

# **Seeker Project Documentation**

Jaroslav Hirniak

Last updated on: *July 4, 2014*

# Contents

<b>1</b>	<b>Project</b>	<b>2</b>
1.1	Motivation . . . . .	2
1.2	Stakeholders . . . . .	2
1.3	Project Name . . . . .	2
<b>2</b>	<b>Requirements</b>	<b>3</b>
2.1	Minimal requirement . . . . .	3
2.2	Data format . . . . .	3
2.3	Use cases . . . . .	3
<b>3</b>	<b>Requirements analysis</b>	<b>11</b>
3.1	Objectives . . . . .	11
3.2	Data layer . . . . .	11
3.3	Application layer . . . . .	13
<b>4</b>	<b>Architecture</b>	<b>14</b>
4.1	Overview . . . . .	14
4.2	Full stack . . . . .	16
4.3	System design . . . . .	16
4.4	Query language . . . . .	16

# Chapter 1

## Project

### 1.1 Motivation

The researches in minority languages need to extract comparative information from sources varying in format and the content. These sources of information are usually very rich, but lack common structure. Most often useful data is spread among documents in PDF. Usually the sources describe many languages (more than 50 for European Committee website) and have different organization and approach to presenting information. Therefore, extracting information from such source is very elaborative and repetitive, being exactly what computers are the best in.

The goal of the project is to provide easy access to all the documents using convenient query language.

### 1.2 Stakeholders

- Researchers - people who are using service often for their research;
- Consumers - people who need to access data occasionally and do not have need to register an account in the service;
- Administrators - people having rights to add sources, respond to tickets, etc.;
- Developer - person responsible for the project development and maintenance.

### 1.3 Project Name

The project name *Seeker* comes from “Seeker, Reaper” poem by George Campbell Hay. You can listen to the poem being recited on YouTube [here](#).

## Chapter 2

# Requirements

### 2.1 Minimal requirement

The goal for 8 week summer project was set to

- carry out research and recommend appropriate solution (architecture, full software stack, and persistence layer);
- Implement a proof of concept functionality to prove the possibility of developing efficient solution for the stated problem, this is, ability to query European Committee for Minority Languages and receive combined documents composed of elements specified in the query.

### 2.2 Data format

The main source of the data are going to be the European Charter for Regional or Minority Languages expert reports<sup>1</sup>.

The goal of the project is to account for the future grow and to include also documents from Gaelic Plans produced by various bodies in Scotland<sup>2</sup>.

The data needs to be assumed to do not have any common format or fields.

### 2.3 Use cases

During the initial meeting on 26 May 2014 the below use cases had been identified.

---

<sup>1</sup>[http://www.coe.int/t/dg4/education/minlang/Report/default\\_en.asp](http://www.coe.int/t/dg4/education/minlang/Report/default_en.asp)

<sup>2</sup><http://www.gaidhlig.org.uk/bord/en/our-work/gaelic-language-plans.php>

Use Case ID:	0
Use Case Name:	Database query
Date Created:	26 May 2014
Date Last Updated:	26 May 2014

Actors:	Consumer
Description:	Actor wants to extract information from the database relating one or multiple documents that were present on source listed web pages.
Preconditions:	<ul style="list-style-type: none"> <li>• Actor on the query page (being also the landing page)</li> </ul>
Postconditions:	<ul style="list-style-type: none"> <li>• The result is displayed in the browser and can be retrieved in preferred format (pdf</li> <li>• txt)</li> <li>• the search refined</li> <li>• or modified and then retrieved.</li> </ul>
Priority:	100
Frequency of Use:	100
Basic Flow:	<ol style="list-style-type: none"> <li>1. Go to the search page</li> <li>2. Insert query in the preferred format (library index, query language, etc.)</li> <li>3. Refine search</li> </ol>
Exceptions:	
Special Requirements:	
Assumptions:	
Notes and Issues:	

Use Case ID:	1
Use Case Name:	Information Retrieval
Date Created:	26 May 2014
Date Last Updated:	26 May 2014

