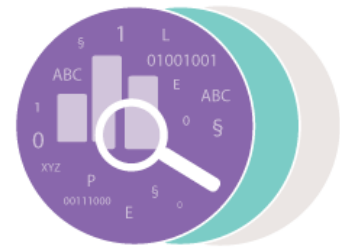




# FutureTDM

Explore . Analyse . Improve



REDUCING BARRIERS AND INCREASING UPTAKE OF TEXT AND DATA MINING FOR RESEARCH ENVIRONMENTS USING A COLLABORATIVE KNOWLEDGE AND OPEN INFORMATION APPROACH

## Deliverable D4.4

### Scientific Review Paper

## Project

Acronym: **FutureTDM**

Title: Reducing Barriers and Increasing Uptake of Text and Data Mining for Research Environments using a Collaborative Knowledge and Open Information Approach

Coordinator: SYNNO GmbH

Reference: 665940

Type: Collaborative project

Programme: HORIZON 2020

Theme: GARRI-3-2014 - Scientific Information in the Digital Age: Text and Data Mining (TDM)

Start: 01. September, 2015

Duration: 24 months

Website: <http://www.futuretdm.eu/>

E-Mail: [office@futuretdm.eu](mailto:office@futuretdm.eu)

Consortium: **SYNNO GmbH**, Research & Development Department, Austria, (SYNNO)  
**Stichting LIBER**, The Netherlands, (LIBER)  
**Open Knowledge**, UK, (OK/CM)  
**Radboud University**, Centre for Language Studies The Netherlands, (RU)  
**The British Library Board**, UK, (BL)  
**Universiteit van Amsterdam**, Inst. for Information Law, The Netherlands, (UVA)  
**Athena Research and Innovation Centre in Information, Communication and Knowledge Technologies**, Inst. for Language and Speech Processing, Greece, (ARC)  
**Ubiquity Press Limited**, UK, (UP)  
**Fundacja Projekt: Polska**, Poland, (FPP)

## Deliverable

Number:	<b>D4.4</b>
Title:	<b>Scientific Review Paper</b>
Lead beneficiary:	<b>RU</b>
Work package:	WP4: Fields of Application, Projects, Best Practices and Resources
Dissemination level:	Public (PU)
Nature:	Report (RE)
Due date:	30.04.2017
Submission date:	02.05.2017
Authors:	<b>Maria Eskevich, RU</b> <b>Antal van den Bosch, RU</b>
Contributors:	<b>Stelios Piperidis, ARC</b> <b>Kanella Pouli, ARC</b> <b>Maria Gavriilidou, ARC</b> <b>Dimitris Galanis, ARC</b>
Review:	<b>Burcu Akinci, SYNYO</b>

**Acknowledgement:** This project has received funding from the European Union's Horizon 2020 Research and Innovation Programme under Grant Agreement No 665940.

**Disclaimer:** The content of this publication is the sole responsibility of the authors, and does not in any way represent the view of the European Commission or its services.

This report by FutureTDM Consortium members can be reused under the [CC-BY 4.0](https://creativecommons.org/licenses/by/4.0/) license (<https://creativecommons.org/licenses/by/4.0/>).

# Table of Contents

- 1 Summary..... 5
- 2 Work on the scientific paper ..... 6
  - 2.1 Criteria for selection of potential publication venues..... 6
  - 2.2 Information on the submitted paper ..... 8
- 3 Work on the FutureTDM knowledge base enrichment..... 9
  - 3.1 Collection of Projects..... 9
  - 3.2 Collection of Organisations ..... 9
  - 3.3 Collection of Methods ..... 10
  - 3.4 Collection of Tools ..... 11
- 4 Conclusions ..... 15
- 5 Annex..... 16

# List of Tables

- Table 1: Total number of projects in the FutureTDM Knowledge Base Collection structured by economic sector ..... 9
- Table 2: Total number of projects in the FutureTDM Knowledge Base Collection structured by country ..... 10
- Table 3: FutureTDM text and data mining methods collection ..... 11
- Table 4: Total list of tools in the FutureTDM Knowledge Base Collection ..... 14

## 1 SUMMARY

This report provides information about the work on the scientific review paper that gives a critical review of the past, present and future of the state-of-the-art of the field of text and data mining (TDM). This publication aims to broaden perspective on the past, current, and future trends in TDM amongst audience, outlined via desktop research, consultation with the practitioners (Deliverables D.4.3. and D4.5.) and scientific communities. Report includes brief description of publication venue selection, text of the paper that can be further updated with the open accesses reviewers comments.

