



Project Number 732223

# **D3.4 Natural Language Components**

Version 1.0 27 December 2017 Final

**Public Distribution** 

## **Edge Hill University**

Project Partners: Athens University of Economics & Business, Bitergia, Castalia Solutions, Centrum Wiskunde & Informatica, Eclipse Foundation Europe, Edge Hill University, FrontEndART, OW2, SOFTEAM, The Open Group, University of L'Aquila, University of York, Unparallel Innovation

Every effort has been made to ensure that all statements and information contained herein are accurate, however the CROSSMINER Project Partners accept no liability for any error or omission in the same.



# **Project Partner Contact Information**

Athens University of Economics & Business	Bitergia
Diomidis Spinellis	José Manrique Lopez de la Fuente
Patision 76	Calle Navarra 5, 4D
104-34 Athens	28921 Alcorcón Madrid
Greece	Spain
Tel: +30 210 820 3370	Tel: +34 6 999 279 58
E-mail: dds@aueb.gr	E-mail: jsmanrique@bitergia.com
Castalia Solutions	Centrum Wiskunde & Informatica
Boris Baldassari	Jurgen J. Vinju
10 Rue de Penthièvre	Science Park 123
75008 Paris	1098 XG Amsterdam
France	Netherlands
Tel: +33 6 48 03 82 89	Tel: +31 20 592 4102
E-mail: boris.baldassari@castalia.solutions	E-mail: jurgen.vinju@cwi.nl
<b>Eclipse Foundation Europe</b>	Edge Hill University
Philippe Krief	Yannis Korkontzelos
Annastrasse 46	St Helens Road
64673 Zwingenberg	Ormskirk L39 4QP
Germany	United Kingdom
Tel: +33 62 101 0681	Tel: +44 1695 654393
E-mail: philippe.krief@eclipse.org	E-mail: yannis.korkontzelos@edgehill.ac.uk
FrontEndART	OW2 Consortium
Rudolf Ferenc	Cedric Thomas
Zászló u. 3 I./5	114 Boulevard Haussmann
H-6722 Szeged	75008 Paris
Hungary	France
Tel: +36 62 319 372	Tel: +33 6 45 81 62 02
E-mail: ferenc@frontendart.com	E-mail: cedric.thomas@ow2.org
SOFTEAM	The Open Group
Alessandra Bagnato	Scott Hansen
21 Avenue Victor Hugo	Rond Point Schuman 6, 5 <sup>th</sup> Floor
75016 Paris	1040 Brussels
France	Belgium
Tel: +33 1 30 12 16 60	Tel: +32 2 675 1136
E-mail: alessandra.bagnato@softeam.fr	E-mail: s.hansen@opengroup.org
University of L'Aquila	University of York
Davide Di Ruscio	Dimitris Kolovos
Piazza Vincenzo Rivera 1	Deramore Lane
67100 L'Aquila	York YO10 5GH
Italy	United Kingdom
Tel: +39 0862 433735	Tel: +44 1904 325167
E-mail: davide.diruscio@univaq.it	E-mail: dimitris.kolovos@york.ac.uk
Unparallel Innovation	
Bruno Almeida	
Rua das Lendas Algarvias, Lote 123	
8500-794 Portimão	
Portugal	
Tel: +351 282 485052	
E-mail: bruno.almeida@unparallel.pt	



# **Document Control**

Version	Status	Date
0.1	Document outline	18 September 2018
0.2	Document outline Restructure	20 September 2018
0.3	Background Section Completed	23 October 2018
0.4	Readers Section Completed	25 October 2018
0.5	Metric Providers and Tools sections completed	14 November 2018
0.6	Introduction, Risks and Limitations, Conclusion sections completed	20 November 2018
0.7	WP3 Progress section completed	26 November 2018
0.8	Internal first draft completed	6 December 2018
0.9	Internal Partner Review Version	10 December 2018
1.0	Final QA version of the deliverable	27 December 2018



# **Table of Contents**

1	Intr	oduction	1
	1.1	Overview	1
	1.2	Intentions	1
	1.3	Outcome	1
2	CRO	OSSMINER Component Overview	3
	2.1	Natural Language Components	3
		2.1.1 Readers	3
		2.1.2 Metric providers	5
		2.1.3 Tools	6
		2.1.4 Knowledge Base	6
	2.2	Summary	6
3	Rea	ders and Reader Related Components	7
	3.1	Delta Based Readers	7
		3.1.1 Bug Tracking Systems	8
		3.1.2 Communication Channels	2
	3.2	Non-Delta Based Readers	6
		3.2.1 Social Media	6
		3.2.2 Question & Answer Websites	17
	3.3	Summary	8
4	Met	ric Providers	19
	4.1	Transient Metrics	9
		4.1.1 Natural Language Processing Transient Metrics	9
		4.1.2 Bug Tracking Systems	20
		4.1.3 Newsgroups	20
		4.1.4 Forums	21
		4.1.5 Document Preparation	21
	4.2	Historic Metrics	22
		4.2.1 Bug Tracking Systems	22
		4.2.2 Newsgroups	23
		4.2.3 Forums	24
		4.2.4 Document Preparation	24
	4.3	Summary	24



5	Tool	S		25
	5.1	Text P	re-Processing Tools	25
		5.1.1	HTML Parser	26
		5.1.2	Markdown Parser	26
		5.1.3	Code Detector	26
	5.2	Natura	l Language Processing Tools	26
		5.2.1	Core Natural Language Processing Tools	27
		5.2.2	Severity Classifier	27
		5.2.3	Emotion Classifier	27
		5.2.4	Request Reply Classifier	27
		5.2.5	Content Classifier	28
		5.2.6	Sentiment Analyser	28
		5.2.7	Thematic Clustering	28
	5.3	Index	Manager	28
	5.4	Summ	ary	29
6	Risk	s and L	imitations	30
7	Con	clusions	5	31
	7.1	Work I	Package 3: Progress	31
	7.2	Future	Work	35



# **Executive Summary**

Work package 3 (WP3) aims to contribute to CROSSMINER a collection of text mining components that allows users and developers of Open source Software to analyse a diverse range of natural language sources: bug tracking systems, communication channels, social media and question and answer websites. The components developed in this work package will contribute to the enrichment of knowledge contained within the GitHub enabling open source users and developers to gain an insight into the tools they use and improve the performance and quality of the code they produce.

WP3 is constructed of 4 tasks that produce 5 deliverables. In Task 3.1 we have analysed text-representation methods and developed new state-of-the-art machine learners for processing natural language. In Task 3.2 we explored methods for reading, searching and storing natural language sources associated with open source projects.

This deliverable is the product built upon the foundations of work presented in D3.1, D3.2 and D3.3. This document focuses on the work associated with Task 3.3 - the development of natural language components. It presents details regarding all of the natural language (NL) components, readers, metric providers and tools that are integrated (newly developed, updated or migrated) into CROSSMINER, discussing in detail their role, importance and inter-relationship with other NL components. The document also provides information concerning risks and limitations of the components developed. The report concludes by presenting the overall progress of WP3, work completed, requirements that have been fulfilled and outstanding work.

# List of Abbreviations

Term	Definition
API	Application Programmable Interface
BTS	Bug Tracking System
CC	Communication Channel
PDE	Plug-in Development Environment
JSON	JavaScript Object Notation
NL	Natural Language
NLP	Natural Language Processing
NNTP	Network News Transfer Protocol
OSS	Open Source Software
REST	Representational State Transfer



# **1** Introduction

Content sources related to open source software, such as GitHub, NNTP news groups and StackOverflow, provide developers with a diverse range of useful information relating to specific open source projects. These sources are rich with natural language. However, key knowledge relating to a project is often overlooked as it is hidden away in vast amounts of text. The key to unlocking this hidden knowledge contained within text is natural language processing (NLP).

With respect to CROSSMINER, this deliverable follows on from work completed in D3.2 and D3.3. It focuses specifically on the development and integration of components related to the natural language processing capabilities of the CROSSMINER platform. Developed components, such as readers, historic and transient metrics and text processing tools are discussed in detail. Consequently, work completed as part of this deliverable lays the foundations for Task 3.4, which is currently in progress and aims at developing a recommender that is capable of suggesting code snippets and documentation relevant to the code being developed.

## 1.1 Overview

The remaining of this deliverable is organised into six sections. Section 2 provides an overview of CROSS-MINER components, focusing on natural language. Sections 3, 4 and 5 discuss the readers, metric providers and tools developed and integrated into CROSSMINER, respectively. Section 6 identifies the risks and limitations of the work presented in this deliverable. Finally, Section 7 presents the conclusion of this deliverable as well as reports on the current progress of WP3 and future & outstanding work.

## **1.2 Intentions**

The aim of Task 3.3 is to develop and integrate natural language related components into CROSSMINER. Fulfilment of the aim requires addressing the following use case and technological requirements presented in Tables 1 and 2, respectively. Each requirement has an ID, a description and priority<sup>1</sup>.

## 1.3 Outcome

The outcomes of this deliverable are the components related to natural language that have been implemented and integrated into CROSSMINER. The developed natural language components, such as readers, text processing tools and metrics, are designed to be compatible with the CROSSMINER platform, the knowledge base and can also be used to run bespoke workflows for knowledge extraction.

