

# COMPETITIVENESS AND INNOVATION FRAMEWORK PROGRAMME

CIP-ICT-PSP-2013-7 Pilot Type B



WP3 – Service platform integration and deployment in  
cloud infrastructure

D3.3.2: Sematic tagging and open data publication tools

Deliverable Lead: PSNC

Deliverable due date: 29/02/2016

Actual submission date: 29/02/2016

Version: 1.2

This project is partially funded under the ICT Policy Support Programme (ICT PSP) as part of the Competitiveness and Innovation Framework Programme by the European Commission under grant agreement no. 621074



**Document Control Page**

<b>Title</b>	D3.3.2 Semantic tagging and open data publication tools
<b>Creator</b>	Marcin Krystek (PSNC)
<b>Description</b>	This document introduces the current prototype of the semantic annotation service
<b>Publisher</b>	FOODIE Consortium
<b>Contributors</b>	Raul Palma (PSNC), Emilio Rubiera (CTIC)
<b>Creation date</b>	01/02/2016
<b>Type</b>	Text
<b>Language</b>	en-GB
<b>Rights</b>	copyright "FOODIE Consortium"
<b>Audience</b>	<input type="checkbox"/> internal <input checked="" type="checkbox"/> public <input type="checkbox"/> restricted
<b>Review status</b>	<input type="checkbox"/> Draft <input type="checkbox"/> WP leader accepted <input type="checkbox"/> Technical Manager accepted <input checked="" type="checkbox"/> Coordinator accepted
<b>Action requested</b>	<input type="checkbox"/> to be revised by Partners <input type="checkbox"/> for approval by the WP leader <input type="checkbox"/> for approval by the Technical Committee <input type="checkbox"/> for approval by the Project Coordinator
<b>Requested deadline</b>	

**STATEMENT FOR OPEN DOCUMENTS**

(c) 2016 FOODIE Consortium

The *FOODIE* Consortium (<http://www.foodie-project.eu>) grants third parties the right to use and distribute all or parts of this document, provided that the *FOODIE* project and the document are properly referenced.

THIS DOCUMENT IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. EXCEPT WHAT SET FORTH BY MANDATORY PROVISIONS OF LAW IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS DOCUMENT, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

**About the project**

FOODIE project aims at creating a platform hub on the cloud where spatial and non-spatial data related to agricultural sector is available for agri-food stakeholders groups and interoperable. It will offer: an infrastructure for the building of an interacting and collaborative network; the integration of existing open datasets related to agriculture; data publication and data linking of external agriculture data sources, providing specific and high-value applications and services for the support of planning and decision-making processes.

FOODIE project is addressed to four basic groups of users: a) stakeholders from the agriculture sector as end-users of final applications, b) public sector for communication with farmers about taxation, subsidies and regulation, c) researchers for large scale experimentation on real data and d) ICT companies for the development of new applications for agriculture and food sector, mainly using implemented tools

FOODIE specifically works on three pilots:

- Pilot 1: Precision Viticulture (Spain) will focus on the appropriate management of the inherent variability of crops,
- Pilot 2: Open Data for Strategic and Tactical Planning (Czech Republic) will focus on improving future management of agricultural companies (farms) by introducing new tools and management methods,
- Pilot 3: Technology allows integration of logistics via service providers and farm management including traceability (Germany).

**Contact information**

Miguel Angel Esbri

Project Coordinator

Atos Spain, Madrid, Spain

E-mail: [miguel.esbri@atos.net](mailto:miguel.esbri@atos.net)

URL: <http://www.foodie-project.eu>

Twitter: [https://twitter.com/FOODIE\\_Project](https://twitter.com/FOODIE_Project)

## Table of Contents

<b>Glossary</b> .....	<b>5</b>
<b>Abbreviations and Acronyms</b> .....	<b>6</b>
<b>Executive Summary</b> .....	<b>7</b>
<b>1 Introduction</b> .....	<b>8</b>
<b>2 Components Description</b> .....	<b>8</b>
<b>2.1 Semantic Annotation Service</b> .....	<b>8</b>
<b>2.2 Annotation Tools</b> .....	<b>8</b>
2.2.1 AgroTagger.....	8
2.2.2 Babelfy.....	9
<b>2.3 Data Semantization Service</b> .....	<b>9</b>
<b>2.4 Data Access and Publishing Components</b> .....	<b>9</b>
2.4.1 Silk.....	9
2.4.2 Sparql.....	9
2.4.3 Faced Search and Browser.....	9
2.4.4 Semantic Annotation Service.....	10
<b>2.5 Storage for Annotation Service</b> .....	<b>10</b>
2.5.1 Semantic Store.....	10
2.5.2 RDBMS.....	10
2.5.3 File system.....	10
<b>3 System Design</b> .....	<b>10</b>
<b>4 Data Model for Annotations</b> .....	<b>11</b>
<b>5 Implementation Details</b> .....	<b>12</b>
<b>5.1 Semantic Annotation Service</b> .....	<b>12</b>
5.1.1 REST API.....	12
5.1.2 Development and frameworks.....	12
<b>5.2 Data Semantization Service</b> .....	<b>13</b>
5.2.1 Tabels Core.....	13
5.2.2 Tabels Web.....	13
<b>6 Deployment</b> .....	<b>13</b>
<b>6.1 Service Location</b> .....	<b>13</b>
<b>6.2 Semantic Annotation Demo Client</b> .....	<b>13</b>
<b>References</b> .....	<b>15</b>