Additionally, the information about the update of FutureTDM knowledge base collection of relevant projects, organisations, methods and tools is provided including the current list of the collections in the Annex.

## 2 WORK ON THE SCIENTIFIC PAPER

### 2.1 Criteria for selection of potential publication venues

We have defined the two main criteria for the target publication venues:

- Diverse type of TDM applications already present in the journal
- The journal has to be Open Access, or at least to have the open access publication option

As the review of TDM work from the perspective of overall scientific development and insights on its visibility across different economic fields, that are formulated within the FutureTDM consortium, differ from the traditional type of publication content that is expected by journals, we took an initiative to contact a number of targeted journals in order to find a venue that would fit our criteria and to ensure that our content would be suitable. These were the journals to contact:

#### 1. Data Science Journal

- Link: <http://datascience.codata.org>
- Aims and scope of the journal: The CODATA Data Science Journal is a peer-reviewed, open access, electronic journal, publishing papers on the management, dissemination, use and reuse of research data and databases across all research domains, including science, technology, the humanities and the arts. The scope of the journal includes descriptions of data systems, their implementations and their publication, applications, infrastructures, software, legal, reproducibility and transparency issues, the availability and usability of complex datasets, and with a particular focus on the principles, policies and practices for open data. All data is in scope, whether born digital or converted from other sources.
  - *Reviews* can cover topics such as current controversies, the current “state of the art” or the historical development of studies as well as issues of regional or temporal focus. Papers should critically engage with the relevant body of extant literature. Review articles should be no longer than 3,000 words in length.
- Open access policy: multidisciplinary Open Access journal

#### 2. PLOSOne

- Link: <http://journals.plos.org/plosone/>
- Aims and scope of the journal: PLOS ONE features reports of original research from all disciplines within science and medicine. By not excluding research on the basis of subject area, PLOS ONE facilitates the discovery of connections between research whether within or between disciplines.

- We consider publishing systematic reviews only if the methods ensure the comprehensive and unbiased sampling of existing literature.
- Submissions describing methods, software, databases, or other tools. We consider submissions describing methods, software, databases, or other tools if they follow the appropriate reporting guidelines.
- Qualitative research. We consider publishing qualitative research only if it adheres to appropriate study design and reporting guidelines.
- Studies reporting negative results.
- Open access policy: multidisciplinary Open Access journal

### 3. DSH (Digital Scholarship for Humanities)

- Link: <https://academic.oup.com/dsh>
- Aims and scope of the journal: *DSH* or *Digital Scholarship in the Humanities* is an international, peer reviewed journal which publishes original contributions on all aspects of digital scholarship in the Humanities including, but not limited to, the field of what is currently called the Digital Humanities. Long and short papers report on theoretical, methodological, experimental, and applied research and include results of research projects, descriptions and evaluations of tools, techniques, and methodologies, and reports on work in progress. *DSH* also publishes reviews of books and resources.
- Open access policy: multidisciplinary Open Access journal

### 4. Science Communication (SAGE)

- Web-Link: <http://journals.sagepub.com/home/scx>
- Aims and scope of the journal: Science Communication is an international, interdisciplinary social science journal that examines such topics as the nature of scientific expertise as represented through communication and the processes or effects characterizing the communication of science in any context. Science is broadly defined to include environmental science, health science, and technology. Science Communication welcomes submissions of empirical research from authors in all relevant disciplines (including the social sciences, the humanities, and science itself). Both qualitative and quantitative research papers with a basis in theory are acceptable. Preference is given to articles that bridge the gap between theory and practice and that will be of interest across disciplines. In addition to peer-reviewed research, Science Communication publishes commentaries that analyse issues and trends in the field – whether scholarly, professional, or policy-related – and a periodic summary of new books in the field.
- Open access policy: OA option is available