<sup>&</sup>lt;sup>1</sup>Shall - must be fulfilled, Should - should be fulfilled, May - may be fulfilled



Ref	Description	Priority
U31	Able to search mailing-lists	Shall
U32	Able to search forums	Shall
U33	Able to search issues	Shall
U37	Able to search IRC	Shall
U38	Able to search simultaneously across all the data sources	Shall
U39	Able to present the results by data source or globally for the project	Shall
U40	Able to order the results by relevance	May
U41	Able to search using regular expression patterns	Shall
U42	Able to detect in the data sources text referring to one or several bugs	Shall
U43	Able to detect in the data sources text referring to one or several commits	Shall
U44	Able to extract sentiment analysis from the communication data sources	Shall
U46	Able to detect similarities in messages so as to suggest answers	Should
U47	Able to propose recommendations to improve community management	Should
U48	Able to identify recurring problems identified by users	Should
U52	Able to list unattended bugs	Should
U53	Able to list unattended messages	Should
U54	Able to create a cluster map for messages using Lingo or similar clustering technology	Should
U55	Able to create a cluster map for bugs contents using Lingo or similar clustering technology	Should
U60	Able to analyse level of user activity in issue trackers	Shall
U61	Able to identify if a thread reached a conclusion and if it was positive	Should

### Table 1: WP3 Use Case Requirements related to Task 3.3

Ref	Description	Priority
D32	The NLP analysis component should record the sources used by the developer and use them	Should
	as a form of automatic feedback to improve its suggestions.	
D72	The other mining tools developed in WPs 2-4 shall expose a REST API that the cross-project	Shall
	relationship miners can consume	
D73	The other mining tools developed in WPs 2-4 shall expose project level metrics that cross-	Shall
	project relationship miners can consume	

### Table 2: WP3 Technology Requirements related to Task 3.3



# 2 CROSSMINER Component Overview

The CROSSMINER platform, as shown in Figure 1, consists of numerous components that have been or will be developed by CROSSMINER research and use case partners. When combined, the components provide software developers with the capabilities to monitor, perform in-depth analysis and evidence-based selection of open source software, and facilitate knowledge extraction from large open source software repositories. The platform is capable of extracting knowledge associated with:

- Source Code
- System Configurations
- Cross-Project Relationships
- Communication in Natural Language

Combined, these sources provide CROSSMINER with a diverse range of knowledge that is related to various aspects of the software development life-cycle, allowing developers to make informed decisions based upon the results they are presented. The remaining of this document focuses specifically on components related to natural language, which are marked with an orange border in Figure 1.

### 2.1 Natural Language Components

Natural language contains vital and potentially hidden information that can be exploited to assist developers in making vital decisions surrounding open source software development. Consider the following scenario. A developer is developing a new parser for XML and would like to compare two libraries, SAX and JAXB. The developer is relatively inexperienced with this type of work and would benefit greatly by using a library that has an active and helpful community. Both of these libraries are hosted on GitHub as open source projects and are associated numerous entries on StackOverflow. To provide an insight into these two projects and present the developer with relevant heuristics requires a series of interconnected components designed specifically to perform tasks associated with reading textual information, processing it, computing suitable metrics and enriching the original text with extra knowledge.

The objectives of natural language components developed for CROSSMINER are three-fold. Firstly, they should be compatible with the CROSSMINER platform, so that they can be used to analyse the history of given Open Source software projects to compute metrics that summarise the quality of support offered to users over time. Secondly, to contribute to the CROSSMINER knowledge base, by enriching documents with extra information that facilitates efficient faceted search. Thirdly, to be used as plug-ad-play parts of bespoke knowledge extraction workflows, in the context of Work Package 5. CROSSMINER NLP components are of three categories: readers, metric providers, and tools. A detailed description of each component type and its role in CROSSMINER is presented in the remaining of this section.

#### 2.1.1 Readers

A reader is the first component in a typical natural language processing workflow. The aim of a natural language reader in CROSSMINER is to retrieve textual information related to an open source software project, along with metadata associated with it. CROSSMINER supports four types of sources of textual information: Bug Tracking Systems, Communication Channels, Social Media and Question and Answer websites.





Figure 1: CROSSMINER Component Diagram

Readers have been developed for a variety of different sources, meaning that data retrieved is heterogeneous. A vital function within each reader is to restructure the data in such a way that it can be consumed by various other components within the natural language workflow in CROSSMINER. Depending on the source type, information retrieved by the readers can be passed either to metric providers or to natural language tools. Both paths result in enrichment of the knowledge base. A detailed description of the readers in this deliverable is found in Section 3.

### 2.1.2 Metric providers

Natural Language metric providers compute a diverse range of heuristics based upon the origin of the textual information<sup>2</sup>. In CROSSMINER, metrics are computed for every day that a project is being evaluated. For every day of the project, readers produce a delta, i.e. the set of changes that occurred in a particular communication means on that day. Changes may among others refer to new messages, updated bug reports, comments and thread status changes. During computation, a metric may rely upon information contained within the delta, results from other metrics and historical information from previous dates or days.

There are two types of metric providers in CROSSMINER. For a given project, transient metrics, as the name suggests, compute metrics for a particular day, overwriting the outcome that was computed previously, for the previous day of the project. Historic metrics are persistent, keeping a record of the changes for each day the project is evaluated. From a functional perspective, metric providers in CROSSMINER are organised into two tier stack as shown in Figure 2. In general, historic metrics enrich the knowledge base with knowledge relating to a project's history, and have dependencies on transient metrics in order to do this. This dependency relationship between the two types of metric providers is vital during the execution of the platform. A transient metric is only triggered during execution of the analysis components in CROSSMINER, if it is associated as a dependency with a historic metric. Without this dependency transient metrics will not be triggered during runtime. As a result of this architectural requirement, there are cases where historical metrics have been integrated into CROSSMINER with the sole purpose to trigger specific transient metrics. A detailed description of the metric providers associated with the processing of natural language and enrichment of the knowledge base in CROSSMINER is available in Section 4.



Figure 2: Metric Provider Stack

<sup>&</sup>lt;sup>2</sup>There are some exclusions to this definition, however this is the general rule.



#### 2.1.3 Tools

Natural language tools provide various capabilities that relate to either processing or indexing of knowledge in the CROSSMINER Knowledge Base. Once again depending on the origin of source data, tools can be triggered either by metric providers or by readers directly. An in-depth presentation of the tools developed and integrated in CROSSMINER is given in Section 5.

#### 2.1.4 Knowledge Base

The knowledge base represents the final stage of the natural language processing workflow in CROSSMINER. The Knowledge Base consists of components that store knowledge that has been previously computed by metric providers and specific tools. For simplicity, the knowledge base consists of an instance of a MongoDB database and an ElasticSearch cluster. Knowledge contained within the knowledge base is presented to the users of the CROSSMINER platform, assisting them in the development of open source software. Although, the knowledge base and its development is not part of Work Package 3, the indexing capabilities for CROSS-MINER are related to this work, as stated among the work package requirements, in Section 1.2.

### 2.2 Summary

This section provides an overview of the CROSSMINER platform. Since Work Package 3 is responsible for natural language processing in CROSSMINER, the latter part of this section discussed the role of various types of natural language components developed, the relationship of these components to the rest of the CROSSMINER platform and the importance of natural language processing for CROSSMINER. The following sections provide further details about each natural language component.



## **3** Readers and Reader Related Components

As discussed in D3.2, OSSMETER utilised a delta-based architecture for its readers and retrieved information from two distinct types of sources: *bug tracking systems* and *communication channels*. As part of Task 3.3, we introduced an additional, *non-delta based* reader architecture. This enables the CROSSMINER platform to support a further two types of sources, i.e. *social media* and *questions & answer websites*, for mining natural language, as shown in Figure 3. Furthermore, eight new readers have been developed, three were updated and a further six were migrated from OSSMETER. Our efforts now allow CROSSMINER to ingest natural language information from a total of 17 sources<sup>3</sup>.



Figure 3: Readers in CROSSMINER

## 3.1 Delta Based Readers

All delta-based readers in CROSSMINER follow the same workflow illustrated in Figure 4. Information about an open source software project is retrieved from bug tracking systems or communication channels. The distinct functionality of delta-based readers is the ability to compute metrics based upon deltas. A delta is the sum of all activity that occurred on a particular day on a communication means, forming a time-line of historic activity as shown in Figure 5. Deltas are used to compute various metrics that enrich the knowledge contained in the knowledge base.

Regardless of the source type, each reader produces deltas in the same way. It begins by calculating the first date of activity for a particular project on a particular communication means. The calculation varies slightly depending on the source type. For example, bug tracking systems use the date of the first issue, whereas communication channels may use the first post in a forum or message posted in a newsgroup. The first date provides the CROSSMINER platform with the point in time, which can be incremented sequentially to produce deltas until the current day of execution.

Deltas contain lists of objects that represent events. For example, a delta can contain a new issue or update to an existing issue or comment on a bug tracking system. The information that is contained within the delta is then utilised by numerous metric providers to compute various heuristics associated to the source type, which in turn enriches the knowledge base. The remaining of this subsection provides details about particular delta-based readers in CROSSMINER and identifies those that are new, those that were updated and those that were migrated from OSSMETER.

<sup>&</sup>lt;sup>3</sup>details of each of the 17 sources supported by CROSSMINER can be found throughout the remainder of this section





Figure 4: Delta Based Reader Workflow



Figure 5: Delta

#### 3.1.1 Bug Tracking Systems

Bug tracking systems (BTS), also known as issue trackers, assist developers in managing "bugs" throughout the lifetime of a software project. They are also used as channels that allow users of a software product to express ideas for feature enhancements or ask for technical support. This type of information is processed and utilised by metrics to measure characteristics such as how active development is, the number of bugs and the response time. The metrics enrich the knowledge base to provide suggestions such as solutions to issues faced by OSS developers.

In CROSSMINER, a BTS reader mines data related to issues and comments posted by users, as a minimum. Some BTS readers may also include additional elements unique to that particular system. These will be discussed in further detail where relevant.

As part of Task 3.3, we introduced support for two new bug tracking systems, updated three existing ones and migrated three from OSSMETER. The CROSSMINER platform now supports nine bug tracking systems. A summary of each bug tracking system supported by CROSSMINER is presented in Table 3.