## Glossary

The glossary of terms used in this deliverable can be found in the public document “FOODIE\_Glossary.pdf” available at: <http://www.foodie-project.eu>

## Abbreviations and Acronyms

Abbreviation / Acronym	Description
API	Application Programming Interface
CPU	Central Processing Unit
DBA	Database Administrator
DBaaS	Database as a Service
DDD	Domain-Driven Design
DNS	Domain Name System
GTM	Global Transaction Manager
HA	High Availability
HDD	Hard Disk Drive
HTTP	Hypertext Transfer Protocol
IaaS	Infrastructure as a Service
IDE	Integrated Development Environment
JAR	Java ARchive
MPP	Massive Parallel Processing
OS	Operating System
POM	Project Object Model
RAM	Random Access Memory
RDBMS	Relational Database Management System
RDF	Resource Description Framework
SLA	Service Level Agreement
SQL	Structured Query Language
TCP	Transmission Control Protocol
TDD	Test – Driven Development
VM	Virtual Machine
API	Application Programming Interface

*Table 1: Abbreviations and Acronyms*

## Executive Summary

This document introduces the semantic annotation service.

## 1 Introduction

Semantic tagging and open data publication tools is a system of services which can be used to extract additional knowledge from unstructured data such as plain text or text files in different formats, transform (semi-) structured data into semantic format, and publish the generated data according the Linked Data principles.

In particular, semantic annotations enrich the source data with a context that is linked to some structured knowledge of a domain or application, which can then be exploited by different applications and services. For instance, this knowledge can be used to enrich the description of the resources presented to the user by the Marketplace platform, as described in D3.5.1<sup>1</sup>.

Since the newly discovered knowledge is described by standard ontologies, stored in machine-readable format and accessible through standard APIs and protocols, it can also be used for further machine processing allowing better integration with existing knowledge bases and their publication in the Linked Open Data Cloud<sup>2</sup>, discovering and understanding relations and dependencies between resources (e.g., Marketplace products), as well as the implementation of all other kind of user scenarios.

Similarly, semantic data generated through the transformation of (semi-)structured data sources, such as tabular data or XML, can then be exploited to discover additional knowledge through the integration with other data elements via semantic links.

## 2 Components Description

Semantic tagging and open data publication tools compose together a system in which the following components can be distinguished:

### 2.1 Semantic Annotation Service

Semantic Annotation Service is a core component of the system. It provides a simple REST API designed for creating, updating and retrieving annotations. Annotation Service orchestrates other services in order to control and fully perform data analysis process, creates semantic form of the generated data and persists semantic data using semantic store.

### 2.2 Annotation Tools

These are external tools used as libraries or services that allows keywords extraction, sentence disambiguation and text meaning identification. The results generated by these tools are then used to create semantic representation of the annotation. Current implementation of the Semantic Annotation Service uses the following annotation tools:

#### 2.2.1 AgroTagger

AgroTagger<sup>3</sup> is a keyword extractor that uses the ARGOVOC thesaurus as its set of allowable keywords. It also identifies concepts from AGROVOC vocabulary in terms of Linked Open Data. AGROVOC is a controlled vocabulary covering all areas of interest of the Food and Agriculture Organization (FAO) of the United Nations, including food, nutrition, agriculture, fisheries, forestry, environment etc. It is published by FAO and edited by a community of experts.

### 2.2.2 Babelfy

Babelfy<sup>4</sup> is a unified, multilingual, graph-based approach to Entity Linking and Word Sense Disambiguation based on a loose identification of candidate meanings coupled with a densest subgraph heuristic, which selects high-coherence semantic interpretations.

## 2.3 Data Semantization Service

The data semantization service is a tool implemented to carry out chiefly conversions from a heterogeneous range of data into RDF.

