# Text Analytics in the EU Fusepool project
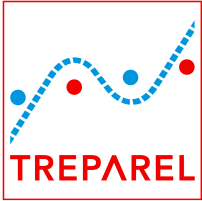
II-SDV '13 - Nice

**TREPAREL**

**Treparel**
Delftechpark 26
2628 XH Delft
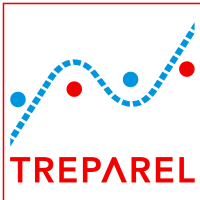The Netherlands
www.treparel.com

**Anton Heijs**
CTO
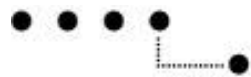anton@treparel.com

April 15, 2013

# Agenda

- Introduction of the EU Fusepool project, Treparel and the consortium partners

- Background on objectives of the Fusepool project

- User adaptive system

- Data pooling and linking

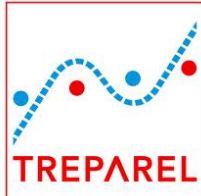- Large scale machine learning and text analytics

- Summary

# The EU Fusepool Consortium

**Fusepool**

**Bern University of Applied Sciences**

**searchbox**

**European Network of Living Labs**
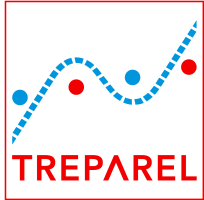
**SEVENTH FRAMEWORK PROGRAMME**

**xerox**

**Geo**

Treparel (Delft, The Netherlands) is a global software provider in Big Data Text Analytics and Visualization.

Global companies, government agencies, software vendors or data publishers are using Treparel KMX text analysis software to gain faster, reliable, precise insights in large complex unstructured data sets (like application notes, blogs, email and patents) allowing them to make better informed decisions.

As part of the Fusepool consortium Treparel integrates the KMX Text Analytics technology for advanced classification, clustering and visualization of large complex document collections.

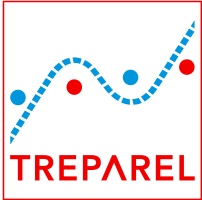www.treparel.com  3

# The Big Picture

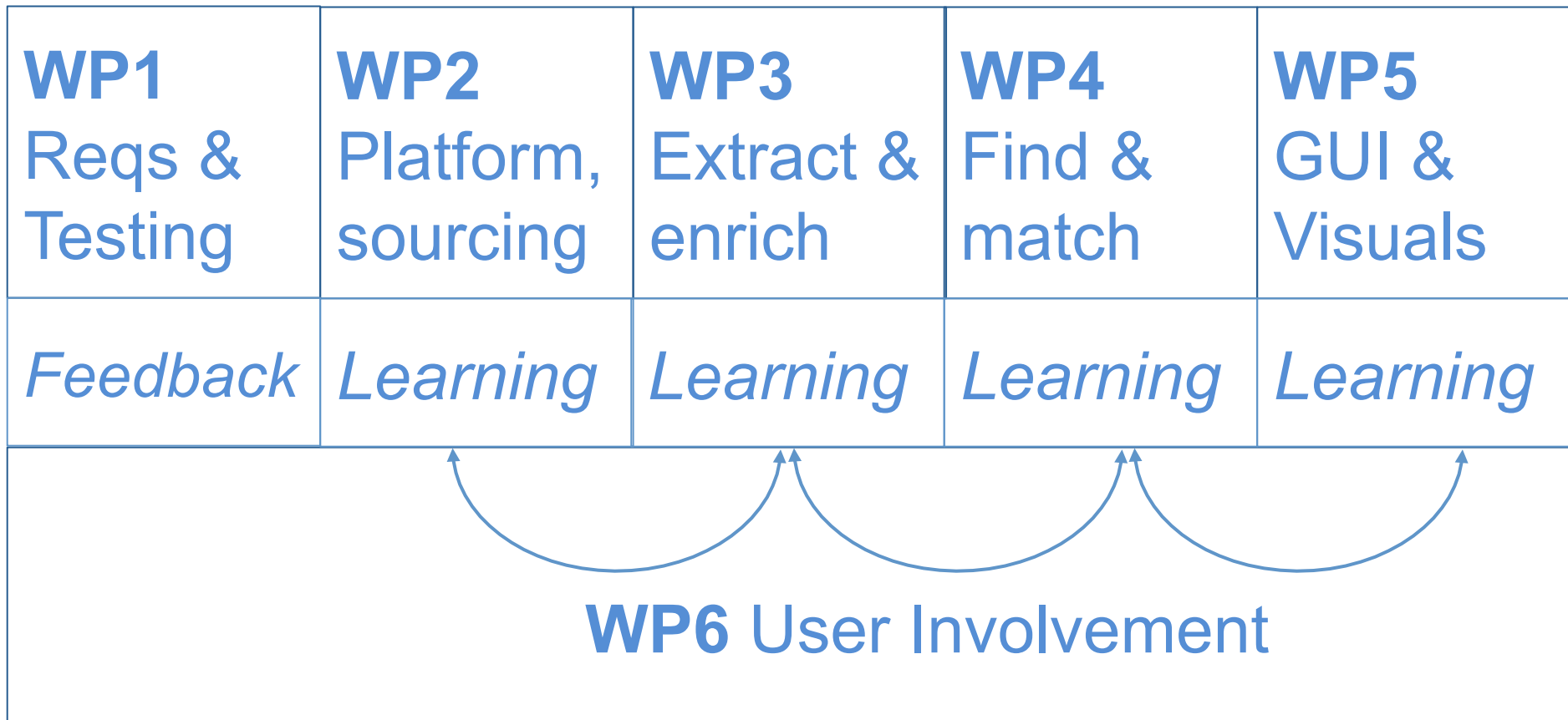Each partner provides critical **building blocks:**

