

# Chapter 1

## Software requirement specification

### 1.1 Introduction

#### 1.1.1 Purpose

The software requirement specification should provide all needed information to develop the context extraction framework and define all delivery objects. All interfaces to external components, input and output data, deployment considerations and quality attribute should be well defined within this document.

#### 1.1.2 Scope

The context extraction framework will perform automated text extraction on a set of HTML test data with two to three different text extraction algorithms. The performance of each algorithm is measured and an output file with the measured results is generated.

## 1.2 General description

### 1.2.1 Operating Environment

#### 1.2.1.1 Local environment

JDK	1.7.X
Gradle	1.1
Eclipse Keppler	2.X
git	1.9.X
python	2.7.X

#### 1.2.1.2 Continuous Integration Environment

Open JDK	1.6.X
Open JDK	1.7.X
Oracle JDK	1.7.X
Oracle JDK	1.8.X
Gradle	1.1
Travis CI	

### 1.2.2 Design and Implementation Constraints

#### 1.2.2.1 User interface

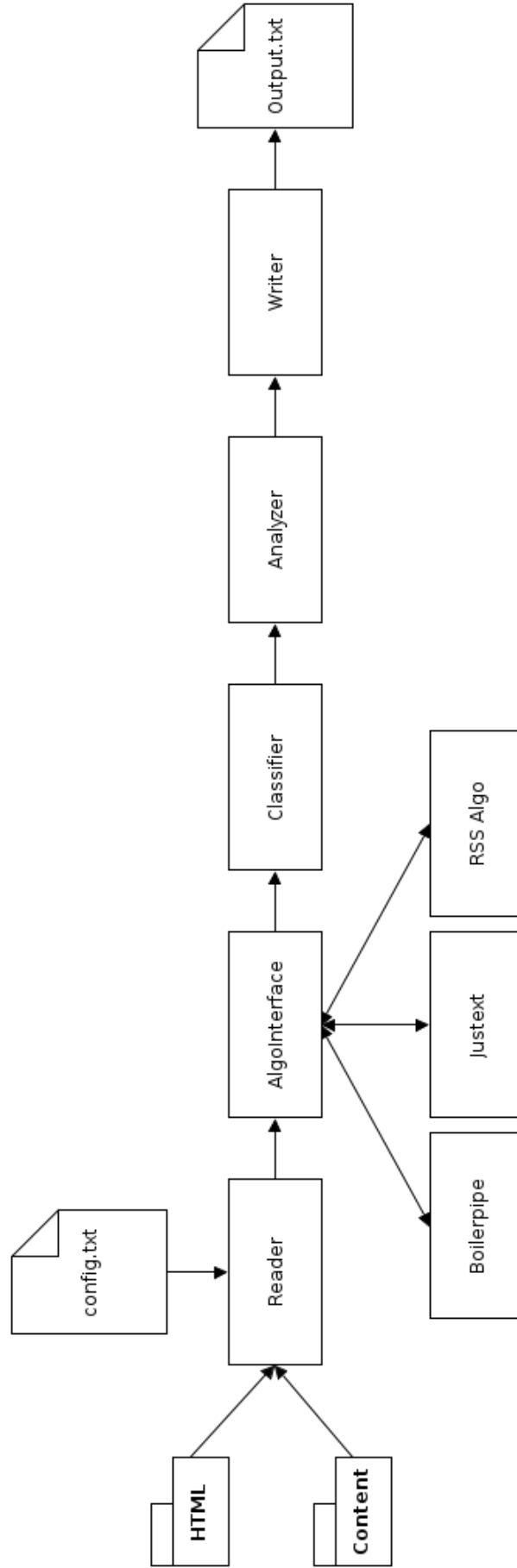
As parts of the text extraction framework may be implemented in a server environment and a user interface is not desired from the client so there will be no graphical user interface. The application is built, deployed and started by gradle. While the application is running, no interaction with the user is needed.

## 1.3 System Features

This section specifies all system features. Each feature is specified more close with multiple user stories but all important information such as external dependencies and output files are defined in this chapter. The related user stories are located in the planning section.

### 1.3.1 Basic functionality

Following diagram and text describes the basic functionality of the application.



There are two folders defined by a configuration files which. The HTML folder contains plain HTML files of web pages. The content folder contains text files with the relevant content of the related HTML files. As soon as a test is started the HTML file and the text file are read and the HTML file is extracted and classified with all the available algorithms. The result of the classification is then compared to the text file with the relevant content and performance data is generated. This performance data is then analyzed with statistical methods.