### 5. IEEE Access

- Web-Link: <http://ieeaccess.ieee.org/learn-more-about-ieee-access/>

- Aims and scope of the journal: IEEE Access publishes articles that are of high interest to readers: original, technically correct, and clearly presented. The scope of this journal comprises all of IEEE's fields of interest, emphasizing applications-oriented and interdisciplinary articles. IEEE Access also accepts traditional technical articles, as well as reviews and surveys.
- Open access policy: Multidisciplinary Open Access Journal, fee \$1,750 per article

Eventually, we have chosen the “Data Science Journal” that is an open access multidisciplinary journal run by an SME publisher “Ubiquity Press”.

## 2.2 Information on the submitted paper

The details of the submission are the following:

- Authors: Maria Eskevich, Stelios Piperidis, Kanella Pouli, Maria Gavriilidou, Dimitris Galanis, Antal van den Bosch.
- Title: Text and Data Mining: Past, present, and future.
- Submission date: 01.05.2017.
- Text of the paper is in the Appendix



## 3 WORK ON THE FUTURETDM KNOWLEDGE BASE ENRICHMENT

### 3.1 Collection of Projects

FutureTDM collection of projects by the Deliverable submission date contains 140 items, 37 of which are already integrated into the platform. For 7 of the projects it is not possible to assign an economic sector, such as (+)Spaces , while the other 133 can be assigned to one or even a combination of sectors due to the multidisciplinary nature of the projects. Table 1 contains the overall statistics of the number of projects in the collection.

Sector or combination of Sectors that that the project is assigned to	Number of projects
Primary-Sector, Secondary-Sector, Tertiary-Sector, Quaternary-Sector	1
Primary-Sector, Quaternary-Sector	1
Secondary-Sector	8
Secondary-Sector, Tertiary-Sector	3
Secondary-Sector, Quaternary-Sector	38
Secondary-Sector, Tertiary-Sector, Quaternary-Sector	6
Tertiary-Sector	13
Tertiary-Sector, Quaternary-Sector	24
Quaternary-Sector	38
Quinary-Sector	1
<b>TOTAL</b>	<b>133</b>

**Table 1: Total number of projects in the FutureTDM Knowledge Base Collection structured by economic sector**

### 3.2 Collection of Organisations

FutureTDM collection of organisations by the Deliverable submission date contains 72 items, 59 of which are already integrated into the platform. 55 of these organisations have only one EU-country location, while three of them spread across more than one country within the EU, and 6 have locations outside EU in USA.

Company locations details	Country	Number
Company is based in one country	Austria	1
	Belgium	2
	Czech Republic	2
	Denmark	2
	Finland	1
	France	4
	Germany	10
	Greece	3
	Ireland	1
	Italy	3
	Norway	1
	Slovenia	1
	Spain	4
	Sweden	2
	Switzerland	4
The Netherlands	5	
UK	9	
Company has offices internationally	Belgium, The Netherlands, Germany, USA	1
	The Netherlands,	1
	Germany/USA	1
	Ireland/Greece	1
	Spain/UK	1
	UK/USA	2
	USA, many	2
Companies outside EU	Japan	1
	USA	7
<b>TOTAL</b>		<b>72</b>

Table 2: Total number of projects in the FutureTDM Knowledge Base Collection structured by country

### 3.3 Collection of Methods

FutureTDM collection of projects splits scientific methods into 29 items, and all of them are integrated into the platform, as they are used to specify the organisations focus.