- ✓ **BUAS**: Dynamic user interfaces
- ✓ **ENOLL**: End users with specific needs
- ✓ **GEOX**: NER, UI design
- ✓ **SEARCHBOX**: text search and similarity semantic matching
- ✓ **TREPAREL**: text/patent classification engines
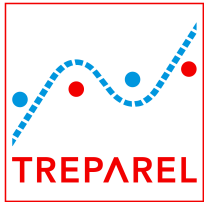- ✓ **XEROX**: user-adaptive learning know-how

**BUT BREAK-THROUGH BY INTEGRATION**

**Fusepool**

| WP1 Reqs & Testing | WP2 Platform, sourcing | WP3 Extract & enrich | WP4 Find & match | WP5 GUI & Visuals |
|---|---|---|---|---|
| *Feedback* | *Learning* | *Learning* | *Learning* | *Learning* |

**WP6** User Involvement

The EC Fusepool project is a 2 year EU project and started in July 2012

www.treparel.com

Vision:                 User-adaptive system
Living Lab:             Rapid app development
Data processing:   Sourcing & interlinking
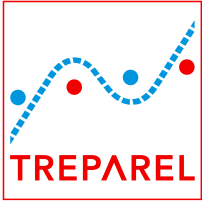Machine learning:  Matching & optimizing
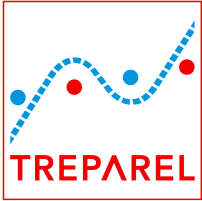




**Integrated Use Cases**

# Background

- SMEs have a need for technology intelligence for detecting and responding to opportunities and threats
    - This a partly driven by growth and complexity of patents and lawsuits
- Consumer intelligence to detect opinions and needs of consumers for product development
- Open innovation requiring cooperation (links between data, e.g. finding business partners)
- Focus: Machine Learning algorithms to improve matching
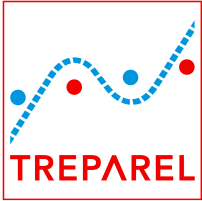
# User-adaptive system

- **Focus**: monitor and learn specific needs and preferences of a user to align features, functionalities, and graphical interfaces

- **Adaptive**: machine learning from crowd-sourcing (rather than rule-based)

- **User-aligned prioritization**: more usable and customized interfaces, suggestions based on activity & user feedback

www.treparel.com
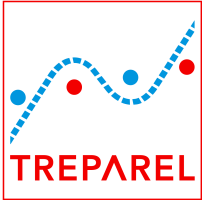
# User-adaptive matching

- **Main goal**: automated user-adaptive matching of users to:
  - Patent analysis
  - Finding funding opportunities
  - Partner matching
- **Key asset**: information provided by the user
- **User Data Credo**: accuracy improves with quantity and quality of user data while variety (breadth) increases with number of users
- **Living lab**: Co-creation between creators and consumers of the Fusepool platform
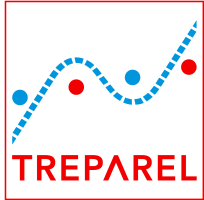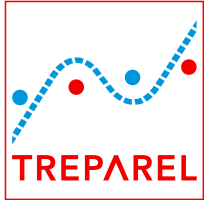
# Data sourcing

- **Sources**: internal & external content from web harvesting and structured data sources
  - using content databases and linked open data
- **Scope**: initial data corpus includes all explicitly in- and excluded sources
- **Gained value from Information**: recommendations based on machine learning from feedback
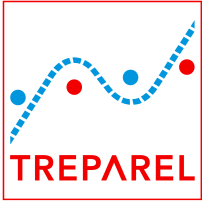
1.  **Text analysis and feature extraction**: ML & NLP methods for categorizing, named entity extraction, etc.