### 1.3.2 Overview

ID	Name	Relevance
f1	Read configuration	needed
f2	Create test	needed
f3	Integration Justext algorithm	needed
f4	Integration Boilerpipe algorithm	needed
f5	Evaluation and Implementation RSS feed algorithm	nice to have
f6	Evaluation of Classification Text	needed
f7	Evaluation of Classification Block	nice to have
f8	Analyze Data	needed

#### 1.3.2.1 Read configuration

<b>Name</b>	Read configuration
<b>Feature id</b>	f1
<b>Description</b>	<p>The text extraction framework is configurable with an external text file. The configuration file will contain following items:</p> <ul style="list-style-type: none"><li>• Path to folder with HTML files</li><li>• Path to folder with text files</li><li>• Path to folder with output files</li><li>• Configuration for algorithms</li><li>• etc.</li></ul> <p>The configuration file location is defined as a relative path to the source directory and structured in a key value list:</p> <pre>key:value; key:value; key:value;</pre>
<b>Relevance</b>	needed
<b>Related stories</b>	tbd

### 1.3.2.2 Create test

<b>Name</b>	Create test
<b>Feature id</b>	f2
<b>Description</b>	A test contains two input files which are a HTML file and a text file. They are located in the directories defined by the configuration. As soon as the test framework finds an HTML and a text file with the same name, a new test is created, the files are read and the test is started.
<b>Relevance</b>	needed
<b>Related stories</b>	tbd

### 1.3.2.3 Integration Justext algorithm

<b>Name</b>	Integration Justext algorithm
<b>Feature id</b>	f3
<b>Description</b>	Justext is implemented in python so a service is needed to call the python script and get the extracted text or the extracted blocks.
<b>Relevance</b>	needed
<b>Related stories</b>	tbd

### 1.3.2.4 Integrate Boilerpipe algorithm

<b>Name</b>	Integration Boilerpipe algorithm
<b>Feature id</b>	f4
<b>Description</b>	Boilerplate is implemented in Java so an interface is needed to call the Boilerplate component and get the extracted text or the extracted blocks.
<b>Relevance</b>	needed
<b>Related stories</b>	tbd

### 1.3.2.5 Integrate RSS feed algorithm

<b>Name</b>	Evaluation and implementation RSS feed algorithm
<b>Feature id</b>	f5
<b>Description</b>	The basic idea of the RSS feed algorithm is to match the content of a HTML document with the related RSS feed and define the relevant content like that. This need to be evaluated, implemented and integrated into the text extraction framework
<b>Relevance</b>	nice to have
<b>Related stories</b>	tbd

### 1.3.2.6 Evaluation of classification text

<b>Name</b>	Evaluation of classification
<b>Feature id</b>	f6
<b>Description</b>	<p>All the text extraction algorithms return an extracted document as text. This document needs to be checked for correctness. To do so the result from the algorithms is compared with the predefined content. This evaluation and classification is defined in more detail here: <a href="#">1.3.3</a></p> <ul style="list-style-type: none"><li>• Check each classified block from the algorithms if it's content can be found in the content file</li><li>• Categorize all blocks as true negative and false positive</li><li>• Put the results into an output text file (structure output file: tbd)</li></ul>
<b>Relevance</b>	needed
<b>Related stories</b>	tbd

### 1.3.2.7 Evaluation of classification blocks

<b>Name</b>	Evaluation of classification blocks
<b>Feature id</b>	f6
<b>Description</b>	<p>A more detailed evaluation of the algorithms could be done if not only the text is classified but each block of an HTML file. To do so, the implementation of Justext and Boilerpipe have to be adapted that they return classified blocks instead of the extracted text. These blocks are then compared with the predefined content and classified. This evaluation and classification is defined in more detail in following section: <a href="#">1.3.3</a></p> <ul style="list-style-type: none"><li>• Check each classified block from the algorithms if it's content can be found in the content file</li><li>• Categorize all blocks as true negative and false positive</li><li>• Put the results into an output text file (structure output file: tbd)</li></ul>
<b>Relevance</b>	nice to have
<b>Related stories</b>	tbd

### 1.3.2.8 Analyze data

<b>Name</b>	Analyze data
<b>Feature id</b>	f7
<b>Description</b>	From the results of the comparison several further values can be evaluated for a better understanding of the results. These values are described in more detail in following section: <a href="#">1.3.3</a>
<b>Relevance</b>	needed
<b>Related stories</b>	tbd

<b>Name</b>	Visualize data
<b>Feature id</b>	f8
<b>Description</b>	The calculated values from feature f8 are visualized in diagrams. (tbd: which tool)
<b>Relevance</b>	tbd
<b>Related stories</b>	tbd

### 1.3.3 Further explanation for evaluation of classification

The general meaning of the expressions true positive, true negative, false positive and false negative related to the text extraction topic is shown in following table:

	<b>Classified as content</b>	<b>Classified as boilerplate</b>
<b>Actual content</b>	True positive (TP)	True negative (TN)
<b>Actual boilerplate</b>	True negative (TN)	False negative (FN)

#### 1.3.3.1 Evaluation the results as text

When the results are compared based on the text, the expressions are interpreted as follow:

	<b>Classified as content</b>	<b>Classified as boilerplate</b>
<b>Actual content</b>	Number of words classified as content and are content	Number of words classified as content but are boilerplate
<b>Actual boilerplate</b>	Number of words classified as content but are boilerplate	Number of words classified as boilerplate and are boilerplate

#### 1.3.3.2 Evaluation the results as blocks

When the results are compared based on HTML blocks, the expressions are interpreted as follow:

When the results are compared based on the text the expressions are interpreted as follow:

	<b>Classified as content</b>	<b>Classified as boilerplate</b>
<b>Actual content</b>	Number of blocks classified as content and are content	Number of blocks classified as content but are boilerplate
<b>Actual boilerplate</b>	Number of blocks classified as content but are boilerplate	Number of blocks classified as boilerplate and are boilerplate

#### 1.3.3.3 Analytical values

The results of the comparison deliver basic characteristics which can be used to calculate statistical values which help to analyze the test outcome.