Source	Status	Package name	
Bitbucket	Updated	org.eclipse.scava.platform.bugtrackingsystem.bitbucket	
Bugzilla	Updated	org.eclipse.scava.platform.bugtrackingsystem.bugzilla	
GitHub	Updated	org.eclipse.scava.platform.bugtrackingsystem.github	
GitLab	New	org.eclipse.scava.platform.bugtrackingsystem.gitlab	
JIRA	Migrated	org.eclipse.scava.platform.bugtrackingsystem.jira	
Mantis BT	New	org.eclipse.scava.platform.bugtrackingsystem.mantis	
Redmine	Migrated	org.eclipse.scava.platform.bugtrackingsystem.redmine	
Sourceforge	Migrated	org.eclipse.scava.platform.bugtrackingsystem.sourceforge	

Table 3: Bug Tracking System Readers

#### 3.1.1.1 New Bug Tracking System Readers

As mentioned earlier, CROSSMINER introduced support for two popular bug tracking systems: *GitLab* and *Mantis Bug Tracker*. With regards to the use case partners, both these new readers are to be utilised during the evaluation of several use cases.

**GitLab** is a web based application for all aspects of DevOps lifecycle and includes a bug monitoring (tracking) system. We have designed a delta-based reader in collaboration with YORK, who have developed *RESTMule*, an automatic client generator for sources with Open API specifications. Currently *RESTMule* is not ready to support *GitLab*, thus the final implementation of this reader will be completed shortly. The reader will mine information related to issues and comments. The GitLab reader will also support two forms of authentication: OAuth2 tokens<sup>4</sup> and Personal Access Tokens<sup>5</sup>. It is expected that the GitLab reader will require access to the project parameters shown in Table 4. It should be noted that they may change in future versions of the GitLab API.

Parameter	Туре	Description	Example
OAuth2 Token	String	OAuth2 Token used for authen-	6Ji3tcrZfNQ3gna7CiE8JSLSiO6FiOIi
		tication	
Personal Access Token	String	Personal Access Token used for	6Ji3tcrZfNQ3gna7CiE8JSLSiO6FiOIi
		authentication	
Project	String	Name of the GitLab Project	Crossminer

Table 4: GitLab: Project Mining Parameters

**Mantis Bug Tracker** is a web-based, open source bug tracking system developed in PHP. It has been made available in several Linux distributions, Mac OSX and Windows. The reader integrated into CROSS-MINER was developed using OKHTTP and Jackson<sup>6</sup>. As with many other bug tracking system readers, this reader can request and process information in issues and comments posted by users that relate to a given open source project. The reader accesses this information via authenticated requests made to the Mantis Bug Tracker REST API. Authentication requires generating an API token, which users can do via "Create API Token" located under "My account section" on the website of Mantis. The level of access a user has to information is subjective to the privileges they were assigned by the system administrator. Please note that in order for closed issues to also be included and processed by CROSSMINER a configuration setting located in

<sup>&</sup>lt;sup>4</sup>GitLab OAuth2 Tokens - docs.gitlab.com/ee/api/#oauth2-tokens

<sup>&</sup>lt;sup>5</sup>GitLab Personal Access Tokens - docs.gitlab.com/ee/api/#personal-access-tokens

<sup>&</sup>lt;sup>6</sup>More information regarding these libraries can be found in D3.2



<path\_to\_install\_directory>\config\config\_inc.php had to be modified. The following line had to be added to
the file or replace any previously assigned value: \$g\_hide\_status\_default = META\_FILTER\_NONE;

The project information shown in table 5 has to be provided to the CROSSMINER platform to mine information related to a specific project contained within a Mantis Bug Tracker instance.

Parameter	Туре	Description	Example
project_id	String	the numerical identifier for the	1
		project	
host	String	URL that points to the Mantis	http://localhost:8888/mantis/
		Bug Tracker instance	
token	String	A generated string consisting of	6Ji3tcrZfNQ3gna7CiE8JSLSiO6FiOIi
		32 characters	

Table 5: Mantis Bug Tracker: Project Mining Parameters

#### 3.1.1.2 Updated Bug Tracker Readers

As stated in D3.2, some OSSMETER readers required updating to address issues related to outdated or depreciated APIs. Three readers, *Bitbucket*, *Bugzilla* and *GitHub*, received significant modifications to their source code from their OSSMETER counterpart to maintain their mining functionality and compatibility with CROSSMINER.

**Bitbucket:** In deliverable 3.2, we stated that the Bitbucket reader would be migrated since our initial testing showed that it was operational. However, since then, Bitbucket has been going through a transitional phase with regards to its REST API. In OSSMETER, this reader utilised both version 1 and version 2 of the Bitbucket REST API. It appears that Bitbucket is discontinuing support for version 1 of API, meaning that eventually the end points of version 1 will be no longer available. Since the reader's functionality depended on certain version 1 endpoints that were no longer available, the reader would crash as data it was expecting was no longer being returned. To address these problems, we updated all the methods that made calls to version 1 endpoints to their version 2 equivalent and made modifications to the serialisation and deserialisation functionality to support the new information that was introduced.

For the reader to make requests to the Bitbucket REST API, a user must be authenticated. The Bitbucket reader in CROSSMINER supports basic authentication (username and password). The following project information shown in Table 6 has to be provided to CROSSMINER to mine information about a specific project contained in a Bitbucket repository.

Parameter	Туре	Description	Example
User	String	The name of the user	Bob
Repository	String	The name of the repository	CROSSMINER
Login	String	The user name used to log into	Cross
		Bitbucket	
Password	String	The password used to log into	M1ner
		Bitbucket	

Table 6: Bitbucket: Project Parameters

**Bugzilla:** in D3.2, we mentioned that the reader developed for Bugzilla in OSSMETER used a depreciated version of the Bugzilla API, meaning that some functionality of the reader would no longer work and therefore required updating. Since Bugzilla 5.0, the Mozilla Foundation has switched to using a REST API. The Bugzilla REST API also adheres to the Open API specification<sup>7</sup>. This enables the use of *RESTMule* developed by our partners from the University of York, to provide the foundations for the updated Bugzilla reader in CROSS-MINER. However, currently RESTMule is not ready to support Bugzilla, thus the final implementation of this reader will be completed shortly. In the meantime the original Bugzilla reader developed in OSSMETER has been made available in CROSSMINER.

The Bugzilla reader in CROSSMINER is capable of mining information related to issues, comments and attachments. This will remain the same in the RESTMule variant to maintain the consistency across both implementations. The reader currently supports two authentication methods: basic authentication (username or password) or token authentication.

Presented in Table 7 are the project parameters related to Bugzilla. Note that at least one authentication method has to be provided in order to make authenticated requests to the repository.

Parameter	Туре	Description	Example
Username	String	The user name used to log into	Cross
		Bugzilla	
Password	String	The password used to log into	M1ner
		Bugzilla	
Product	String	Refers to a product contained	Crossminer
		within Bugzilla	
Component	String	Refers to a specific component	Bugzilla Reader
		within a project	

Table 7: Bugzilla: Project Parameters

**GitHub:** Similar to the Bugzilla reader, the API used in the OSSMETER version of the GitHub reader used eGit API version 4.1, which has undergone several improvements since its use in OSSMETER. Furthermore, GitHub has also revised its REST API (currently version 3) with several improvements and also made it adhere to the Open API specification. Hence why the resilient client generator was used to provide the foundations of the GitHub reader in CROSSMINER. The GitHub reader in CROSSMINER mines information related to issues, pull requests and comments posted by users; this was to maintain consistency with the data mined by the reader implemented in OSSMETER. The GitHub reader supports two types of authentication: basic authentication (username or password) or token authentication.

Although GitHub enforces a restriction on the number of calls that can be made per hour, the capabilities of the resilient client generator ensures that this is managed correctly and will not negatively impact CROSSMINER. The GitHub reader is associated with the following project parameters presented in Table 8.

<sup>&</sup>lt;sup>7</sup>More information regarding Open API specification and the benefits it provides can be found in D3.2



Parameter	Туре	Description	Example
User	String	The name of the user	Bob
Login	String	Username used to access	Cross
		GitHub (basic authentication)	
Password	String	Password used to access GitHub	Miner
		(basic authentication)	
Token	String	A unique token used to make	6ea3a349d921f6d33dc745df44d2487db3c
		authenticated requests token au-	
		thentication	
Repository	String	The name of the repository	crossminer
Owner	String	the username of the user that	Bob2020
		owns the repository	

Table 8: GitHub: Project Parameters

#### 3.1.1.3 Refactored Bug Tracker Readers

The readers listed below received minor refactoring and were tested to ensure compatibility with CROSS-MINER during their migration from OSSMETER. The modifications did not concern their original functionality. Each of the readers mines information about issues and comments. The project parameters required by CROSSMINER for each source are presented in Table 9.

Source	Parameter	Туре	Description	Example
	Login	String	User-name used to access Jira	Cross
Jira	Password	String	Password used to access Jira	M1ner
	Project	String	The name of the project	crossminer
Redmine	Login	String	User-name used to access Jira	Cross
	Password	String	Password used to access Jira	M1ner
	Name	String	Name of the user	Bob
	Project	String	The name of the project	crossminer
Sourceforge	URL	String	A url that represents the bugs of	http://sourceforge.net/rest/p/ soapui/bugs
			a project	

 Table 9: Migrated Bug Tracking Systems Readers: Project Parameters

#### **3.1.2** Communication Channels

Within the context of CROSSMINER, a communication channel refers to an online community where various discussions related to OSS take place. Communication Channels are used to provide information for calculating metrics that measure various factors relating to a community of users, such as how active and reactive a community is, or estimate the general sentiment about specific discussions, or provide historic information concerning bugs and their respective solutions.