Actors:	System
Description:	Actor process all source listed documents and updates the database accordingly.
Preconditions:	<ul style="list-style-type: none"> <li>• Regular update in time of small server activity (e.g. 2am)</li> <li>• Update requested by administrator</li> </ul>
Postconditions:	<ul style="list-style-type: none"> <li>• The database contains the latest document information</li> <li>• If data accessed during update user notified that data may be inconsistent by a warning as well warning was saved in all generated documents</li> </ul>
Priority:	100
Frequency of Use:	10
Basic Flow:	<ol style="list-style-type: none"> <li>1. Activate at predefined time</li> <li>2. Fetch all documents</li> <li>3. Extract information</li> <li>4. Update aggregates accordingly</li> </ol>
Exceptions:	Update already in progress
Special Requirements:	Cannot be requested one after another
Assumptions:	
Notes and Issues:	

Use Case ID:	2
Use Case Name:	Save generated output to a file
Date Created:	26 May 2014
Date Last Updated:	26 May 2014

Actors:	Consumer
Description:	Actor saves generated data to a file in specified format.
Preconditions:	<ul style="list-style-type: none"> <li>• Legal format specified</li> <li>• Data generated</li> </ul>
Postconditions:	<ul style="list-style-type: none"> <li>• The download process of a generated document started</li> </ul>
Priority:	75
Frequency of Use:	50
Basic Flow:	<ol style="list-style-type: none"> <li>1. Generate the desired output</li> <li>2. Select file format</li> <li>3. Save to disk</li> </ol>
Exceptions:	
Special Requirements:	
Assumptions:	
Notes and Issues:	

Use Case ID:	3
Use Case Name:	Save generated output to a file
Date Created:	26 May 2014
Date Last Updated:	26 May 2014

Actors:	Consumer
Description:	Actor wants to register on the website.
Preconditions:	<ul style="list-style-type: none"> <li>• User name not registered</li> <li>• All required data provided</li> </ul>
Postconditions:	<ul style="list-style-type: none"> <li>• The actor account exist in the system</li> </ul>
Priority:	25
Frequency of Use:	25
Basic Flow:	<ol style="list-style-type: none"> <li>1. Go to the sign up page</li> <li>2. Fill the form</li> <li>3. Sign up</li> <li>4. Verify</li> <li>5. Log in</li> </ol>
Exceptions:	
Special Requirements:	
Assumptions:	
Notes and Issues:	



Use Case ID:	4
Use Case Name:	Update the source listing
Date Created:	26 May 2014
Date Last Updated:	26 May 2014

Actors:	Administrator
Description:	Actor wants to update/delete existing listing or insert a new one.
Preconditions:	<ul style="list-style-type: none"> <li>• Listing already not specified</li> <li>• Listing not being subdomain of existing listing</li> </ul>
Postconditions:	<ul style="list-style-type: none"> <li>• Data can be fetched from the listing source</li> </ul>
Priority:	50
Frequency of Use:	5
Basic Flow:	<ol style="list-style-type: none"> <li>1. Log in</li> <li>2. Go to the source listing page</li> <li>3. Request an action</li> <li>4. Fill the data</li> </ol>
Exceptions:	
Special Requirements:	
Assumptions:	
Notes and Issues:	

Use Case ID:	5
Use Case Name:	Request listing to be included in the source listing
Date Created:	26 May 2014
Date Last Updated:	26 May 2014

Actors:	Consumer
Description:	Actor wants that his listing would be included in the retrieval process.
Preconditions:	<ul style="list-style-type: none"> <li>• Listing already not specified</li> <li>• Listing not being subdomain of existing listing</li> </ul>
Postconditions:	<ul style="list-style-type: none"> <li>• Data can be fetched from the listing source</li> </ul>
Priority:	40
Frequency of Use:	2
Basic Flow:	<ol style="list-style-type: none"> <li>1. Go to the suggest listing page</li> <li>2. Insert link</li> <li>3. Normalize link</li> <li>4. Test if valid</li> <li>5. Notify about the result</li> <li>6. Notify administrator to take action</li> </ol>
Exceptions:	
Special Requirements:	
Assumptions:	
Notes and Issues:	

Use Case ID:	6
Use Case Name:	Report Bug
Date Created:	26 May 2014
Date Last Updated:	26 May 2014

