e.g. from a dictionary or from an online resource such as WordNet.

If you are interested, [http://en.wikipedia.org/wiki/Natural\\_language\\_processing](http://en.wikipedia.org/wiki/Natural_language_processing) give a good overview of the field.

---

<sup>1</sup>See [http://en.wikipedia.org/wiki/Brown\\_Corpus](http://en.wikipedia.org/wiki/Brown_Corpus) for further reading, if you are interested.

## 2.2 Context-free Grammars

Context-free Grammars (CFGs) have their roots in formal language theory and were originally invented in an attempt to formalise the grammatical structure of natural languages. Since it has been shown that not all natural languages can be described by CFGs they have somewhat lost their importance in the NLP field. However, they retain immense importance in computer science as they form the basis of all programming languages (and also query languages, which is why we introduce them here).

**Definition** Formally a CFG is defined by a set of *production rules* of the form:

$$S \rightarrow s_1 s_2 \dots s_n \text{ for } n \in \mathbb{N}_0$$

where  $S$  is a *nonterminal* symbol and the  $s_i$  are either nonterminal or *terminal*, meaning that  $s_i$  is a *token* of the alphabet of the language. Terminal symbols may never appear on the left-hand side of a production rule and each nonterminal symbol has to appear on at least one left-hand side of some production rule.

**Example** Consider the following grammar for simple arithmetical expressions:

- 1 :  $E \rightarrow num$  where  $num \in \mathbb{R}$  for example
- 2 :  $E \rightarrow E \cdot E$
- 3 :  $E \rightarrow \frac{E}{E}$
- 4 :  $E \rightarrow E + E$
- 5 :  $E \rightarrow E - E$
- 6 :  $E \rightarrow (E)$

Also consider the following example sentence  $s$  in the language of arithmetical expressions  $s = 2 \cdot (3 + 4)$ .

**Definition** In order to show that a sentence is part of the language of a CFG we have to find a *derivation*, that starting with the start symbol and applying production rules to the right hand side, we arrive at the sentence.

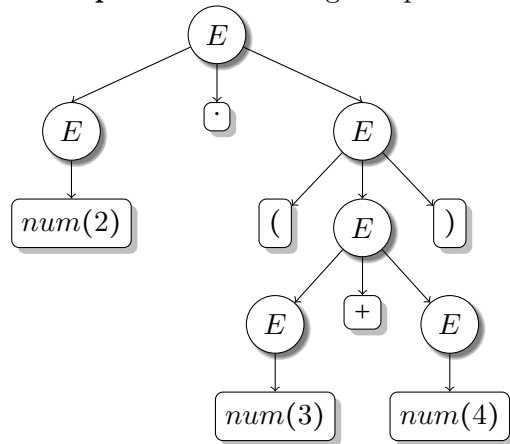
**Example** For the grammar above and sentence  $s$  the following would be a possible derivation:

$E$	start symbol
$E \cdot E$	rule 2
$num(2) \cdot E$	rule 1
$num(2) \cdot (E)$	rule 6
$num(2) \cdot (E + E)$	rule 4
$num(2) \cdot (num(3) + E)$	rule 1
$num(2) \cdot (num(3) + num(4))$	rule 1

Since we have derived  $s$  from the start symbol we have shown that  $s$  is part of the given grammar.

**Definition** While it is not difficult for a human to find a derivation that matches a sentence, a program would have try a (possibly large number) of different derivation until it by chance finds one that matches the input sentence. Since this is not feasible, programs usually apply the opposite process to derivation, which is called *parsing*. The result of the parsing process is a *parse tree* with the start symbol on top and the syntactical structure of the sentence represented in the structure of the tree. If parsing consumed all the symbols of a sentence that sentence is part of the language of the grammar.

**Example** The following is a possible parse tree for  $s$ :



**Note** The quick introduction to CFGs presented here is certainly far from exhaustive. There are important issues of ambiguity and associativity of operations we have not considered at all so far. However, it should be sufficient to perform the required tasks. If you require more information or are simply interested feel free to look at [http://en.wikipedia.org/wiki/Context-free\\_grammar](http://en.wikipedia.org/wiki/Context-free_grammar) and follow references from there.

### 3 Tasks

#### 3.1 Indexing

Write description

TODO

#### 3.2 Simple Search

Write description

TODO

#### 3.3 Sorting

Write description

TODO

#### 3.4 Queries

Your query language should support *union* and *ordering*. It is described by the following CFG:

$E \rightarrow T$	simple query
$E \rightarrow T \text{ orderby } Property \ Direction$	ordering
$T \rightarrow word$	a search term
$T \rightarrow T \ T$	union
$Property \rightarrow relevance$	how often do the search terms appear in a document?
$Property \rightarrow popularity$	how popular is the document?
$Direction \rightarrow asc$	increasing order
$Direction \rightarrow desc$	decreasing order

Write the rest of the description

TODO