In OSSMETER, communication channels existed in two forms. Network News Transfer Protocol (NNTP) articles and SourceForge discussions. Due to the requirements from our use case partners, we extended the capabilities of CROSSMINER to also support forums as a new type of communication channel. From the platform's perspective, forums are handled differently due to their "threaded" nature and the information they

contain. This required additional functionality to be added to the Communication Channel Delta<sup>8</sup> and the addition of new metrics, which will be discussed in further details in section 4.

As shown in Table 10, CROSSMINER is capable of ingesting information from six Communication Channels; three readers are newly developed, two were refactored and one was migrated from OSSMETER.

Source	Status	Package name
Eclipse Forums	New	org.eclipse.scava.platform.communicationchannel.eclipseforums
IRC	New	org.eclipse.scava.platform.communicationchannel.irc
NNTP	Refactored	org.eclipse.scava.platform.communicationchannel.nntp
Sourceforge	Refactored	org.eclipse.scava.platform.communicationchannel.sourceforge
Sympa Emails	New	org.eclipse.scava.platform.communicationchannel.sympa
Zendesk	Migrated	org.eclipse.scava.platform.communicationchannel.zendesk

Table 10: Communication Channel Readers

#### 3.1.2.1 New Communication Channel Readers

The following readers for communication channels were developed for CROSSMINER.

**Eclipse Forums:** Eclipse Forums are a centralised repository for people who develop and use Eclipse-based tools to discuss technical matters. They are organised in the following structure:

- *Category* is a collection of related forums.
- *Forum* is a collection of related topics
- *Topic* is a collection of related posts.
- *Post* is a message associated with a specific topic.

The reader integrated into CROSSMINER was developed using OKHTTP and JACKSON libraries. Given a specific forum, the Eclipse Forums reader is capable of extracting all of the posts associated with topics that are related to the forum. Accessing this information requires requests to be made to the Eclipse Foundation REST API. Eclipse enforces a rate limit of 1,000 requests per hour for both authenticated<sup>9</sup> and anonymous requests.

The reader, from an authentication perspective, supports both anonymous and authenticated requests. For authenticated requests, Eclipse requires a OAuth2 token. This can be generated following the instructions provided by the Eclipse Foundation<sup>10</sup>. The reader in its current form does not request information from protected resources, therefore its functionality is not impacted regardless of the availability of a token. However, authenticated users may request to increase the number of calls they are able to submit per hour, subject to Eclipse approval, allowing more data to be consumed by CROSSMINER per hour.

The project parameters associated with the Eclipse Forums reader are presented in Table 11.

<sup>&</sup>lt;sup>8</sup>A component responsible for managing Communication Channel Deltas in CROSSMINER

<sup>&</sup>lt;sup>9</sup>Note the rate limit can be extended, by contacting the appropriate eclipse representative. See https://api.eclipse.org/docs for more information.

<sup>&</sup>lt;sup>10</sup>Eclipse Forums OAuth2 token generation - https://api.eclipse.org/docs/auth



Parameter	Туре	Description	Example
Client Id	String	Client Id (optional)	G12DrqtW6qPuXC88Oorqc4rt4Xw5
Client Secret	String	Client secret (optional)	wpl9NBtLn2KxA4234mPxWoDbb5I
Forum Id	String	The unique id of the eclipse forum	305

Table 11:	Eclipse	Forums:	Project	Parameters
-----------	---------	---------	---------	------------

**Internet Relay Chat (IRC)** is an application for user communication in textual format. The chat is based on a server and client model. Each user uses an IRC client, either locally installed or web-based. All client applications are connected to a server. A user message is sent through the client application to the server and in turn to other user clients. IRC is mainly designed for group communication organised in channels, but also supports user to user private messaging. Communication messages can be logged by the intervening server. In the log, each message is timed and marked with the sender's username.

IRC used to be very popular in the 1990s and early 2000s. Nowadays, most users have moved to modern social media. However, IRC is still used for communication in software communities and organisations. Analysing IRC communication messages about OSS can inform about the quality of support offered to users. As analysing IRC logs is of interest to our use case partners, we developed a delta based reader of archived IRC logs. The reader assumes that an archive of IRC logs stored locally and does not use specialised libraries.

Parameter	Туре	Description	Example
Log Archive	String	The Log archive to be processed	Crossminer_irc.log
Channel	String	The specific IRC channel which will be the focus	Readers

Table 12: IRC: Project Parameters

**Sympa:** is a mailing list management software. Similarly to most email servers, Sympa stores all emails sent by users in archives. The archive contains a folder for each year of activity, which contains a folder for each month, which in turn contains a folder for each day of that month. These day folders contain a text file for every email sent on that day. Emails are stored following the Multipurpose Internet Mail Extensions (MIME) format.

Parsing and analysing Sympa email archives is of interest to our use case partners. Thus, we developed a reader for Sympa archives, that can read all emails in an archive, or just the emails of a particular day. In succession, emails can be processed using any of the natural language processing components of the CROSS-MINER platform to compute quality metrics or searched by the users via the Eclipse IDE plugin, as part of the knowledge base. The reader uses *javax-mail, apache commons mail* and *apache commons net*.

Parameter	Туре	Description	Example
Log Archive	String	The Log archive to be processed	Crossminer_emails.log

Table 13: Sympa: Project Parameters

#### 3.1.2.2 Refactored Communication Channel Readers

As mentioned earlier, three readers were migrated from OSSMETER. The first reader is capable of accessing newsgroups that utilise the Network News Transfer protocol (NNTP), the second processes information from Sourceforge discussions and the third one from ZenDesk. Each of these readers was subjected to testing after refactoring and was not modified to it's source code for use within CROSSMINER. Presented in the table below are the project parameters for each of the readers.

Source	Parameter	Туре	Description	Example
	Authentication Re-	String	Represents a flag (this is	true
NNTP	quired		dependent on the NNTP	
	T T		server)	
	Username	String	User-name used to access	Cross
	D 1	<u></u>	the news group	
	Password	String	Password used to access the	MIner
	Devet	Lat	news group	1224
	Port	Int	Port number used to access	1234
	Description	Chains	Description of the news	This is a serie series for the
	Description	String	Description of the news-	This is a news group for the
	NT	Ct. in a	group	CROSSMINER project
	Iname	String	Name of the user	Bob
		Int	Polling interval in seconds "	5
	URL	String	A url which points to the	nntp://host.domain/newsgroup
	NT	Ct.	NNTP server	CDOSCIMINED
0	News group name	String	Name of the news group	
Sourceforge	URL	String	A Url that points to discus-	http://sourceforge.net/rest/p/
	A dissident Di	Ct.	sions hosted on Sourceforge	gimponosx/discussion/
Zendesk	Authentication Re-	String	Represents a flag (this is de-	true
	quired		f and	
	I.I.e.	Chains	liguration)	Crass
	Username	String	Zandash	Cross
	Decouverd	String	Zelluesk	Minor
	Password	Sumg	Zandash	WITTER
	Dout	Int	Dort number used to occore	1224
	Pon	Int	Zendesk	1234
	Description	String	Description of the project	This is the CROSSMINER
	Description	Sung	bested in Zendesk	project
	Nama	String	Name of the user	Project Bob
	Interval	Int	Polling interval in seconds <sup>12</sup>	5
		String	A LIPL which points to the	J
		Sung	Zandesk server	http://2endesk.com/Crossiiiiiei
	News group name	String	Name of the news group in	CROSSMINER
	incws group name	Sumg	Zandask	CROSSIVIIINER
			Zenuesk	

Table 14: Migrated Communication Channel Systems Readers: Project Parameters

<sup>&</sup>lt;sup>12</sup>Polling interval is the period of time between the end of the time-out period of the last retry or completion of a network request, and the next request for data on the network.



## 3.2 Non-Delta Based Readers

Non-delta based readers aim to enrich the knowledge base with further contextual information regarding a specific open source project. Due to the large variety of data available on "Social Media" and "Question & Answer websites" it is difficult to process information in daily deltas. Non-delta based readers utilise queries to focus on information relevant to the requirements of the end user. The information is then processed by the workflow shown in Figure 6. Indexed information is searchable by the developers at the time of development through a specialised Eclipse plugin.

A potential issue that non-delta based readers face is that the size of data returned by the query can be large and may take a considerable amount of time to process. To alleviate this delay, CROSSMINER users can use computational resources available on other nodes of their network. For this purpose we utilised Crossflow, a tool that enables parallel and distributed execution of workflows, where nodes within the workflow can be tailored to meet specific requirements. Crossflow was developed by the University of York. We have worked closely with our partners to design a suitable workflow capable of enriching the ElasticSearch cluster contained within the knowledge base and to ensure that our natural language components are compatible with the requirements of Crossflow.



Figure 6: Non-Delta Based Workflow

#### 3.2.1 Social Media

From a software perspective, social media can provide an insight into how a community perceives a project. This perception is an important factor that developers may wish to take into consideration when selecting which projects to use. As part of this deliverable we integrated support for Hacker News into CROSSMINER.

**Hacker News**<sup>13</sup> is a social media website mainly focusing on sharing news related to technology and related discussions. In Hacker News, users can also post questions, job offers, polls, among others; these can be discussed too through comments. As any social media website, it allows users to upvote and downvote publications. Access to posts and comments is free to all users, either subscribed or not.

Hacker News has two REST APIs: one developed by the owners of Hacker News and one provided by the search engine Algolia. The former<sup>14</sup> offers access to the data posted in Hacker News almost on real time by using Firebase<sup>15</sup>, Google's mobile platform. This API allows to access all information available in Hacker News by publication ID, where publications are either posts or comments. Without knowing the ID of a publication, only a limited number of publications is accessible, such as the most recent posts or the most

<sup>&</sup>lt;sup>13</sup>news.ycombinator.com

<sup>&</sup>lt;sup>14</sup>github.com/HackerNews/API

<sup>&</sup>lt;sup>15</sup>firebase.google.com



upvoted ones. The API does not require authentication and, to the best of our knowledge, it does not have any query quota.