Mainly focused on the transformation of tabular data in different formats, it is expected to add as well means to transform XML files into RDF, in future implementations.

It is an open source tool with a simple Web interface where users upload the file to be transformed and thanks to the mediation of a small lightweight piece of code that identifies the structure of the added data, it is capable of generating an output file fully compliant with the RDF family language.

This early version offers only a manual tool with no guidelines or best practices for the end user; however, in the next releases, it is planned that the conversion become more user-friendly and even show semi-automatic capabilities. Besides, the number of input format supported will be increased.

Thanks to this tool, users will be able to transform their unstructured or semi-structured data into more friendly formats and pave the way to share this data in an easier way or simply get that data in a format fully compliant with the linked open data guidelines.

## 2.4 Data Access and Publishing Components

Semantic annotations and other semantic data generated and stored by the semantic services are accessible using the following frameworks and API's:

### 2.4.1 Silk

Silk<sup>5</sup> is a linked data integration framework for integrating heterogeneous data sources. It allows generating links between related data items within different Linked Data sources and applying data transformations to structured data sources. Linked Data publishers can use Silk to set RDF links from their data sources to other data sources on the Web. Accordingly, in FOODIE, Silk is used to discover (and store) links between the semantic data in the semantic store and other key knowledge bases in the Linked Open Data Cloud.

### 2.4.2 Sparql

Sparql<sup>6</sup> is a standard query language for RDF data. FOODIE provides access to RDF data through a SPARQL endpoint that is exposed by the semantic store.

### 2.4.3 Faced Search and Browser

Faced Search and Browser is based on Virtuoso Faceted Browser<sup>7</sup>, which allows simple text search and results browsing. In FOODIE, we use this interface mainly to enable simple navigation of RDF data through their links.

#### 2.4.4 Semantic Annotation Service

Semantic Annotation Service provides dedicated API which allows annotation retrieval in the simplified JSON form. Currently this API is being used by the Marketplace portal<sup>8</sup>, which showcases its usage through external applications.

### 2.5 Storage for Annotation Service

The Semantic Annotation service processes and produces various types of data. In order to meet specificity of the data type as well as different requirements and constraints regarding data processing procedure, three different storage components are used.

#### 2.5.1 Semantic Store

For storing semantic data the OpenLink Virtuoso<sup>9</sup> server is used. Virtuoso has built in support for SPARQL query language and inference engine. It also provides good performance, geospatial support and capabilities to generate linked data from different data sources types via Sponger. Sponger is also transparently integrated into Virtuoso's SPARQL Query Processor delivering URI de-referencing within SPARQL query patterns and is directly accessible via REST interfaces.

#### 2.5.2 RDBMS

For storing simple relational data such as mappings, resource processing status and statistics the relational database management system is used. The current implementation of Semantic Annotation Service uses sqlite<sup>10</sup> implementation as simple, lite and providing satisfying performance solution. Since sqlite provides standard SQL interface it can be easily replaced with other, more advanced solution if required.

#### 2.5.3 File system

Semantic Annotation Service allows annotating unstructured files. All files send to the service are stored temporarily in the file system. They are available only for the annotation procedure and will be removed after procedure successfully finishes.

## 3 System Design

Semantic tagging and open data publication tools compose together a system for which three main layers can be identified:

- API layer – implemented as set of REST services. Each service provides different functionality for the client application or final user. API's denotes overall system functionality.
- Service layer – it is set of different web services or libraries directly used by Semantic Annotation Service or operating on data generated by Semantic Annotation Service
- Storage layer – provides persistence for all data processed or generated by the Semantic Annotation Service

