# Text Mining: the next data frontier

# Repositories in the centre of new scientific knowledge

openM1N7ED

Natalia Manola
Athena Research &
Innovation Centre

#openminted_eu, #or2016, #tdm

# Some facts About scientific literature

The global research community generates over 1.5 million new scholarly articles per annum.

<div align="right">The STM report (2009)</div>

… some 90% of papers … are never cited.
… 50% of papers are never read by anyone other than their authors, referees and journal editors

<div align="right">Lokman I. Meho, The rise and rise of citation analysis, 2007</div>

… one paper published every 30 seconds

… 70,000 papers published on a single protein, the tumor suppressor p53

<div align="right">Spangler et al, Automated Hypothesis Generation based on Mining Scientific Literature, 2014</div>

openM1N7ED

# Emerging solution(S)

## Machine reading

process textual sources, organise and classify in various dimensions, extract main (indexical) information items,

## … and "understanding"

identify and extract entities and relations between entities, facilitate the transformation of unstructured textual sources into structured data

## … and predicting

enable the multidimensional analysis of structured data to extract meaningful insights and improve the ability to predict

openM1N7ED

# What OpenMinted is About

# MAIN Objectives

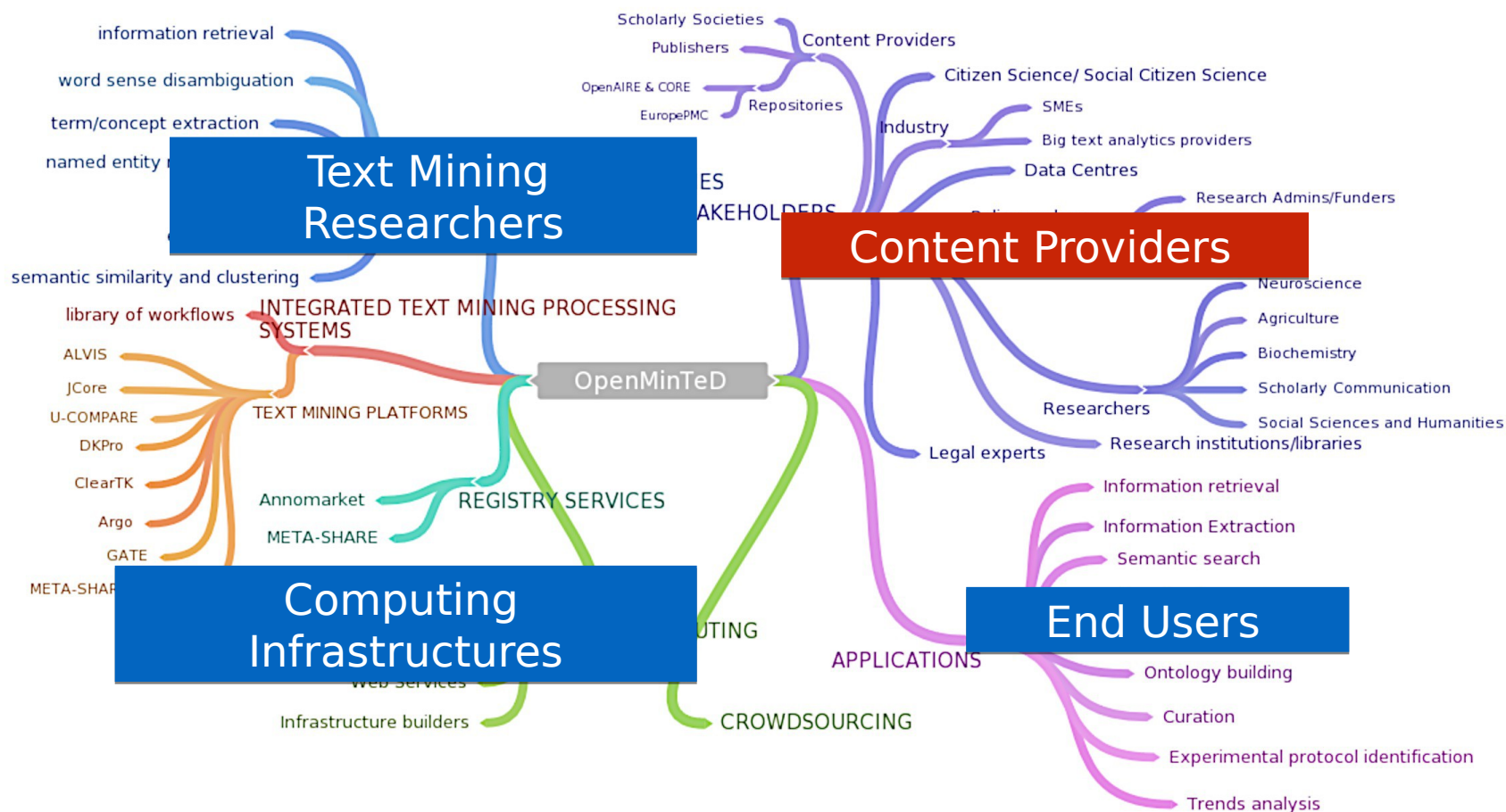Establish an **open** and **sustainable** Text and Data Mining (TDM) **platform** and **infrastructure** where researchers can discover, collaboratively create, share and re-use knowledge from a wide range of text based **scientific and scholarly**