The second API<sup>16</sup>, developed by the search engine Algolia, is similar to the one designed by the owners of Hacker News, with extra functionality that allows to search publications. As mentioned previously, to access a publication in Hacker News using the first API, it is necessary to know its ID. This is not necessary on Algolia's API. As it indexes Hacker News on real time, natural language queries are possible. This API to Hacker News is back-ended on the original API, therefore it allows using all its functionality in addition to natural language queries. Algolia's API does not require authentication. However, there is a quota of 10,000 queries per hour per IP.

We have developed a reader for Hacker News using OKHTTP and Jackson libraries. The REST API used is the one provided by Algolia, as it is more advanced that the original one and offers extra features. The reader has three modes. The first mode allows users to access the Hacker News API simply and easily, through a *Builder Pattern*. Users can submit queries, retrieve specific elements and filter them. The second mode allows to retrieve all elements on Hacker News in sequential order. The third mode allows to download the publications that match a specific query and their comments. The last two models are useful in populating CROSSMINER's knowledge database.

#### 3.2.2 Question & Answer Websites

Unlike social media sources, Question & Answer (Q&A) websites consist of questions and answers derived from a community of users. Users of these websites seek and/or disseminate knowledge to a community to help others. For CROSSMINER this information can be exploited to help developers start new projects faster, by seeking answers to questions about the use of particular libraries. As part of Task 3.3, we integrated support for StackOverflow into CROSSMINER.

**StackOverflow**<sup>17</sup> is a crowd-sourced community, governed by rules similar to those of social media, where users can post questions regarding software development and also reply to the questions of other users. Users can comment, upvote and downvote questions and answers. StackOverflow covers a wide range of programming languages and libraries used by programmers.

StackOverflow offers a REST API to access the information posted on its website. As we expect CROSS-MINER users to submit many queries to StackOverflow, we have decided to create a reader that will use dumps provided by StackOverflow's parent company, StackExchange<sup>18</sup>. In this way, CROSSMINER can after index StackOverflow locally and offer real time access to questions and answers without any restrictions, such as query quotas.

The reader has been developed using a *Builder Pattern*, meaning that CROSSMINER users can define which StackOverflow elements to index from the dump. For example, users can decide to index questions, answers, accepted answers, posts posted between certain dates, or posts regarding a specific tags<sup>19</sup>, only. As Stack-Overflow dumps are very large XML files, we have decided to back-end the reader with Java's SAX parser. Specially designed for large files, SAX efficiently accesses a file line by line, instead of loading the entire file in memory.

<sup>&</sup>lt;sup>16</sup>hn.algolia.com/api

<sup>&</sup>lt;sup>17</sup>stackoverflow.com

<sup>&</sup>lt;sup>18</sup>Dumps can be found in: archive.org/details/stackexchange. A new dump is released every three months, approximately.

<sup>&</sup>lt;sup>19</sup>Tags are textual elements that denote a topic on StackOverflow, e.g. *Java*, *C*++ and *parsing*.



It should be indicated, that the new reader can be used on other dumps provided by StackExchange, as it is using a standardised format. Using this tool, CROSSMINER users can read dumps from any StackExchange website, such as  $SuperUser^{20}$  and  $ServerFault^{21}$ .

### 3.3 Summary

This section presented a comprehensive discussion about all readers that were developed and migrated as part of Task 3.3 to access natural language information for use within CROSSMINER. We have developed 8 new readers in total, updated 3 readers and migrated 6 from OSSMETER. Additionally, we extended the functionality of the Communication Channel delta to support forums. With the introduction of the non-delta based reader architecture, the CROSSMINER platform now has the ability to enrich the knowledge base with further contextual information relating to open source projects, in addition to the previously existing delta-based enrichment. We also introduced support for two new sources of information in the form of social media and question & answer websites. Overall CROSSMINER is capable of mining textual information from a total of 17 different sources.

<sup>&</sup>lt;sup>20</sup>superuser.com

<sup>&</sup>lt;sup>21</sup>serverfault.com



## 4 Metric Providers

As discussed in Section 2, metrics compute heuristics that enrich the knowledge base with useful information, that OSS developers can use to make informed decisions throughout the stages of software development. The purpose of this section is to present all metrics integrated into CROSSMINER as part of Task 3.3. The remaining of this section is organised into 2 subsections each discussing transient metrics and historic metrics, respectively. Each metric is marked with a status label from the following:

- New: A new development for CROSSMINER
- Updated: Metric received modifications to include preprocessing of text.
- Upgraded: Metric received upgrades as a consequence of developments since D3.3
- **Migrated:** Migrated from OSSMETER

### 4.1 Transient Metrics

As their name suggests, transient metrics are used to calculate heuristics that are associated with a particular period in time, i.e. a day in the case of CROSSMINER. They are stored temporarily within the knowledge base and their output is passed as parameters in the calculation of other transient and historic metrics. Transient metrics, with regards to this deliverable, are related to five areas: Bug Tracking Systems, Newsgroups, Forums, Natural Language Processing Tools and Document Preparation.

#### 4.1.1 Natural Language Processing Transient Metrics

The transient metrics presented in Table 15 are associated with various natural language tools integrated into CROSSMINER. These metrics process the input text and return either a classification value or a conversion of text, which is in turn useful for the transient metrics that compute heuristics from Bug Tracking Systems, Newsgroups and Forums.

Metric	Status	Description
Code detector	New	This metric determines which parts of a bug comment or a newsgroup
		article are code and natural language
Plain Text Processing	New	This metric pre-processes each bug comment, article or forum post into
		a split plain text format
Request Reply Classifica-	Upgraded	This metric computes if each bug comment, newsgroup article or forum
tion		post is a request of a reply. It is based on the corresponding classifica-
		tion tool described in D3.3
Sentiment Classification	Upgraded	This metric computes the sentiment of each bug comment, newsgroup
		article or forum post. Currently it uses an upgraded version presented
		in D3.3
Emotion Classification	Upgraded	This metric indicates the emotions expressed in each bug comment,
		newsgroup article or forum post. We have replaced the tool that was
		used in OSSMETER with the one presented in D3.3
Severity Classification	Updated	This metric computes the severity level of each bug or thread.
Topic Classification	Updated	This metric computes a summary of topics dicussed in a collection of
		documents and the number of documents related to each topic

#### Table 15: CROSSMINER Transient Metrics related to Natural Language Processing



#### 4.1.2 Bug Tracking Systems

Metric	Status	Description
Active Users	Migrated	This metric keeps track of the users that submitted bug comments in the
		last 15 days
Bug Metadata	Migrated	This metric holds various metadata of a bug header, i.e. priority, status,
		operation system and resolution
Number of Comments	Migrated	The number of bug comments over time
Content Classes	Migrated	Identifies content classes in bug comments
Daily Requests & Replies	Migrated	This metric stores the number of comments, requests and replies for
		each day of the week.
Emotions	Migrated	Identifies emotions expressed in bug comments
Hourly Requests & Replies	Migrated	This metric stores the number of articles, requests and replies for each
		hour of the day.
Number of New Bugs	Migrated	Tracks the number of new bugs over time
Number Request & Replies	Migrated	Computes for each bug whether it is answered. If yes, it also computes
		the response time

Presented in Table 16 are the transient metrics available in CROSSMINER.

Table 16: CROSSMINER Transient Metrics related to Bug Tracking Systems

#### 4.1.3 Newsgroups

The transient metrics shown in Table 17 are the heuristics supported by CROSSMINER for Newsgroups.

Metric	Status	Description
Active users	Migrated	This metric keeps track of the users that submitted news in the last 15
		days.
Articles	Migrated	Tracks the number of articles
Content Classes	Migrated	Identifies content classes in Newsgroup Articles
Daily Requests & Replies	Migrated	This metric stores the number of articles, requests and replies for each
		day of the week
Emotions	Migrated	Identifies the emotions expressed in Newsgroup articles
Hourly Requests & Replies	Migrated	This metric stores the number of articles, requests and replies for each
		hour of the day
Sentiment	Migrated	Computes sentiment at the beginning of each thread, at its end, and on
		average
Threads	Migrated	This metric holds information for assigning newsgroup articles to
		threads. The threading algorithm is executed from scratch every time.
Thread Request & Replies	Migrated	The metric identifies for each thread whether it is answered. If yes, it
		also computes the response time

### Table 17: CROSSMINER Transient Metrics related to Newsgroups



#### 4.1.4 Forums

Due to the unique nature of Communication Channel sources we decided to handle and compute transient metrics associated with forums differently than Newsgroups. In particular, forums posts are already threaded and there is no guarantee that they directly relate to a bug. The integration of these components has been delayed due to problems discovered with the data returned by Eclipse Forums not matching what was displayed on-line. Therefore, we have chosen not to integrate the metrics shown in Table 18 before testing them properly to ensure accuracy. The metrics are subject to change depending on the outcome of testing.

Metric	Status	Description
Active Users	New	This metric keeps track of the users that submitted news comments in
		the last 15 days
Number of posts	New	The number of posts over time
Content Classes	New	Identifies content classes in posts
Daily Requests & Replies	New	This metric stores the number of comments, requests and replies for
		each day of the week
Emotions	New	Identifies the emotions expressed in forum posts
Hourly Requests & Replies	New	This metric stores the number of forum posts, requests and replies for
		each hour of the day
Number of New Posts	New	Tracks the number of new forum posts over time
Number of New Forum Top-	New	Tracks the number of new forum topics
ics		
Number Request & Replies	New	Indicates for each forum post whether it is answered. If yes, it also
		computes the response time

#### Table 18: CROSSMINER Transient Metrics related to Forums

#### 4.1.5 Document Preparation

The transient metrics presented in Table 15 are associated with the new Indexing capabilities of CROSS-MINER. These metrics prepare documents for indexing from bug tracking systems and communication channels. Each document contains information from delta-based readers and other metrics and is represented in JSON format that is compatible with indexing.