2.  **Shared metadata models**: mapping text features to existing/custom ontologies and generation of semantic triplets

    → High-level abstraction & persistence for reuse

    → Lightweight storage: mostly metadata only, text indexing and abstraction uses schema-free key-value (enabling actionable facets)

- **Contextualize**: terms are interlinked with same and similar terms across sources:
  - Enrich the extracted content with existing information available in the Internet
  - Interlink as much information as possible to increase the value of knowledge extraction
  - Use available public sector resources in Semantic Web and LOD format

- **Metadata**: when a user uploads texts to be matched with other content, only the metadata descriptors are transmitted

- **Data privacy**: data fusion from diverse sources without endangering user privacy
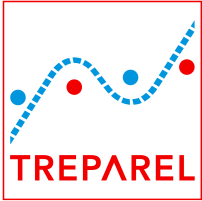
# Searching & finding

- **Key search-oriented features**:
  - Search through all content in the data pool
  - Faceted search (categories, metadata, entities)
  - Integration of Linked Open Data (LOD) results
  - Cross-lingual indexing and cross-referencing
  - "Did you mean?"-functionality in case of typos and auto-completion of search queries

- **User-adaptive**: indexing and integration based on users needs (e.g. user profiling)
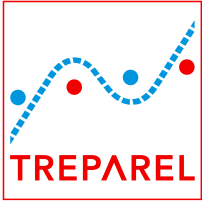
# Adaptation & refinement

- **Adaptive search**: results are aligned to user preferences based on analysis of user implicit and explicit feedback

- **Multi-task ranking**: good trade-off between user-independent search (high coverage but lower precision) and a very customized approach

- **Query intent discovery**: analysis of the query structure and interlinking of queries

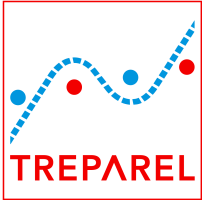www.treparel.com

# Correlating & matching

- **Search guided navigation**: semantic matching extracts contextual relationships to list related content
  - suggestions organized by categories
  - exposing facets within related content
- **Distributed rule and event model**: defines states, actions, and consequences (e.g. notifications, visualizations) for reasoning based on light-weight ontologies
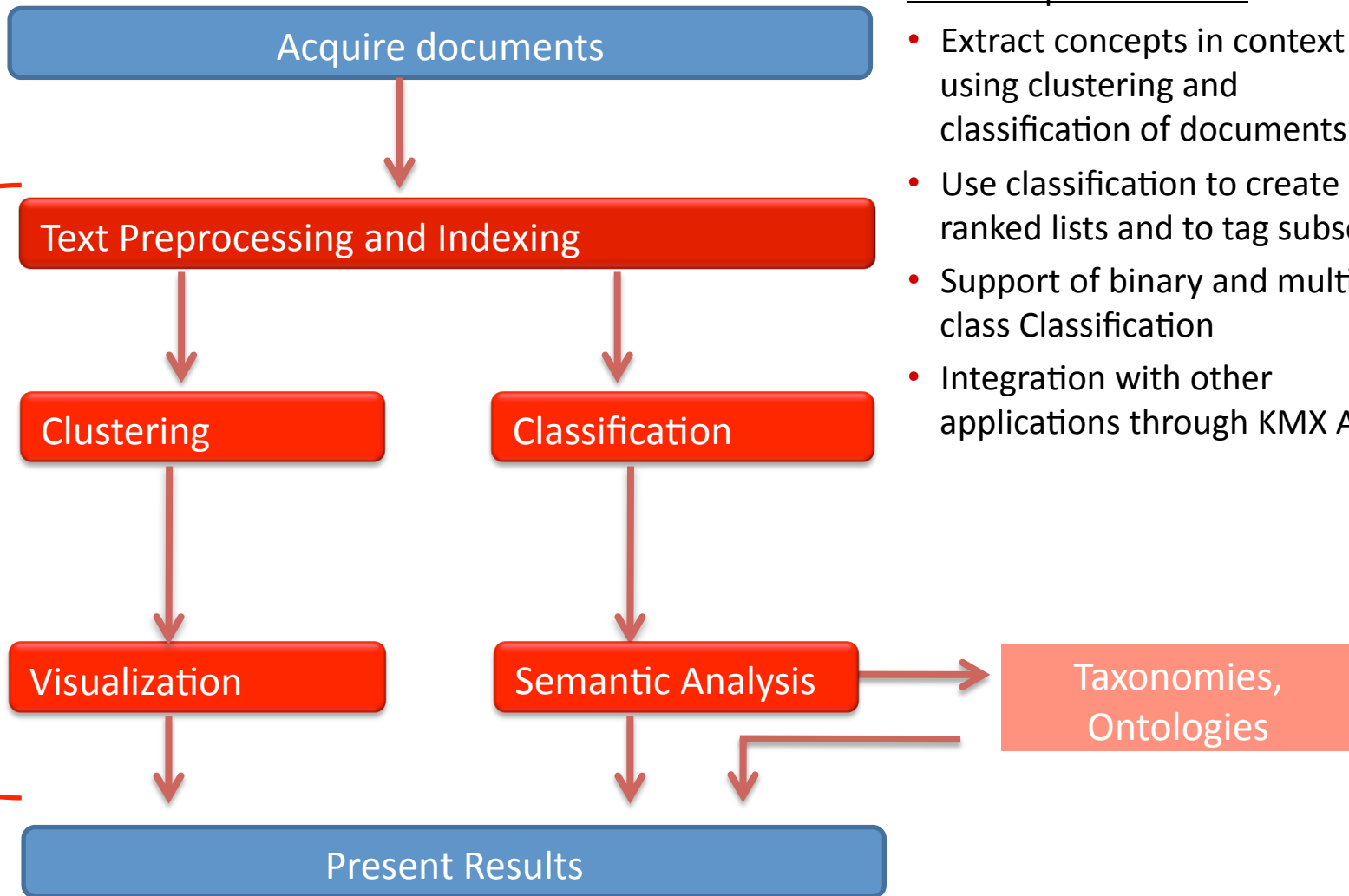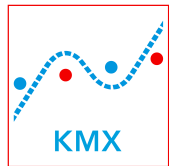
# Crowd sourcing & supervised automation