#### Recall / Sensitivity / True positive rate / TPR

Recall is the probability that a relevant document is retrieved in a search.

$$Recall = \frac{TP}{TP + FN} \quad (1.1)$$



**Precision / True negative rate / TNR**

Precision is the probability that a retrieved document is relevant.

$$Precision = \frac{TP}{TP + FP} \quad (1.2)$$

**F-measure / F1-score / F-score**

F-measure is a measure of the test's accuracy so it is the weighted harmonic mean of precision and recall.

$$Fmeasure = 2 * \frac{precision * recall}{precision + recall} \quad (1.3)$$

**Fallout / False positive rate / FPR**

Fallout is the proportion of non-relevant documents that are retrieved, out of all non-relevant documents available.

$$Fallout = 2 * \frac{FP}{FP + TN} \quad (1.4)$$

**Reciever Operation Characteristics (ROC)**

The ROC curve is the sensitivity as a function of fall-out which illustrates the performance of a classifier. It can be used to find the best results using different parameters.

$$ROC = TPR = f(FPR) \quad (1.5)$$

## 1.4 External Interface Requirements

### 1.4.1 Boilerpipe

The boilerpipe algorithm is already implemented in Java so it is easy to integrate. The API can be found under following link. <https://code.google.com/p/boilerpipe/>

Other useful links:

Getting started: <http://code.google.com/p/boilerpipe/wiki/QuickStart>

javadoc extractor: <http://boilerpipe.googlecode.com/svn/trunk/boilerpipe-core/javadoc/1.0/de/l3s/boilerpipe/extractors/ExtractorBase.HTML>

### **1.4.2 justext**

The justext algorithm is implemented in python and it is not yet defined how it will be integrated into the text extraction framework. See risk analysis for further information. The documentation can be found under following link:

<https://code.google.com/p/justext/>

jython: <http://www.jython.org/>

## Chapter 2

# Risk analysis

### 2.1 Introduction

#### 2.1.1 Purpose

This document evaluates and weights all possible risks and defines actions to minimize them as good as possible.

### 2.2 Risk evaluation

#### 2.2.1 Unclear requirements

Requirements are somehow vague at the beginning of each project and if they are not well defined as soon as possible, they stay vague trough out the whole project and this can lead to a disaster.

#### 2.2.2 New technologies

The new technologies which are present in this project are:

- Gradle
- Travis CI
- Python

Each of them brings his own risk.

### 2.2.3 Interface Boilerpipe

The Boilerplate algorithm needs to be integrated into the text extraction framework. Every interface of an external component is a possible risk factor.

### 2.2.4 Interface Justext

The Justext algorithm needs to be integrated into the text extraction framework. Every interface of an external component is a possible risk factor.

### 2.2.5 Implementation RSS algorithm

The development and implementation of a new algorithm is predestined to generate risks.

## 2.3 Assessment of risks

Risk	Impact	Probability of occurrence	Risk factor
Unclear requirements	2	4	8
New technologies	3	3	9
Interface Boilerpipe	5	1	5
Interface Justext	5	5	20
Implementation RSS algorithm	1	5	5

## 2.4 Consequences

### 2.4.1 Unclear requirements

As I am working with the client each and every day, it is very easy prevent misunderstandings with asking the client at once. Even though misunderstandings can occur between student and expert. To prevent this, it is necessary to have a document to define the requirements as soon and as exact as possible. This will be done in the form of the system requirement specification in the first mile stone. Possible ambiguities can be clarified at the first mile stone meeting.

### 2.4.2 New technologies

It is important to do prototyping with new technologies in the first phase of the project to eliminate these risks as soon as possible.

Gralde and Travic CI are needed in the first mile stone to set up the programming environment. So if there is any problem it will occur in a very early stage of the project and a possible solution can be found.

The risks about python are related to the chapter [2.4.4](#).

### 2.4.3 Interface Boilerpipe

This risk is rated much lower than the Justext interface because it's implementation is in Java and it provides a Java API. Never then less a prototype should be done as soon as possible to prevent any nasty surprises with the interface.

### 2.4.4 Interface Justext

This point is classified as the highest risk of all. This is because the implementation is in Python and it is not clarified yet how it will be integrated into the text extraction framework. An analysis of possible solution with prototypes needs to be done as soon as possible.

Possible solution are:

- jython (<http://www.jython.org/>)
- Implementation in Java
- Java Processor Interface

### 2.4.5 Implementation RSS algorithm

This risk has a very high probability of occurrence because it is very likely that a development and an implementation of a new algorithm is going to cause problems. There is no real solution to that risk. But because of this requirement is nice to have, the impact on the outcome of the project is very low. Further more Patrik Lengacher, the tutor of this project, is very experienced in this subject area and will be able to help out if any problems occur.