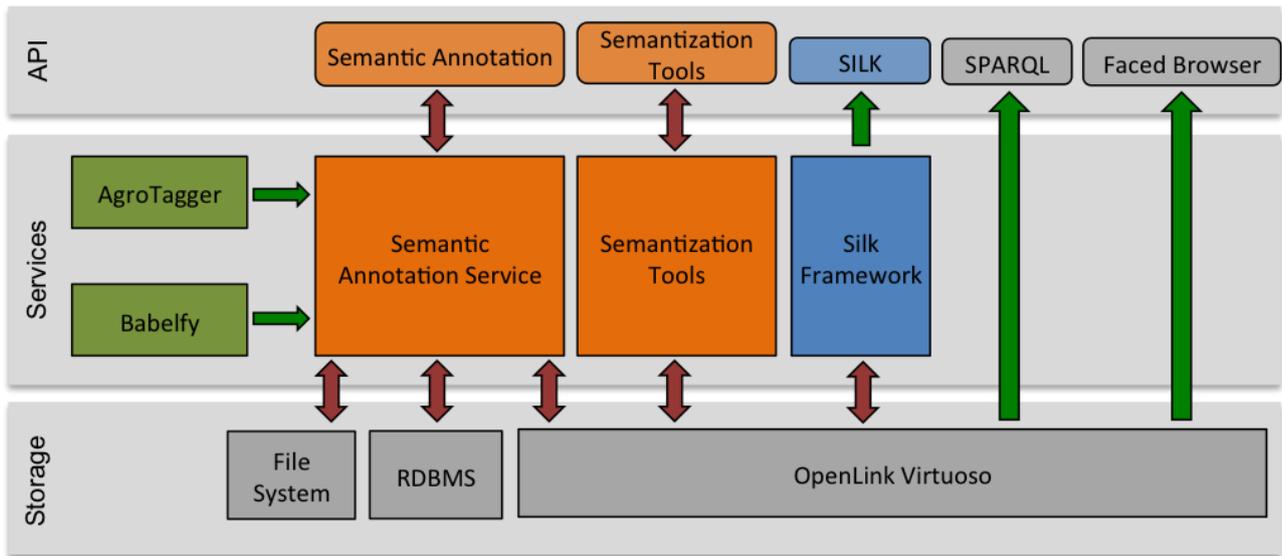


Figure 1 Sematic tagging and open data publication tools system design

Arrows on the Figure 1 denotes direction of the data flow direction.

#### 4 Data Model for Annotations

Annotations created by the Semantic Annotation Service are modeled using Modular Unified Tagging Ontology (MUTO)<sup>11</sup>. It is designed specifically for tagging and folksonomies. It allows representing public and private tagging, simple and auto generated tags and others. It is also easily extensible since all concepts defined in MUTO ontology inherits from other more general ontologies like SKOS<sup>12</sup>, SIOC<sup>13</sup> or vocabularies as RDFS<sup>14</sup>.

The general structure of the MUTO ontology and main concepts are presented in the following diagram.

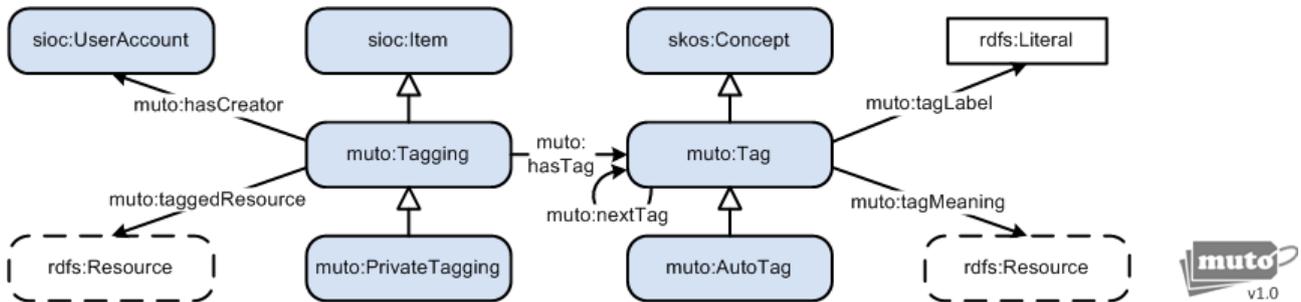


Figure 2 MUTO ontology general design

The central concept of the MUTO ontology is muto:Tagging. muto:Tagging instance is created as response for create annotation request. For each tagging a tagged resource is also created. Tagged resource instance represents an object that is actually annotated. It can be anything (file, set of files, text description) that can be addressed or identified by the client. For each tagging instance there is a set of muto:Tag instances created. Each muto:Tag instance represents a single concept identified by the annotation tools. muto:Tag instance has its default label defined as muto:tagLabel. If the identified concept has its representation in the Linked Open Data Cloud, an additional skos:Concept instance is created and related with tag using muto:tagMeaning. This object can be also enriched with labels in different languages if such are available in LOD Cloud.

Tags identified by automatic annotation tools are created as `muto:AutoTag` instances. Tags defined by user manually are stored as `muto:Tag` instances. `muto:PrivateTagging` and `sioc:UserAccount` are currently not used.

## 5 Implementation Details

### 5.1 Semantic Annotation Service

#### 5.1.1 REST API

Semantic Annotation Service REST API provides the following functionality:

Method	Description
<b>POST /tagging</b>	Creates new tagging for requested list of files and text description. Tagging will consist of set of tags (annotations) identified by annotation tools. It also creates tagged resource instance. If resource instance URI is provided in the request, the method updates existing tagging.
<b>GET /tagging</b>	Retrieves tag labels for the resource and language specified in the request.
<b>POST /tag</b>	Adds new tag to the existing tagging.
<b>GET /statistics/meaning</b>	Retrieves statistics showing which concepts are most commonly identified in the annotated resources.