- **Relational learning**: related instances are used to reason about the focal instance
  - Rationality of content (links to other content, people, etc.) provide rich information
  - Similarities/dissimilarities to other content is established purely on relational properties

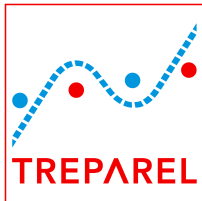# Document level analysis using the KMX technology

Query & Search Tools

**Acquire documents**

**Text Preprocessing and Indexing**

**Clustering**

**Classification**

**Visualization**

**Semantic Analysis**

**Taxonomies, Ontologies**

**Present Results**

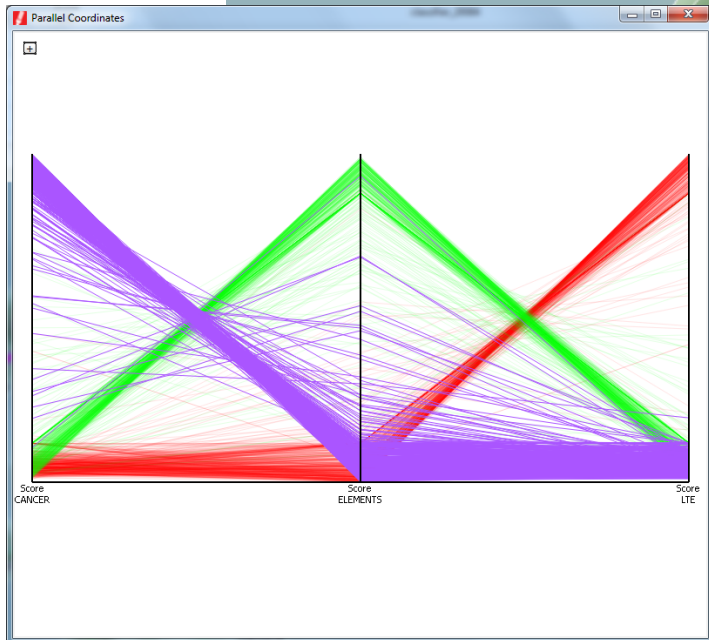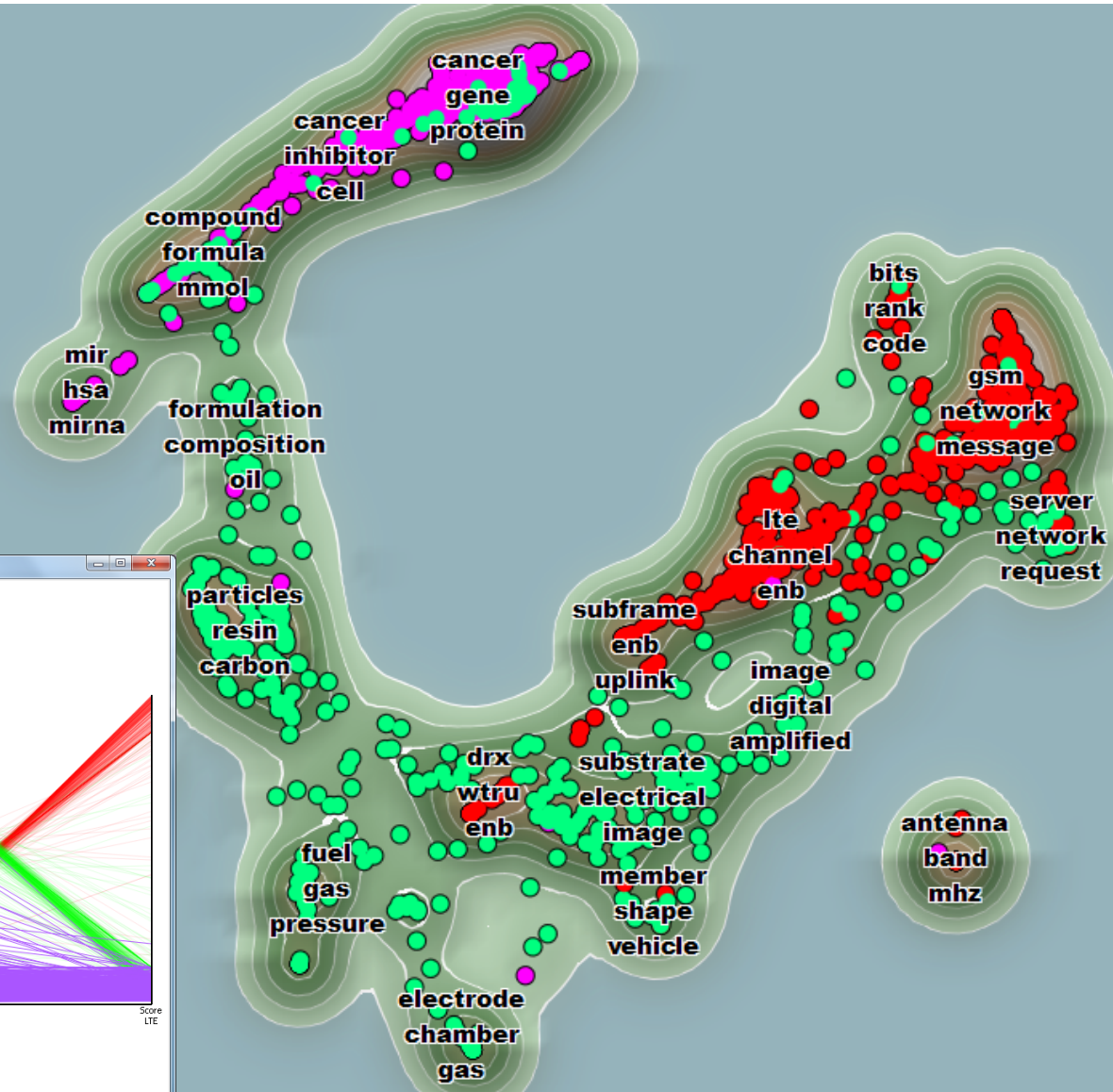KMX