TDM Methods			
Artificial Intelligence	Machine Learning	Statistics	Kernel Methods
Association rules	Machine Translation	Summarisation	Ensemble Learning
Classification	Multimedia Processing	Term/Concept Extraction	Dimensionality Reduction
Clustering	Natural Language Processing	Text Mining	Decision Trees
Data Mining	Negation and Modality Detection	Textual Entailment	Deep Learning
Graph Mining	Predictive Analytics	Artificial Neural Networks	
Information Extraction	Regression	Bayesian Inference	
Information Retrieval	Sentiment Analysis/Opinion Mining	Instance-Based Learning	

**Table 3: FutureTDM text and data mining methods collection**

### 3.4 Collection of Tools

FutureTDM collection of projects by the Deliverable submission date contains 72 items, 59 of which are already integrated into the platform. Table contains the complete list of tools collected so far.

Name	URL
ALVEO	<a href="http://alveo.edu.au/">http://alveo.edu.au/</a>
Alvis	<a href="https://migale.jouy.inra.fr/redmine/projects/alvisnlp">https://migale.jouy.inra.fr/redmine/projects/alvisnlp</a>
Angoss Knowledge STUDIO	<a href="http://www.angoss.com/">http://www.angoss.com/</a>
AnnoMarket	<a href="https://annomarket.eu/">https://annomarket.eu/</a>
Apache cTAKES	<a href="http://ctakes.apache.org">http://ctakes.apache.org</a>
Apache OpenNLP	<a href="http://opennlp.apache.org">http://opennlp.apache.org</a>
Apache UIMA	<a href="https://uima.apache.org">https://uima.apache.org</a>
Argo	<a href="http://argo.nactem.ac.uk">http://argo.nactem.ac.uk</a>

BioCatalogue	<a href="https://www.biocatalogue.org/">https://www.biocatalogue.org/</a>
BiodiversityCatalogue	<a href="https://www.biodiversitycatalogue.org">https://www.biodiversitycatalogue.org</a>
Bluima	<a href="https://github.com/BlueBrain/bluima">https://github.com/BlueBrain/bluima</a>
CLARIN-DK	<a href="http://clarin.dk/">http://clarin.dk/</a>
ClearTK	<a href="https://cleartk.github.io/cleartk">https://cleartk.github.io/cleartk</a>
Clementine	<a href="http://datamining.togaware.com/survivor/Clementine.html">http://datamining.togaware.com/survivor/Clementine.html</a>
ContentMine	<a href="http://www.contentmine.org/">http://www.contentmine.org/</a>
DKPro Core	<a href="https://dkpro.github.io/dkpro-core">https://dkpro.github.io/dkpro-core</a>
ELKI	<a href="http://elki.dbs.ifi.lmu.de/">http://elki.dbs.ifi.lmu.de/</a>
FICO Data Management	<a href="http://www.fico.com/">http://www.fico.com/</a>
Galaxy	<a href="https://galaxyproject.org/">https://galaxyproject.org/</a>
GATE Embedded	<a href="https://gate.ac.uk">https://gate.ac.uk</a>
Heart of Gold	<a href="http://heartofgold.dfki.de">http://heartofgold.dfki.de</a>
IBM SPSS Predictive Analysis	<a href="http://www-03.ibm.com/software/products/en/spss-predictive-analytics-enterprise">http://www-03.ibm.com/software/products/en/spss-predictive-analytics-enterprise</a>
JCoRe	<a href="http://julielab.github.io">http://julielab.github.io</a>
Jpylyzer	<a href="http://jpylyzer.openpreservation.org/">http://jpylyzer.openpreservation.org/</a>
KAF	<a href="http://kyoto-project.eu/xmlgroup.iit.cnr.it/kyoto/index2091.html?option=com_content&amp;view=article&amp;id=504&amp;Itemid=173">http://kyoto-project.eu/xmlgroup.iit.cnr.it/kyoto/index2091.html?option=com_content&amp;view=article&amp;id=504&amp;Itemid=173</a>
Kepler	<a href="https://kepler-project.org/">https://kepler-project.org/</a>
KNIME	<a href="http://www.knime.org/knime">http://www.knime.org/knime</a>
Language Grid	<a href="http://langrid.org/en/index.html">http://langrid.org/en/index.html</a>
LAPPS Grid	<a href="http://www.lappsgrid.org/">http://www.lappsgrid.org/</a>
Massive Online Analysis (MOA)	<a href="http://moa.cms.waikato.ac.nz/">http://moa.cms.waikato.ac.nz/</a>
Matchbox	<a href="http://matchbox.openpreservation.org/">http://matchbox.openpreservation.org/</a>
Microsoft SQL Server	<a href="https://technet.microsoft.com/en-">https://technet.microsoft.com/en-</a>