Metric	Status	Description
Bug Tracker Indexer	New	This metric prepares documents from bug tracking systems for indexing
Communication Channel In-	New	This metric prepares documents from communication channels for in-
dexer		dexing

Table 19: CROSSMINER Transient Metrics related to Document Preparation



### 4.2 Historic Metrics

As discussed in Section 2, historic metrics keep track of various heuristics associated with a specific open source project over its lifetime. In CROSSMINER, historic metrics are organised into four areas: Bug Tracking System, Newsgroup, Forum and Document Preparation.

#### 4.2.1 Bug Tracking Systems

Table 20 presents the historic metrics related to bug tracking systems in CROSSMINER:

Metric	Status	Description				
Number of Bugs	Migrated	It indicates how many bugs there are across all the bug tracking systems				
		related to a project over time.				
Number of Comments	Migrated	It indicates the evolution in the quantity of comments for all the bugs				
		related to a project over time				
Emotions	Migrated	Summarises the emotions expressed in the bug tracking systems of a				
		given project				
Number of New Bugs	Migrated	Indicates the number of new bugs created over time				
Number of new users	Migrated	Indicates the number of new users over time				
Open Time	Migrated	Measures the average duration an issue is open				
Number of patches	Migrated	Computes the number of bug patches over time				
Number of Requests &	Migrated	Indicates the number of requests and replies out of all comments in the				
Replies		bug tracking systems for a given project				
Response Time	Migrated	Measures the average time it takes for a request to be replied, for a given				
		project.				
Sentiment	Migrated	Reveals which sentiments are expressed in a bug tracking system for a				
		project				
Severity	Migrated	Estimates the severity levels of bugs found in the bug tracking system				
		of a project				
Status	Migrated	Indicates the status of a bug, e.g. open, closed				
Topics	Migrated	I Indicates the topics that were discussed throughout the project via clus-				
		tering				
Unanswered	Migrated	Indicates how many issues remain unanswered over time				
Users	Migrated	Holds the users that have posted at the bug tracking systems of a project				
		over time				

Table 20: CROSSMINER Historic Metrics related to Bug Tracking Systems

### 4.2.2 Newsgroups

Presented in Table 21 are the Historic metrics associated with Newsgroups in CROSSMINER.

Metric	Status	Description
Number of Articles	Migrated	Computes the number of articles per day for each newsgroup separately
Emotions	Migrated	Computes the number of emotional dimensions in comments submitted
		every day
New Threads	Migrated	Computes the number of new threads per day for each newsgroup sep-
		arately
New Users	Migrated	Computes the number of new users per day for each newsgroup sepa-
		rately
Requests & Replies	Migrated	Computes the number of request and reply newsgroup articles per day
		for each newsgroup separately
Average Requests & Replies	Migrated	Computes the average number of request and reply newsgroup articles
		per day
Response Time	Migrated	Computes the average time in which the community responds to open
		threads per day for each newsgroup separately.
		(Format: dd:HH:mm:ss:SS, where dd=days, HH:hours, mm=minutes,
		ss:seconds, SS=milliseconds)
Sentiment	Migrated	Computes the overall sentiment per repository up to the processing date.
		The overall sentiment score ranges from -1 (negative sentiment) to +1
		(positive sentiment). In the computation, the sentiment score of each
		thread contributes equally, independently of its size.
Severity	Migrated	Computes the number of newsgroup severity levels in threads submitted
		every day
Severity Response Time	Migrated	Computes the average sentiment, the sentiment at the beginning of
		threads and the sentiment at the end of threads per severity level, in
		newsgroup threads submitted every day
Severity Sentiment Time	Migrated	Computes the average sentiment, the sentiment at the beginning of
		threads and the sentiment at the end of threads per severity level, in
		newsgroup threads submitted every day
Threads	Migrated	Computes the number of threads per day for each newsgroup separately
Topics	Migrated	Computes the labels of topics (thematic clusters) in articles submitted
		every day
Unanswered Threads	Migrated	Computes the number of unanswered threads per day for each news-
		group separately
Users	Migrated	Computes the number of active and inactive users per day for each
		newsgroup separately

Table 21: CROSSMINER Historic Metrics related to Newsgroups



#### 4.2.3 Forums

The historic metrics shown in Table 22 were all newly developed for CROSSMINER to account for the new type of communication channel that was implemented. For reasons discussed earlier in Section 4.1.4, we have chosen not to integrate these at this time until we can ensure their accuracy.

Metric	Status	Description				
Number of Posts	New	Indicates how many forum posts there are within a specific forum over				
		time				
Number of Topics	New	Indicates how many forum topics there are within a specific forum over				
		time				
Emotions	New	Summarises the emotions expressed in the bug tracking systems of a				
		given project				
Number of New Posts	New	Indicates the number of new posts created over time				
Number of new users	New	Indicates the number of new users over time				
Number of Requests &	New	Indicates the number of requests and replies out of all forum posts				
Replies		within a specific forum				
Response Time	New	For a given project, it measures on average the time it takes for a reply				
		to a forum post.				
Sentiment	New	Reveals which sentiments are expressed in posts related to a specific				
		forum				
Status	New	Indicates the status of a topic, e.g. open, closed				
Unanswered	New	Indicates how many topics remain unanswered over time				
Users	New	Holds the users that have posted at the bug tracking systems of a project				
		over time				

#### Table 22: CROSSMINER Historic Metrics related to Forums

#### 4.2.4 Document Preparation

Unlike other historic metrics, the metric shown in Table 23 does not contribute any knowledge from delta based sources to the knowledge base<sup>22</sup>. This historic metric was included to enable the transient metrics responsible for document preparation to be triggered during execution, as discussed in Section 2.1.2.

Metric	Status	Description
Index	New	This metric is an architectural requirement associated with the prepara-
		tion of documents for indexing. This does not 'compute' knowledge for
		the knowledge base.

Table 23: A CROSSMINER Historic Metric related to Indexing

### 4.3 Summary

In summary, this section has described all transient and historic metrics associated with natural language processing in CROSSMINER. We have introduced a total of 25 new metric providers, updated a further two, upgraded three and migrated a further 49 from OSSMETER into CROSSMINER.

 $<sup>^{22}</sup>$ It should be noted that this historic metric should be hidden from users of the platform, so that they cannot select to exclude it



# 5 Tools

A number of tools related to natural language have been integrated into CROSSMINER. Most tools were designed and developed from scratch, adding new capabilities to the platform.

Components presented in this deliverable depend on various tool-kits<sup>23</sup> Traditionally these dependencies would be added using a dependency management mechanism such as Maven or Gradle. However, due to the limitations<sup>24</sup> of Eclipse headerless PDE (Plug-in Development Environment), libraries required by each tool-kit had to be added manually as a plug-in package that exposes the libraries at runtime. This ensures that the libraries are accessible by the tools that require them. Additionally, each tool kit was carefully scrutinised to ensure that its licence is compatible with CROSSMINER.

CROSSMINER natural language tools can be divided into three types depending on the functionality they provide to the platform, as shown in Figure 7. In this section, we describe the types of tools and the tools themselves, detailing their role in CROSSMINER, as well as their dependencies.



Figure 7: Types of Tools in CROSSMINER

## 5.1 Text Pre-Processing Tools

Due to the diverse nature of natural language sources used in CROSSMINER, data needs to be pre-processed before natural language processing. In other words, data is transformed into a standard format in a pre-processing stage. For example, GitHub and StackOverflow text comes formatted in markdown syntax. The former always returns text in the original markdown syntax, whereas StackOverflow returns text in HTML format.

To ensure that the input of natural language processing tools is formatted as required, we have implemented and integrated into CROSSMINER a series of pre-processing tools, accessible to CROSSMINER tools and users. The pre-processing tools are:

<sup>&</sup>lt;sup>23</sup>In the context of this document, a tool-kit represents a collection of code libraries (jars).

<sup>&</sup>lt;sup>24</sup>More information regarding this limitation can be found at http://maven.apache.org/plugins/maven-eclipse-plugin/pde.html



- An HTML Parser
- A Markdown Parser
- A Code Detector

#### 5.1.1 HTML Parser

An HTML Parser extracts clean text from HTML formatted input. In other words, it removes tags, such as  $\langle b \rangle$ ,  $\langle font \rangle$  and  $\langle a \rangle$ , and returns only the text, that would be shown if the input was opened in a web browser.

The HTML parser uses an interface to JSOUP HTML parser and has been extended with extra methods, such as splitting text into paragraphs<sup>25</sup>, parsing and removing certain HTML tags, or filtering of text with respect to specific labels, only.

#### 5.1.2 Markdown Parser

Markdown provides syntax for formatting plain text in websites. Markdown is widely used in websites, such as GitHub and StackOverflow, as it can be easily parsed and converted into HTML. We integrated a markdown parser able to convert the enriched text, either into plain text or HTML formatted text. The tool is back-ended by the Atlassian Commonmark parser.

#### 5.1.3 Code Detector

In software communities, messages usually contain code listings often referred to as snippets. From a human perspective, this is helpful for developers for reasons such as solving a problem or declaring a bug. Nonetheless, from a natural language processing perspective, code-related elements can be a source of noise.