## KMX unique functions:

- Extract concepts in context using clustering and classification of documents

- Use classification to create ranked lists and to tag subsets

- Support of binary and multi-class Classification

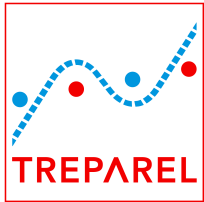- Integration with other applications through KMX API
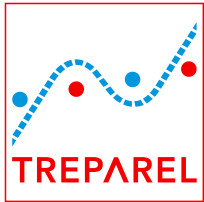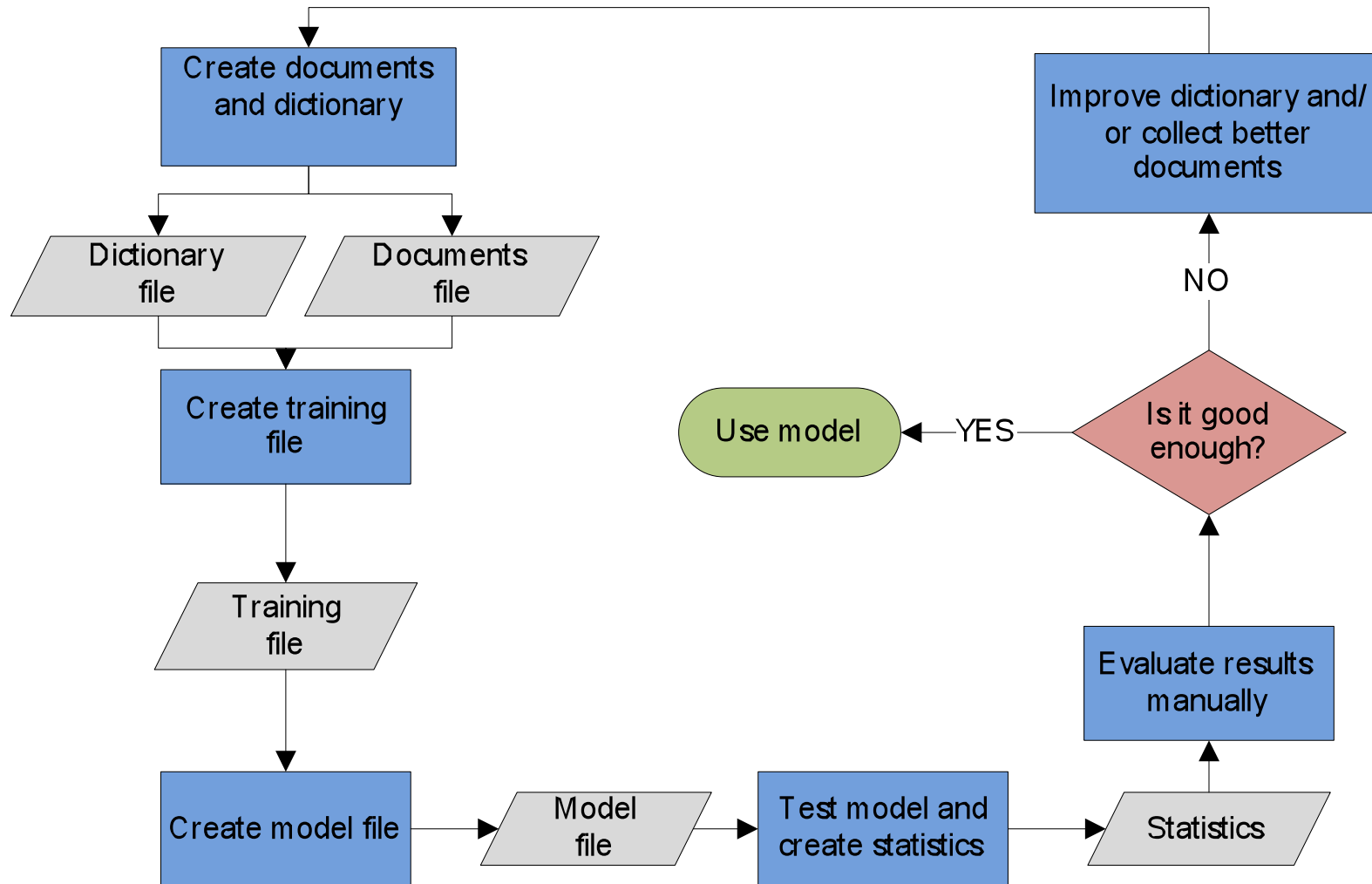
NER in the landscaping

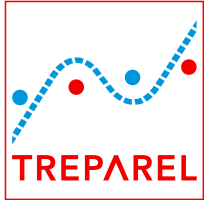# Sentence level analysis using Named Entity Recognition

- The aim of NER is to identify entities in unstructured text documents.
    - To locate : mark-up the entities
    - To classify : into predefined categories/domains
- Aim of usage
    - To recognise trends (trained and new high frequency)
    - To find all „trained" entities
    - To „discover" new entites
- NER approaches
    - Statistics based (supervised machine learning)
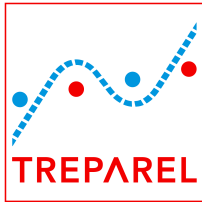    - Rule based (regular expressions)