More technical description of the Semantic Annotation Service REST API can be found in D3.2.2. Open API specification.

The tagging creation request may contain plain text description or set of files. Current implementation of Semantic Annotation Service supports processing of the following file types: PDF, DOC, XLS, HTML, TXT

#### 5.1.2 Development and frameworks

- **Programming language**  
Semantic Annotation Service is a Java web application. It requires at least Java 1.7 for development and for runtime environment.
- **Application framework**  
Semantic Annotation Service uses Spring v.4.1.4 as main application framework, taking advantage of its bean management and event processing functionalities as well as support for REST API implementation.
- **Building tools**  
Semantic Annotation Service is build using Apache Maven v.3.1.3. All project dependencies and deployment process are defined and managed by Maven.
- **Web container**  
Building process creates war file that can be deployed on any web container. Currently Apache Tomcat v.7.0 is used for the development and production deployment.
- **Source code location**

Source code of the application is available on git repository:  
<https://git.man.poznan.pl/stash/projects/FOOD/repos/semantic-annotation/browse/SemanticAnnotationRestApi>

- Artifactory

Dependencies that cannot be resolved using public repositories are defined in the following artifactory:  
<http://maven.man.poznan.pl/repository/simple/foodie-releases-public>

## 5.2 Data Semantization Service

The Data Semantization Service comprises two components: Tabela Core and Tabela Web.

### 5.2.1 Tabela Core

Tabela Core holds the capability to map from supported tabular formats to RDF documents; the mapping between them should be specified by a transformation program by the means of a DSL (domain specific language) developed with the only purpose of these mappings.

Tabela core was developed in Scala language and using sbt to support build definitions. The source code is available at: <https://bitbucket.org/fundacionctic/tabela>

### 5.2.2 Tabela Web

Tabela Web is a demonstrator for the interoperability capabilities of a group of technologies related to Linked Data. It integrates online tabela core, tapinos-ws, tapinos-js, undermaps and a bunch of linked data exploitation tools.

Tabela Web was developed in Grails. The source code is available at: <https://bitbucket.org/fundacionctic/tabela>

## 6 Deployment

### 6.1 Service Location

- Semantic Annotation Service is available at:  
<http://foodie-vm3.man.poznan.pl/semanticAnnotation>
- Data Semantization Service is available at:  
<http://prepro.fundacionctic.org/foodie/> (not yet deployed in FOODIE cloud)
- SPARQL endpoint is available at:  
<http://foodie-vm1.man.poznan.pl/sparql>
- Faced Browser is available at:  
<http://foodie-vm1.man.poznan.pl/fct>
- SILK endpoint is available at:  
<http://foodie-vm3.man.poznan.pl/silk>

### 6.2 Semantic Annotation Demo Client

Demo client is a simple web form. It allows user to simply add text description and select files that he wants to annotate. The form values are then sent to the Semantic Annotation Service where tagging is performed. As a

response, user will receive URI of the created resource. User will use this resource URI to retrieve or update existing tagging. Currently available at: <http://foodie-vm3.man.poznan.pl/semanticAnnotation/index.html>

## References

---

- <sup>1</sup> D3.5.1 Marketplace; Maciej Łabędzki, Patryk Promiński, Adam Rybicki
- <sup>2</sup> <http://lod-cloud.net/>
- <sup>3</sup> <https://github.com/fcproj/agrotagger>
- <sup>4</sup> <http://babelfy.org/about>
- <sup>5</sup> <http://silkframework.org/>
- <sup>6</sup> <https://www.w3.org/TR/rdf-sparql-query/>
- <sup>7</sup> <http://virtuoso.openlinksw.com/dataspace/doc/dav/wiki/Main/VirtFacetBrowserInstallConfig>
- <sup>8</sup> <https://www.foodie-cloud.org/>
- <sup>9</sup> <http://virtuoso.openlinksw.com/dataspace/doc/dav/wiki/Main/>
- <sup>10</sup> <https://www.sqlite.org/>
- <sup>11</sup> <http://muto.socialtagging.org/core/v1.html>
- <sup>12</sup> <https://www.w3.org/2004/02/skos/>
- <sup>13</sup> <http://rdfs.org/sioc/spec/>
- <sup>14</sup> <https://www.w3.org/TR/rdf-schema/>