Although some data sources, such as StackOverflow or GitHub, support special tags to mark code portions, others, such as newsgroups, do not provide such features. Moreover, marking code portions, in most cases, is left to the user decision, therefore the tag is sometimes used inaccurately.

For these reasons, we have developed a code detector and presented it in D3.1. It is identifies if a body of text contains code or natural language in English. The code detector is currently used by the sentiment analyser, the emotion classifier and the severity classifier.

### 5.2 Natural Language Processing Tools

The following tools are utilised for processing natural language in CROSSMINER. For each tool, we present a brief description of its role within CROSSMINER and its input and output. All tools have been designed and developed so that they can be used as standalone tools and also are compatible with both delta-based and non delta-based workflows. More details about the research aspects associated with these tools are presented in D3.3.

<sup>&</sup>lt;sup>25</sup>Paragraphs are originally defined by the HTML tag  $\langle p \rangle$ .



#### 5.2.1 Core Natural Language Processing Tools

In preparation for complex natural language processing tools, such as a sentiment analyser or a severity classifier, it is necessary to apply a series of core natural language processing tools, i.e. tools that perform fundamental text processing tasks. The tools listed below have been developed and integrated in CROSSMINER.

- A text tokeniser
- A text lemmatiser
- A text normaliser
- A Part-of-Speech tagger

The text normaliser was developed in house from scratch, whereas all other tools are back-ended by NLP4J<sup>26</sup>. Depending on the needs of the text processing task that remain to be investigated, more core natural language processing tools may be added in the future.

#### 5.2.2 Severity Classifier

A software bug can impact development in many ways. In bug tracking systems, such as Bugzilla, bugs are often assigned a severity label based on the impact that the bug may have for the project. However, filling in this label is not mandatory for posting a bug and it is often omitted. Within the context of CROSSMINER, the severity of a bug is a vital piece of information that can be used to influence decisions made by open source software developers. The severity classifier is a tool that can predict the severity of a bug found on a communication means. We have designed a severity classifier that considers seven levels of severity: *normal*, *enhancement*, *major*, *minor*, *critical* or *blocker*. The tools uses six *one-vs-rest* Support Vector Machines (SVM) organised in a tree structure.

#### 5.2.3 Emotion Classifier

An emotion classifier is a tool that determines the emotions that are expressed in an input text. Similar to the sentiment analyser, it is a tool that helps to quantify the opinion of users for a software project. Unlike OSSMETER, which used a lexicon-based emotion classifier, in CROSSMINER we have trained a neural network that allows estimating the emotions expressed in a text, e.g. *love*, *joy*, *anger*, *sadness*, *fear* and/or *surprise*.

#### 5.2.4 Request Reply Classifier

Unlike Q&A websites, where, by definition, only one question is allowed per post<sup>27</sup>, in other sources of data explored in CROSSMINER, it is possible to have multiple questions/requests combined with answers/replies. Detection of these requests and replies can improve the management of the communication means of a project. Developers can know whether posts in sources, such as forums or issue trackers, contain multiple request and whether these have been replied. In addition, external developers, potential adopters of the project, can determine how well the community is managed, and whether their requests would be replied. We have developed and integrated into CROSSMINER a request/reply classifier based on a linear Support Vector Machine (SVM).

<sup>&</sup>lt;sup>26</sup>emorynlp.github.io/nlp4j/

<sup>&</sup>lt;sup>27</sup>In some cases, users can ask questions to the author of another question or an answer, in the form of comments. However, these questions usually aim at clarifying the initial one.



#### 5.2.5 Content Classifier

Threaded messages contain indicators that can help developers evaluate an OSS project. Messages can be of different content/functionality types. For example, a message can propose a solution to a previously reported problem, or it may reinforce a previous question that was not answered, or verify that a previously proposed solution actually worked and addressed a problem. As part of OSSMETER, a detailed hierarchy of content types was constructed. In CROSSMINER, we used this hierarchy to develop a content classifier, i.e. an automated tool that assigns one or more class labels to a message based on its content. Applying the content classifier to communication channel articles and bug tracker comments provides with useful statistics about the kind of interaction taking place. The content classifier that was developed and integrated in CROSSMINER utilises a multi-class SVM.

#### 5.2.6 Sentiment Analyser

A sentiment analyser is a tool that determines whether a text expresses positive or negative sentiment, or is neutral, i.e. does not express any sentiment. In the domain of software engineering and software development, sentiment analysis is a useful tool that can provide a general estimation of the feelings of users about a software project. In other words, by applying a sentiment analyser on text from a source, such as Bugzilla or Eclipse Forums, we can determine whether users are satisfied or not through the polarity of sentiment expressed in messages. The sentiment analyser developed for and integrated into CROSSMINER is a multi-class classifier based on VastText, a neural network developed in house.

#### 5.2.7 Thematic Clustering

Due to the large quantity of information available along all the sources, it is sometimes necessary to know, *grosso modo*, which are the topics/subjects discussed. By knowing the topics discussed by users, developers can improve the management of a project. In other words, thematic clustering allows identifying trends of errors, discussions or thoughts. Therefore, in CROSSMINER, we provide a thematic clustering tool. The tool uses *Carrot*<sup>2</sup>, an open source clustering tool that groups documents into labelled clusters according to their topic. Cluster labels are determined automatically by the *Lingo* algorithm, which internally scores terms to find the candidate that bests represents each group of cluster.

### 5.3 Index Manager

The Index Manager is a new component developed for CROSSMINER. It provides the facilities to manage, interact and populate the ElasticSearch cluster contained within the knowledge base. It was developed following the singleton design pattern and using the tool-kit exposed by the *org.eclipse.scava.elasticsearch.dependencies* package. This approach restricts the instantiation of the *IndexManager* class to a single object that is responsible for coordinating all interactions between the CROSSMINER platform and the ElasticSearch cluster, as shown in Figure 8.

Page 28





Figure 8: Index Manager

The Index Manager is designed to interact with an ElasticSearch cluster running version 6.3. The toolkit integrated into CROSSMINER contains all the necessary libraries associated with version 6.3 of the "High-level Rest Client API" and "Java API" (transport client).

Currently, the Index Manager publicly exposes two static methods called *performIndexing* for indexing a document. These method automatically handles the steps for indexing a document including the creation of the index, adding mappings, converting objects to JSON and creating or updating elements within the Elastic-Search cluster. The methods require arguments such as: *indexName*, *mapping* (*optional*), *documentType*, *uid*, and *document*. Each are described in Table 24.

Parameter	Туре	Description			
Index Name	String	The name of the index you wish to add the document to			
Mapping	String	(optional)A string in JSON format that contains informa-			
		tion information regarding the data type of the fields con-			
		tained within a document.			
Document	object	A Java object that contains the information you wish to			
		index			
Document Type	String	The type of the document, e.g. 'github_issue'			
Unique ID	String	A uniquely identifiable string for the document			

Table 24: *performIndexing* Arguments

The Index Manager is also capable of supporting various configurations of ElasticSearch depending on the user's specific requirements. Users of the CROSSMINER platform are required to modify/populate the *elasticsearch.properties* file located in the *elasticsearch.properties* directory to provide the tool with the correct configuration settings for their ElasticSearch cluster contained within the knowledge base. This is a vital step to ensure that the Index Manager is capable of interacting with the user's ElasticSearch cluster.

In relation to the other natural language components developed as part Task 3.3., the Index Manager is employed by the delta-based workflow, as shown in Figure 4, and also by two transient metrics responsible for document preparation. It is also used as a component in the non-delta based workflow, as shown in Figure 6, for enriching the knowledge base with information from "social media" and "question & answer websites". The latter plays a critical role in our future D3.5.

## 5.4 Summary

This section has presented a collection of tools integrated into CROSSMINER that are utilised by numerous metric providers to facilitate access, pre-processing, processing and indexing of natural language sources.



# 6 **Risks and Limitations**

The purpose of this section is to discuss the risks and limitations related to this deliverable. Due to the nature of software development and the design choices that were made, the components discussed in this deliverable have limitations and are associated with risks, as shown below:

#### **Risks:**

- The APIs and methods used to access information may become outdated / deprecated.
- The number of API calls allowed for sources may be reduced by the API policy makers.
- Discontinuation of support for libraries used for natural language processing tools.

#### Limitations:

- For readers that utilise dumps/archives, manual intervention is required when new ones are made available.
- Large volumes of archived data from non-delta based readers may take a long time to process.
- The size of classification models of our tools may increase the overall system requirements for CROSS-MINER.
- Information contained within dumps/archives is not always up-to-date.



# 7 Conclusions

To conclude, this deliverable has discussed in detail the natural language components integrated into CROSS-MINER as a consequence of the requirements provided by our use case partners. Presented below is a summary of the outcomes of this deliverable with respect to each of the natural language components discussed throughout this report:

- Readers:
  - Two methods for enriching the knowledge base: via delta based and non-delta based readers
  - CROSSMINER now has the capabilities to mine information from Bug tracking systems, Communication channels, Social Media and Question & Answer Websites.
  - Integrated a total of 17 readers
- Metric Providers:
  - Developed various new metric providers to support forums, a new type of communication channel
  - Upgraded several natural language processing metric providers to utilise: text pre-processing tools and/or the core set of NLP tools
  - Migrated / integrated metric providers associated with: Bug tracking systems, Newsgroups and Forums
  - Included transient and historic metrics that prepare documents for indexing
- Tools:
  - Developed and integrated a core set of natural language processing tools.
  - Integrated several new and improved tools, developed in Task 3.1.
  - Developed and integrated an Index Manager, i.e. a component that adds support for populating the ElasticSearch cluster contained within the CROSSMINER knowledge base.

## 7.1 Work Package 3: Progress

The objective of this section is to provide an insight into the overall progress of Work Package 3 (WP3). WP3 has a total of four tasks and five deliverables. To date three tasks have been completed and four deliverables have been submitted. Each completed deliverable is summarised below.