# 3 NER model examples
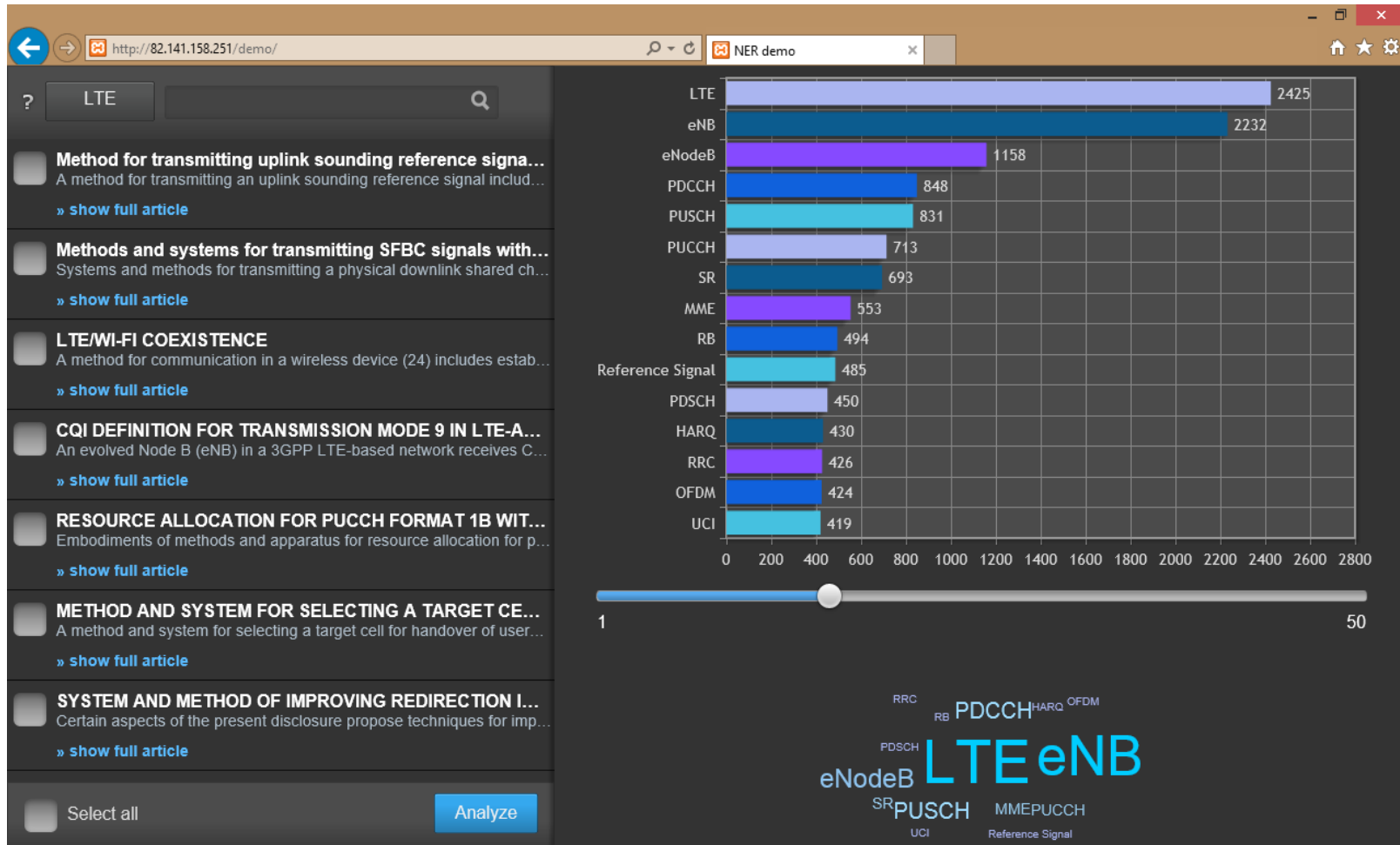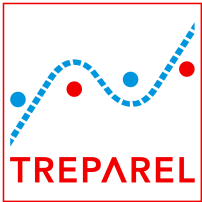
✓ Training text: 500 patents from EPO