Actors:	Consumer
Description:	Actor wants to report a bug which he or she encountered either during usage of the system or observed basing on the experience (can be constructive to improve the quality of service).
Preconditions:	<ul style="list-style-type: none"> <li>• Error occurred</li> <li>• User at suggestions page</li> </ul>
Postconditions:	<ul style="list-style-type: none"> <li>• Ticket saved in the system and present on the ticket board (visible to administrators)</li> </ul>
Priority:	60
Frequency of Use:	10
Basic Flow:	<ol style="list-style-type: none"> <li>1. Error occurred and feedback form was presented</li> <li>2. User filled the form and clicked submit</li> <li>3. Ticket shown in the system</li> </ol>
Exceptions:	
Special Requirements:	Must protect from XSS and code injections
Assumptions:	Input is normalized
Notes and Issues:	

## Chapter 3

# Requirements analysis

### 3.1 Objectives

- Continuable - the project development should be easy to continue
  - Well-documented
  - Popular and well-documented technologies should be used
  - Well-organized code base
- Schemaless - documents (sometimes even from the same source) will have different structure, so it is important to not restrict the database to store one kind (structure) of information
- Adaptability - it should be easy to adapt a project to a new kind of documents/task
- Modularity - the distinguishable project parts should be implemented in modular fashion to increase clarity and make the future development easier
- Open source - all technologies and tools used should be open source

### 3.2 Data layer

#### 3.2.1 Example query

1. Show all Committee of Experts recommendations under article 7b.
2. Show all sections relating to Scottish Gaelic.

### 3.2.2 Options

The following database technologies were taken into consideration:

- Traditional RDMS such as PostgreSQL and MySQL
- Key-value store such as Riak and Redis
- Document database such as MongoDB and CouchDB
- Graph database such as Neo4J and Infinite Graph

**Traditional RDMS** has drawback of dedicated to the outlined schema and every new document (not to mention every new source of documents) has a chance to have different organization than the previous ones.

**Key-value store** has not enough good query languages and

**Document database** has both flexible structure of aggregates (so we are not dedicated to one schema) and powerful query language (which can be implemented using incremental update)

**Graph database** has ability to perform rich queries regarding relationships, but is not as powerful when it comes to queries relating the content, can answer queries like 1 and 2 in  $O(1)$  time, but update is expensive.

### 3.2.3 Choice

MongoDB has been chosen for storing documents because

- it has big community of supporters and developers;
- it is well-documented;
- naturally blends with JavaScript when implementing query (jQuery notation) and data from the database (JSON);
- can store and operate on big documents of varying schema easily
- has rich query language with greater capabilities than the need outlined in the example queries subsection;
- when used with Node.js increases programmers productivity significantly (one language, centralized view)

When developing session and account management functionality it would be beneficial to implement the concept of polyglot persistence and for these two use key-value database, which creates separation of the concerns as well as suits better the case.

## 3.3 Application layer

### 3.3.1 Example usage

### 3.3.2 Options

Three programming languages/frameworks were considered as potentially suitable for the project, namely:

- Java and JSP
- Python and Django
- JavaScript and Node.js

**Java and JSP** does not support well the programmer productivity as to perform even the simplest operations it is very expressive. When used the application needs to be developed and tested as a whole, so it does not support modularity and extendability well. The data layer technology of choice provides API for the language, but it is rather easier and more natural to operate on it in JavaScript than in Java.

**Python and Django** had the same drawbacks, except it provides greater productivity.

**JavaScript and Node.js** means that one language would be used for front-end and back-end, including data querying and manipulation, what highly increases the programmer productivity. Further, JavaScript as front-end language provides rich choice of libraries like jQuery, Angular, etc. that can be used to further increase productivity and presentation quality. When it comes to the back-end all operations are handled efficiently using event-loop. The aggregates are retrieved in JSON format and MongoDB can be naturally queried using jQuery like notation. This all results in complete and efficient solution.