- **D3.1** This deliverable reported the current progress of Task 3.1 and includes an investigation of stateof-the-art methods for text representation as well as presented experimentation with English vs Code classifier. This deliverable provided the part of the foundations for D3.3 and was delivered on time in M12.
- D3.2 The focus of this deliverable was to investigate methods for reading and searching text sources associated with OSS projects. It identified various new sources of information, highlighted issues with existing readers and provided insight to which readers could be migrated from OSSMETER to CROSS-MINER. It also explored how two new source types can be exploited in CROSSMINER. Finally, it presented an analysis of two open source search engines with regards to their suitability for CROSS-MINER. Its purpose was to provide the foundations for deliverables D3.3 and D3.4 and was delivered on time in M18 of the project.
- **D3.3** Since task one spanned over 24 months, this deliverable represents the completion of Task 3.1 and presents the experimentation and results in relation to natural language tools developed for CROSS-MINER. This deliverable was delivered on time in M24 of the project.

• **D3.4** This deliverable is associated with design, development and integration of the various natural language related components into the CROSSMINER platform. This deliverable was also delivered on time in month 24 of the project.

With regards to the current overall progress of WP3, there has been significant advancements made towards the completion of the requirements that its is linked to. Requirements associated to this work package are divided into two areas: those provided by CROSSMINER use case partners and those provided by technology partners. The progress for each area of requirements is presented in Table 25 and 26, respectively.

Ref	Description	Priority	Progress	Notes			
3. Nati	3. Natural Language Searching						
U30	Able to search "part of speech" parts of	Should	Planning	Task 3.4. Involves			
	code			the analysis of com-			
				ments			
U31	Able to search mailing-lists	Shall	Completed	Task 3.3			
U32	Able to search forums	Shall	Completed	Task 3.3			
U33	Able to search issues	Shall	Completed	Task 3.3			
U34	Able to search configuration files	Shall	In Progress	Task 3.4. Involves			
				the analysis of com-			
				ments			
U34	Able to search documentation	Shall	Planning	Task 3.4			
U36	Able to search a blog	Shall	Planning	Task 3.4. Involves			
				the analysis of snip-			
1127		01 11		pets			
U37	Able to search IRC	Shall	Completed	Task 3.3			
038	Able to search simultaneously across	Shall	Completed	Task 3.3			
1120	all the data sources	C1. 11	Constant 1	T. 1.2.2			
039	Able to present the results by data	Shall	Completed	Task 5.5			
1140	Able to order the results by relevence	May	In Drogrado	To be investigated			
040	Able to order the results by relevance	May	III Progress	further in Tests 2.4			
11/1	Able to search using regular expression	Shall	Completed	Turtulet III Task 5.4			
041	natterns	Shan	Completed	148K 5.5			
1142	Able to detect in the data sources text	Shall	In Progress	Just the implementa-			
042	referring to one or several bugs	Shan	III I TOgress	tion stage to be com-			
	referring to one of several bugs			nleted			
U43	Able to detect in the data sources text	Shall	In Progress	In collaboration with			
	referring to one or several commits	~~~~~		CWI. Just the imple-			
				mentation stage to be			
				completed.			
U44	Able to extract sentiment analysis from	Shall	Completed	Task 3.3			
	the communication data sources		_				
U45	Able to extract sentiment analysis from	Should	In Progress	Task 3.4. The Sen-			
	wikis			timent Analysis com-			
				ponent is completed.			
				A reader for wikis to			
				be implemented.			
U46	Able to detect similarities in messages	Should	Completed	Task 3.3			
	so as to suggest answers						
U47	Able to propose recommendations to	Should	Completed	Task 3.3			
	improve community management						



Ref	Description	Priority	Progress	Notes
U48	Able to identify recurring problems	Should	Completed	Task 3.3
	identified by users		1	
U49	Able to detect in the data sources one	Shall	In Progress	Task 3.4. In collabo-
	or several commits hashes		6	ration with CWI. Im-
U50	Able to list commits with bugs	Shall	In Progress	plementation stage to
U51	Able to list bugs with commits	Shall	In Progress	be completed.
U52	Able to list unattended bugs	Should	Completed	Task 3.3
U53	Able to list unattended messages	Should	Completed	Task 3.3
U54	Able to create a cluster map for mes-	Should	Completed	Task 3.3
	sages using Lingo or similar clustering			
	technology			
U55	Able to create a cluster map for bugs	Should	Completed	Task 3.3
	contents using Lingo or similar cluster-		<b>I</b>	
	ing technology			
U56	Able to create a cluster map for com-	Should	In Progress	
	mit messages contents using Lingo or		6	Task 3.4. In collabo-
	similar clustering technology			ration with CWI
U57	Able to identify the list of changed	Should	In Progress	
	third-party API methods from GitHub		U	
	issues			
U58	Able to identify that the topic of a fo-	Should	In Progress	
	rum thread is the migration of a given		C	
	third-party API			
U59	Able to identify code snippets that use	Shall	In Progress	
	old and new third-party API in forum		-	
	threads concerning migration of the us-			
	age of the given third-party API			
U60	Able to analyse level of user activity in	Shall	Completed	Task 3.3
	issue trackers			
U61	Able to identify if a thread reached a	Should	Completed	Task 3.3
	conclusion and if it was positive			
4. Ana	lysis of Docs/Code Snippets			
U62	Able to extract text from HTML and	Shall	Completed	Task 3.4
	markdown to feed natural language			
	analysis and identify code snippets			
U63	Able to extract text from PDF to feed	May	In Progress	Task 3.4
	natural language analysis and identify			
	code snippets			
U64	Provides recommendations to add doc-	Should	Planning	Task 3.4
	umentation commonly found in suc-			
	cessful projects			
U65	Provides recommendations to improve	Should	Planning	Task 3.4
	the structure of the documentation			
U66	Able to analyse Java code snippets	Shall	Planning	Task 3.4
U67	Able to analyse JavaScript code snip-	Should	Planning	Task 3.4
	pets			
U68	Able to analyse C code snippets	Shall	Planning	Task 3.4
U69	Able to analyse PHP code snippets	Shall	Planning	Task 3.4



Ref	Description	Priority	Progress	Notes
U72	Able to determine migration pattern	Shall	Planning	Task 3.4
	from two (or more) code snippets when			
	one of them uses the old third-party			
	API and the other uses the new third-			
	party API			
U73	Able to identify the part of the API that	Shall	Planning	Task 3.4
	the developer is currently using to pro-			
	vide code snippets in relation with cur-			
	rent development activity			
U74	Able to analyse the API documenta-	May	Planning	Task 3.4
	tion (when available) and determine if			
	it matches the current API of the library			
8. Proj	ect Documentation Analysis			
U110	Able to analyse PDF documents	May	Planning	Task 3.4
U111	Able to identify the documentation	Should	Planning	Task 3.4
	contains a Getting Started			
U112	Able to identify if the documentation	Should	Planning	Task 3.4
	contains a User Guide			
U113	Able to identify if the documentation	Should	Planning	Task 3.4
	contains a Developer Guide			
U114	Able to identify if the documentation	Should	Planning	Task 3.4
	contains code snippets			
U115	Able to analyse if the documentation	Shall	Planning	Task 3.4
	has a license			
U116	Able to analyse readability of docu-	Shall	Planning	Task 3.4
	mentation			
U117	Able to identify the list of changed	Should	Planning	Task 3.4
	third-party API methods from docu-			
	mentation			
U118	Able to identify the list of deprecated	Should	Planning	Task 3.4
	third-party API methods from docu-			
	mentation			
U119	Able to look up the documentation of	Shall	Planning	Task 3.4
	the public API provided by a project			
U120	Able to look up the source code exam-	Shall	Planning	Task 3.4
	ples of a public API of project			
U121	Able to analyse documentation to col-	Should	Planning	Task 3.4
	lect information related to the usage of			
	the API provided by the project			

Table 25: WP3 Use Case Requirements - Current progress

Ref	Description	Priority	Progress	Notes
D32	The NLP analysis component should	Should	In Progress	Task 3.3. In collabo-
	record the sources used by the devel-			ration with FEA
	oper and use them as a form of auto-			
	matic feedback to improve its sugges-			
	tions.			



Ref	Description	Priority	Progress	Notes
D72	The other mining tools developed in	Shall	Completed	Task 3.3
	WPs 2-4 shall expose a REST API			
	that the cross-project relationship min-			
	ers can consume			
D73	The other mining tools developed in	Shall	Completed	Task 3.3
	WPs 2-4 shall expose project level met-			
	rics that cross-project relationship min-			
	ers can consume			

Table 26: WP3 Technology Requirements - Current progress

At the time of writing this deliverable, there is one ongoing task in WP 3, Task 3.4. This task focuses on the development of a recommender that suggests code snippets and discussion from social media relevant to the code that is being developed in the Integrated Development Environment (IDE). For the recommendation we shall consider a variety of resources including code, documentation and online discussions. The deliverable associated with this task is due during M30 of the project.

## 7.2 Future Work

With regards to WP3, our efforts will be divided between several areas. As discussed in Section 3, GitLab and Bugzilla readers will be developed and integrated into CROSSMINER once RESTmule has the capabilities to support both sources. Metric providers associated with forums will be tested and integrated once the issue with the forums has been fixed. We will also address any issues that arise due to modifications to other CROSSMINER components that have been developed outside the scope of this deliverable. We need to conduct further investigation and discussions surrounding requirement D32, shown in Table 26. However, the majority of our efforts will focus on Task 3.4, the development of a recommender that suggests code snippets and discussions from social media relevant to the code that is being developed in the Integrated Development Environment (IDE).