✓ Training: model building by Stanford NER

- **LTE (long term evolution)**
  - F1: 88%
  - New entities: „GSM BSS", „LTE TDD"
  - False pos: „Location Area LA"
- **Elements**
  - F1: 98%
  - False pos: „argon/hydrogen"
- **Cancer**
  - F1: 87%
  - New entities: „myeloma cell", „tumor tissue"
  - False pos: „cell mortality", „the test compound"
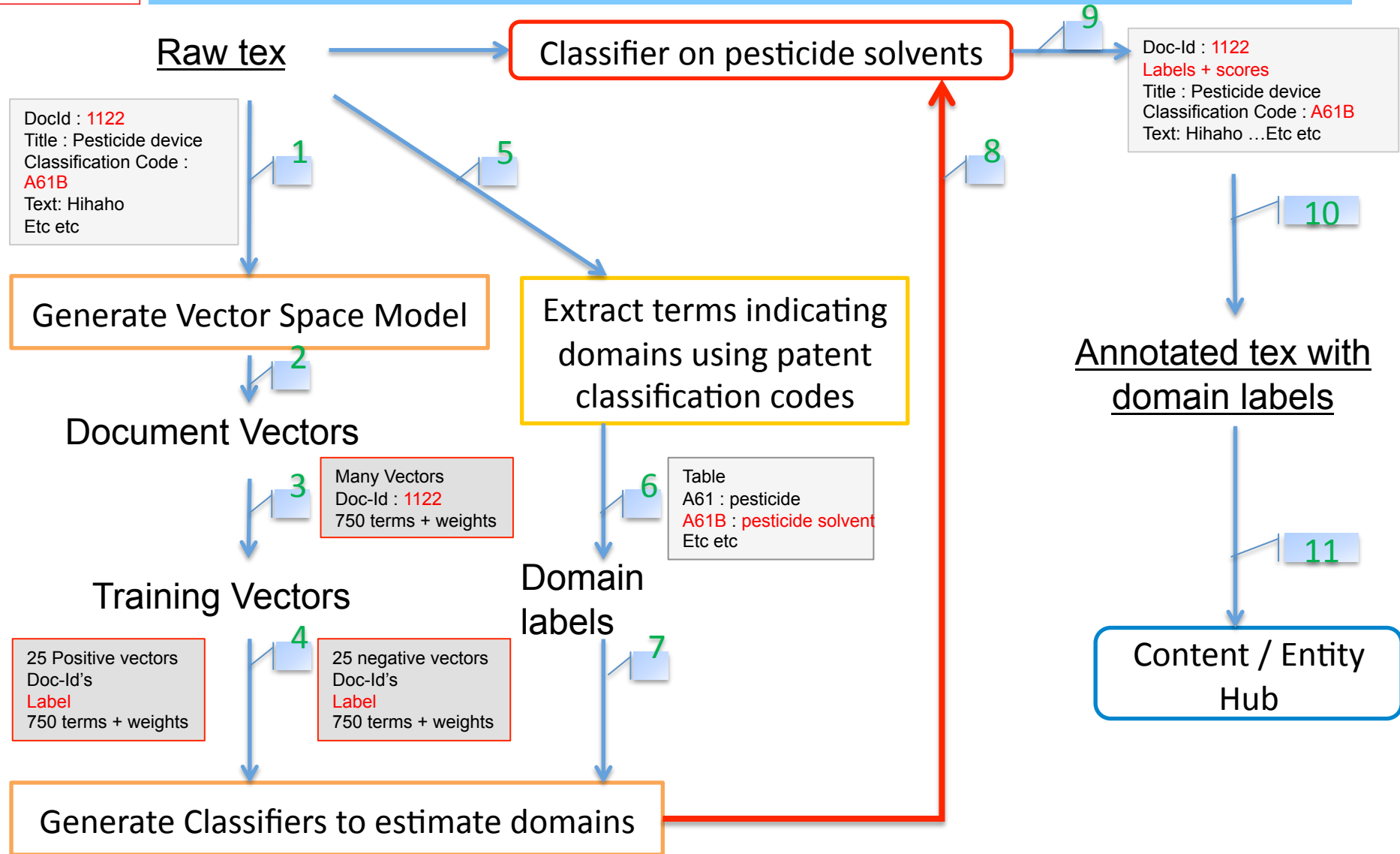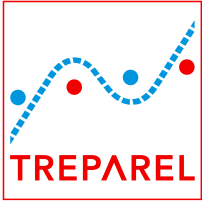
# Using NER in a GUI

# Software architecture

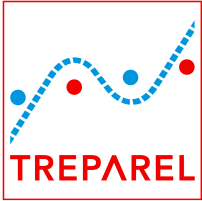- **Data sourcing**: retrieves data from data sources via *data integrator*

- **Storage**: raw (e.g. text) or processed data stored in database or triple store

- **Using ML** to enable learning from the crowd

- **Push & Pull**:
  - interface to consumers (web and mobile apps)
  - supporting quality and access control from portal

- **Portal**: business logic, storage of registered data sources, access control using web frontend
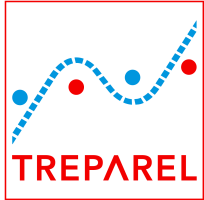
# Open Call for Users

- **Applications received** from Finland, Germany, Spain, Greece, Bulgaria, Hungary, Italy, UK, Belgium, China, Switzerland, France, Denmark, Portugal, Ireland and The Netherlands

- **Business areas covered**:
    - Bio-medical,
    - Pharma and biotech,
    - ICT/Telecommunications,
    - Digital media,
    - Renewable energies,
    - Education,
    - Innovation/consultancy services for SMEs.

- **Mix of profiles** from SMEs, Research, SME intermediaries (incubators, Science parks, Living Labs, etc), developers and more

- **Multiple areas** of research, development and innovation areas are covered.

# EC perspective

- **Data as a service**: scale economies of scale in management of data
- **Data pooling**: processes need timely aggregation and redistribution of diverse data but building own is redundant and prohibitive for SMEs
  - provide services on top of pool with high quality data
  - provide access to services on demand
- **Success criteria**:
  - Early provision of scalable basic Fusepool services
  - SME involvement and uptake of Fusepool services
  - Machine learning for data curation & user adaptation
- **Required steps**:
  - Stepwise integration of existing & new technologies
  - Early and ongoing feedback from end users

**The EC FP 7 program Fusepool is all about:**

- Building a platform with web enabled services for

  - Data pooling

  - Large scale text analysis

  - Large scale machine learning of user input

- Enable SME's with analytics for improving their innovation and competitive strengths using

  - SME's involvement and feedback to the Fusepool services

  - Machine learning for data curation & user adaptation

Welcome to visit Fusepool at:
www.fusepool.net