Analysis Service (SSAS)	<a href="http://us/library/ms175609%28v=sql.90%29.aspx">us/library/ms175609%28v=sql.90%29.aspx</a>
Neural Designer	<a href="https://www.neuraldesigner.com/">https://www.neuraldesigner.com/</a>
NLTK	<a href="http://www.nltk.org">http://www.nltk.org</a>
Open Calais	<a href="http://www.opencalais.com/">www.opencalais.com/</a>
Oracle Data Miner GUI	<a href="http://www.oracle.com/technetwork/database/options/odm/dataminerworkflow-168677.html">http://www.oracle.com/technetwork/database/options/odm/dataminerworkflow-168677.html</a>
Orange	<a href="http://orange.biolab.si/">http://orange.biolab.si/</a>
Pagelyzer	<a href="http://pagelyzer.openpreservation.org/">http://pagelyzer.openpreservation.org/</a>
Pegasus	<a href="https://pegasus.isi.edu">https://pegasus.isi.edu</a>
Pentaho	<a href="http://www.pentaho.com/">http://www.pentaho.com/</a>
Pipeline Pilot	<a href="http://accelrys.com/products/collaborative-science/biovia-pipeline-pilot/">http://accelrys.com/products/collaborative-science/biovia-pipeline-pilot/</a>
QT21	<a href="http://qt21.metashare.ilsp.gr/">http://qt21.metashare.ilsp.gr/</a>
Rapidminer Radoop	<a href="http://www.rapidminer.com">www.rapidminer.com</a>
Rapidminer Server	<a href="http://www.rapidminer.com">www.rapidminer.com</a>
Rapidminer Studio	<a href="http://www.rapidminer.com">www.rapidminer.com</a>
SAP Hana	<a href="https://hana.sap.com">https://hana.sap.com</a>
SAS Data Mining	<a href="http://www.sas.com">http://www.sas.com</a>
SPSS	<a href="http://www.ibm.com/analytics/us/en/technology/spss/">http://www.ibm.com/analytics/us/en/technology/spss/</a>
Taverna	<a href="http://www.taverna.org.uk/">http://www.taverna.org.uk/</a>
Think Analytics	<a href="http://thinkanalytics.com/">http://thinkanalytics.com/</a>
Think Enterprise Data Miner	<a href="http://www.thinkanalytics.com/">www.thinkanalytics.com/</a>
Triana	<a href="http://www.trianacode.org">http://www.trianacode.org</a>
TTNWW	<a href="http://yago.meertens.knaw.nl/apache/TTNWW/">http://yago.meertens.knaw.nl/apache/TTNWW/</a>
Weblicht	<a href="https://weblicht.sfs.uni-tuebingen.de/">https://weblicht.sfs.uni-tuebingen.de/</a>
WEKA	<a href="http://www.cs.waikato.ac.nz/~ml/weka">http://www.cs.waikato.ac.nz/~ml/weka</a>
xcorrSound	<a href="http://xcorrSound.openpreservation.org/">http://xcorrSound.openpreservation.org/</a>