### 3.3.3 Choice

Node.js has been chosen for the front-end and back-end development because

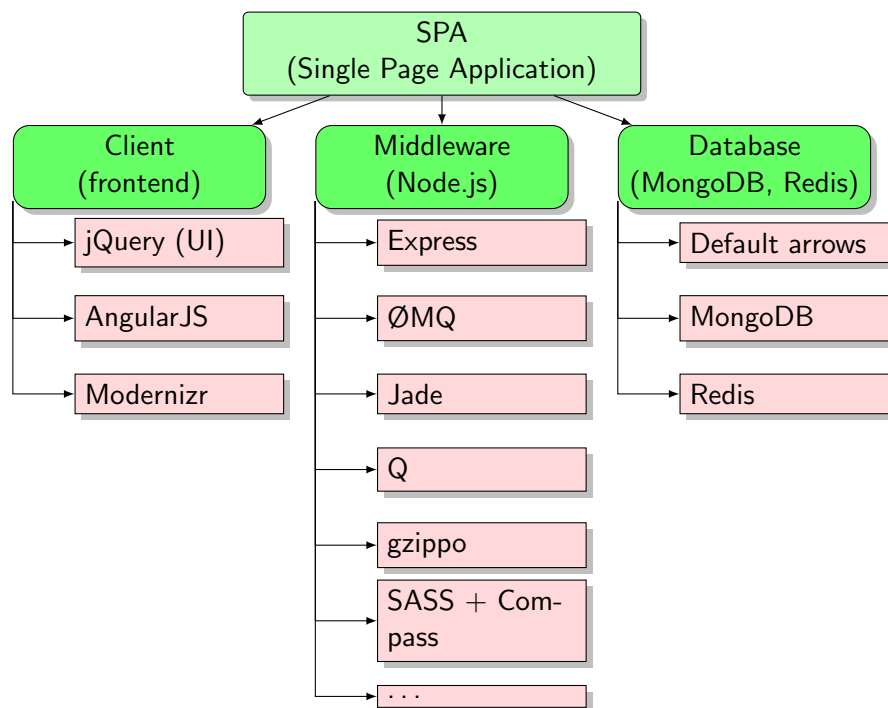
- it integrates well with MongoDB
- enables using one language for everything (front-end, back-end, database);
- is efficient

## Chapter 4

# Architecture

### 4.1 Overview

#### 4.1.1 Application



Description

#### Client:

- jQuery (UI)
- AngularJS
- Modernizr

### Middleware:

- Express - main framework for message passing and service the content to the clients
- ØMQ - for message passing, mainly for efficient Router/Dealer protocol implementation
- Jade - for more productivity in templating
- gzip - for serving gzipped content
- SASS + Compass - styling, plus features like building sprites out of files

### Database:

- MongoDB - document data store for document management *TODO*
- Redis - key-value data store for session and user data management

#### 4.1.2 Full layer

*Picture showing Yeoman, Grunt, Karma testing with Jasmin, etc.* code:  
sudo npm install -g yo # installs yeoman, grunt, bower, etc. globally  
sudo npm install -g generator-angular-fullstack # install angular generator (btw the most popular generator)  
yo angular-fullstack seeker (projectname) # create full stack locally with jade templates [?] Would you like to use Sass (with Compass)? Yes [?] Would you like to include Twitter Bootstrap? Yes [?] Would you like to use the Sass version of Twitter Bootstrap? Yes [?] Which modules would you like to include? angular-resource.js angular-cookies.js angular-sanitize.js angular-route.js mongo with mongoose passport All yeses



### 4.1.3 Code management and deployment process

## 4.2 Full stack

### 4.2.1 Client facing

### 4.2.2 Middleware

### 4.2.3 Persistence data store

## 4.3 System design

## 4.4 Query language

### 4.4.1 Elements

- *source documents* (e.g. Committee of Experts reports)
- *document specifiers* (e.g. cycle, country producing the document)
- *selectors* (i.e., languages, sections, words)
- *organizer* specifying how to structure the output document (e.g. arrange by section numbers, in each section by countries alphabetically)
- *template* choose how to style the output document

### 4.4.2 Dependencies as tree

1. *Source documents* must be root.
2. *Document specifiers* can be children of each another, but can occur only once on any path from root to leaves, can appear multiple times.
3. *Selectors* (i.e. languages, sections, words) can be children of any node, and they are always terminal (leaves).
4. *Organizers*, *templates*, and other similar options which can be added in the future must be children of root, and in case of *organizers* they must be aware of all root's children except itself to display available options.

### 4.4.3 Query evaluation stages

1. Select relevant documents using:
  - producer - group (e.g. committee)
  - producer - origin, i.e. countries (e.g. UK)
  - cycles - time (e.g. last)

2. filter the document content using selectors (i.e. sections, languages, filters- words, phrases, etc.)
3. reorganize resulting JSON using organizers
4. render using template