## A next step from Open Access to Open Science

# A complex Landscape



information retrieval
word sense disambiguation
term/concept extraction
named entity r

**Text Mining Researchers**

semantic similarity and clustering
library of workflows INTEGRATED TEXT MINING PROCESSING SYSTEMS
ALVIS
JCore
U-COMPARE TEXT MINING PLATFORMS
DKPro
ClearTK
Argo
GATE
META-SHAR

**Computing Infrastructures**

Annomarket REGISTRY SERVICES
META-SHARE

Web Services
Infrastructure builders CROWDSOURCING

Scholarly Societies
Publishers Content Providers
OpenAIRE & CORE
EuropePMC Repositories

ES
AKEHOLDERS

Citizen Science/ Social Citizen Science
SMEs
Industry Big text analytics providers
Data Centres

Research Admins/Funders

**Content Providers**

Neuroscience
Agriculture
Biochemistry
Scholarly Communication
Researchers Social Sciences and Humanities
Legal experts Research institutions/libraries

Information retrieval
Information Extraction
Semantic search

**End Users**

Ontology building
Curation
Experimental protocol identification
Trends analysis

OpenMinTeD

APPLICATIONS

openM1N7ED

European Commission

# HIGH LEVEL ARCHITECTURE



Users: researchers, curators, text-miners and new services developers

Registry | Auth2 & Policy management | Workflow Management | Annotator | Accounting

Interoperability of text mining services Compatibility of text-mining components

Mining Platforms — GATE — general architecture for text engineering

Mining Platforms — Unstructured Information Management Architecture

Mining Platforms — Proprietary architectures

Mining Platforms — NLTK 3.0

Interoperability of language resources & corpora

Language resources and corpora registry service

Language resources

Publisher text corpus | Other text corpora | OpenAIRE/CORE text corpus | Other text corpora | PMC text corpus | Other text corpora | Other text corpora

Access Interoperability to shared storage and computing resources

Data centre 1 | Data centre 2 | Data centre 3 | Data centre 4 in public cloud

Policies & guidelines

openM1N7ED

European Commission

# Key Characteristics

**1** service oriented – discovery, re-use of content and tools

**2** build on existing TDM tools - no focus on new algorithms

**3** infrastructure – focus on interoperability

**4** community driven - user centric requirements

**5** open science - openness at all levels

openM1N7ED

European Commission

# Challenges

## Discoverable & accessible content & services

- Document literature cont[ent]
  data categories taxonomies, provenance information
- Document langu[age]
  workflows
- Generic and do[main specific metadata descriptions]

## Interoper[ability]

- Combine services i[n]
- Combine content and language resources with services and workflows
- Combine automat[ed]
  services

## IPR and licensing

- Stu[dy] IPR restrictions for reuse of sources as well as possible

**Starting with repositories and OA publishers
via OpenAIRE and CORE**

**Building on existing language resources repositori[es] and infras (meta-share, clarin)**

**Promoting existing standards and best practice[s] AND technologies**

**In close collaboration with the FUTURETDM proje[ct]
http://project.futuretdm.eu/**

# Community Driven

om the very beginning...
**equirements**, content, barriers, expected outcomes.

to the very end
eate applications, **validate and evaluate** the results.

## Scholarly Comm.

Feature extraction

Data citation

Research analysis

## Life Sciences

Curation of databases and lexica in Chembolomics & neuroinformatic s

## Agriculture

Extracting information from tables for food safety alerts

## Social Sciences

Data citation

openMIN7ED

European Commission

openM1N7ED

# THANK YOU!

Natalia Manola
natalia@di.uoa.gr

twitter.com/openminted_eu

facebook.com/openminted

bit.do/openmintedlinkedin

vimeo.com/openminted

bit.do/openmintedplus