VisTrails	<a href="http://www.vistrails.org">http://www.vistrails.org</a>
Trendminer	<a href="https://www.trendminer.com/">https://www.trendminer.com/</a>
Opendatasoft	<a href="https://www.opendatasoft.com/">https://www.opendatasoft.com/</a>
Meta	<a href="http://meta.com/">http://meta.com/</a>
TensorFlow	<a href="https://www.tensorflow.org/about/">https://www.tensorflow.org/about/</a>
Epinium	<a href="http://epinium.com/">http://epinium.com/</a>
Exclusivi	<a href="http://exclusivi.com/hotels/">http://exclusivi.com/hotels/</a>
Eparkomat	<a href="http://www.eparkomat.com/">http://www.eparkomat.com/</a>
Data Scouts	<a href="https://datascouts.eu/">https://datascouts.eu/</a>
TopPlace	<a href="http://www.avuxi.com/products/topplace">http://www.avuxi.com/products/topplace</a>
Influency	<a href="https://influency.com/en/">https://influency.com/en/</a>
Plazi TreatmentBank	<a href="http://plazi.org/api-tools/api/#Plazi_API">http://plazi.org/api-tools/api/#Plazi_API</a>
SafeToNet	<a href="https://safetonet.com/">https://safetonet.com/</a>
Egas	<a href="https://demo.bmd-software.com/egas/">https://demo.bmd-software.com/egas/</a>
Wordfreak	<a href="http://wordfreak.sourceforge.net/">http://wordfreak.sourceforge.net/</a>
Theano	<a href="https://github.com/Theano/">https://github.com/Theano/</a>

**Table 4: Total list of tools in the FutureTDM Knowledge Base Collection**

## 4 CONCLUSIONS

This Deliverable contains the full text of submitted TDM review paper, together with the explanation on the publication venue selection.

We provide quantitative updates on the progress with filling in the four main categories of the Knowledge Base Collection, while the Deliverables 6.1, 6.2 and 6.3<sup>1</sup> contains the qualitative description of the technology behind the FutureTDM Platform.

---

<sup>1</sup> [www.futuretdm.eu](http://www.futuretdm.eu)

## 5 ANNEX

In the following we present the scientific paper “**Text and Data Mining: Past, present, and future**” and the collections below described in the previous section. The list will be given in this order:

- Collection of Organisations
- Collection of Projects
- Collection of Tools





Collection of Projects

Project ID	Project Name	Client	Start Date	End Date	Status	Phase	Progress (%)	Team Lead	Team Members	Budget (€)	Actual Cost (€)	Revenue (€)	Profit (€)	ROI (%)
001	Project A	Client A	2023-01-01	2023-03-31	Completed	Phase 1	100	John Doe	Team A	100000	95000	120000	25000	25%
002	Project B	Client B	2023-02-01	2023-04-30	In Progress	Phase 2	75	Jane Smith	Team B	150000	140000	180000	40000	27%
003	Project C	Client C	2023-03-01	2023-05-31	On Hold	Phase 1	20	Mike Johnson	Team C	80000	80000	0	-80000	-100%
004	Project D	Client D	2023-04-01	2023-06-30	Completed	Phase 3	100	Sarah Lee	Team D	120000	115000	150000	35000	29%
005	Project E	Client E	2023-05-01	2023-07-31	In Progress	Phase 1	50	David Kim	Team E	90000	85000	110000	25000	28%
006	Project F	Client F	2023-06-01	2023-08-31	On Hold	Phase 2	30	Emily White	Team F	110000	110000	0	-110000	-100%
007	Project G	Client G	2023-07-01	2023-09-30	Completed	Phase 3	100	Chris Brown	Team G	130000	125000	160000	35000	27%
008	Project H	Client H	2023-08-01	2023-10-31	In Progress	Phase 1	60	Alex Green	Team H	100000	95000	130000	35000	35%
009	Project I	Client I	2023-09-01	2023-11-30	On Hold	Phase 2	40	Mia Black	Team I	140000	140000	0	-140000	-100%
010	Project J	Client J	2023-10-01	2023-12-31	Completed	Phase 3	100	Noah Grey	Team J	160000	155000	200000	45000	28%

