**FIZ** Karlsruhe

Leibniz Institute for Information Infrastructure

**SEMANTiCS**
**Karlsruhe 2019**

**ADVANCING SCIENCE**

**From specific problems to a generic solution:**
A scalable framework for analyzing Big Data of Patent Information
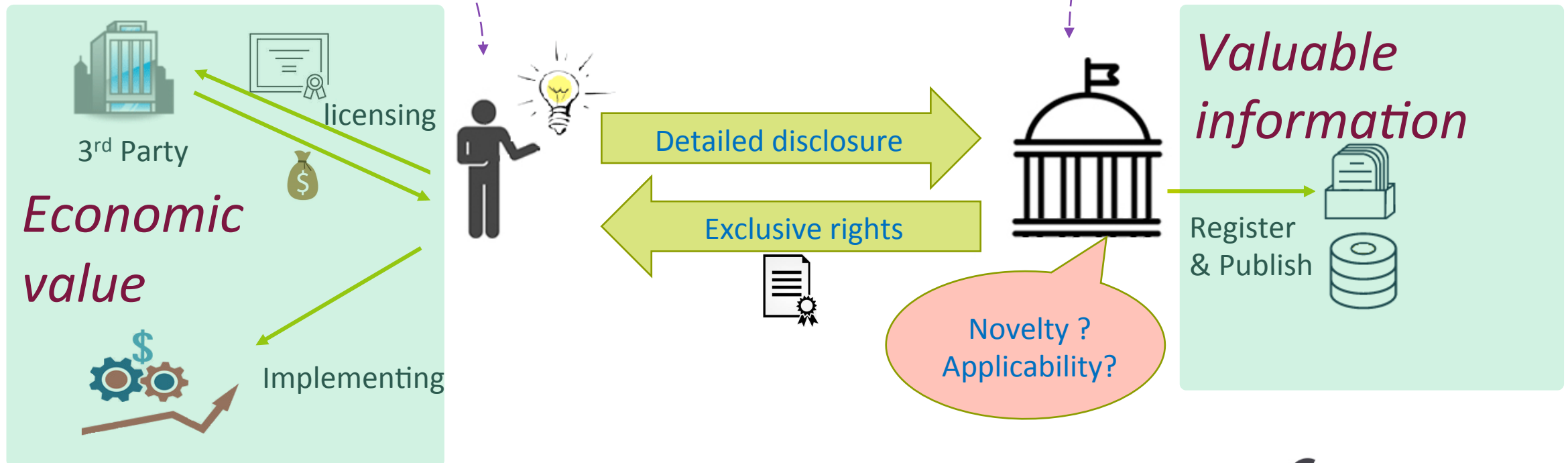
Ahmad Alrifai, September 09-12, 2019

## Outline

➢ Patent information – An overview

➢ Patent search and analysis

➢ Text and Data mining @FIZ

➢ Generic scalable framework

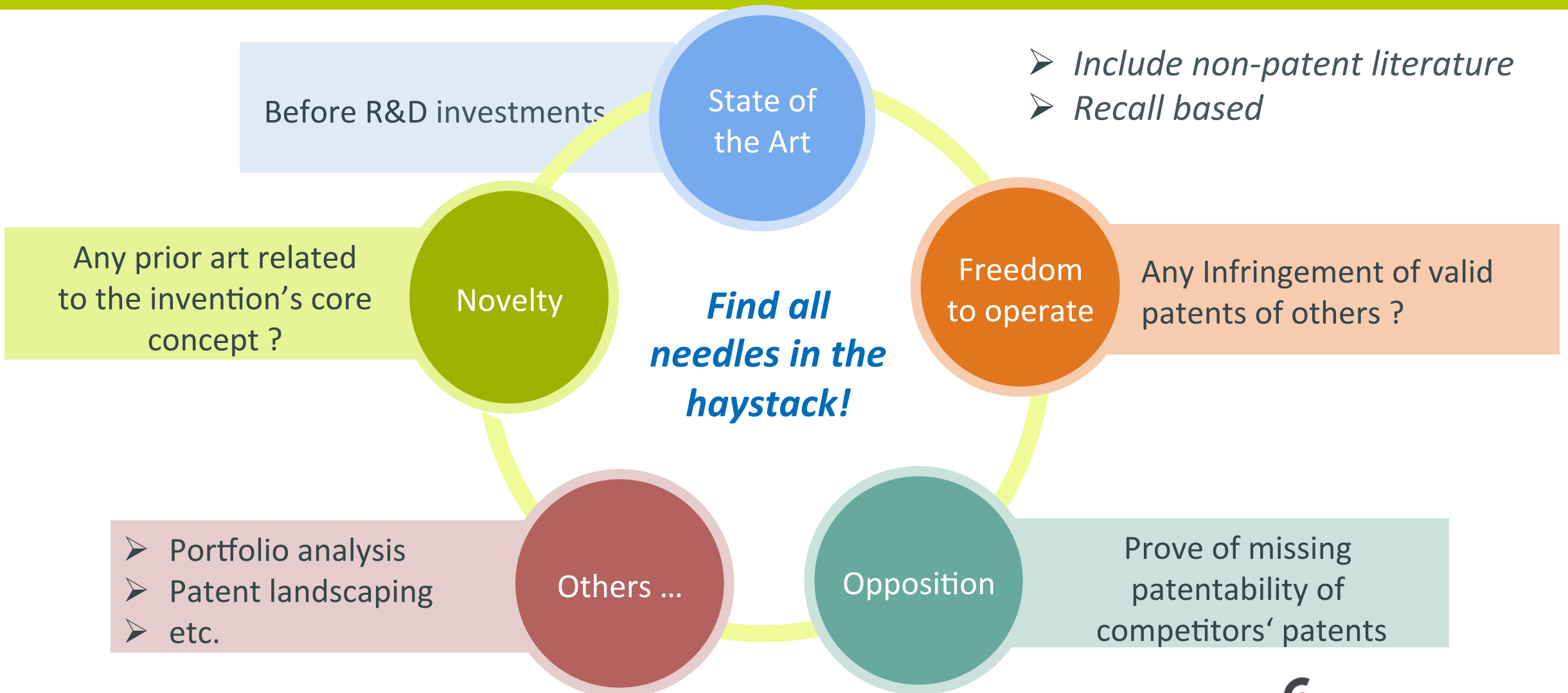➢ Ongoing and future work

➢ Conclusion

*"A patent (/ˈpætənt/ or /ˈpeɪtənt/) is a set of exclusive rights granted by a sovereign state to an inventor or assignee for a limited period of time in exchange for detailed public disclosure of an invention…."*

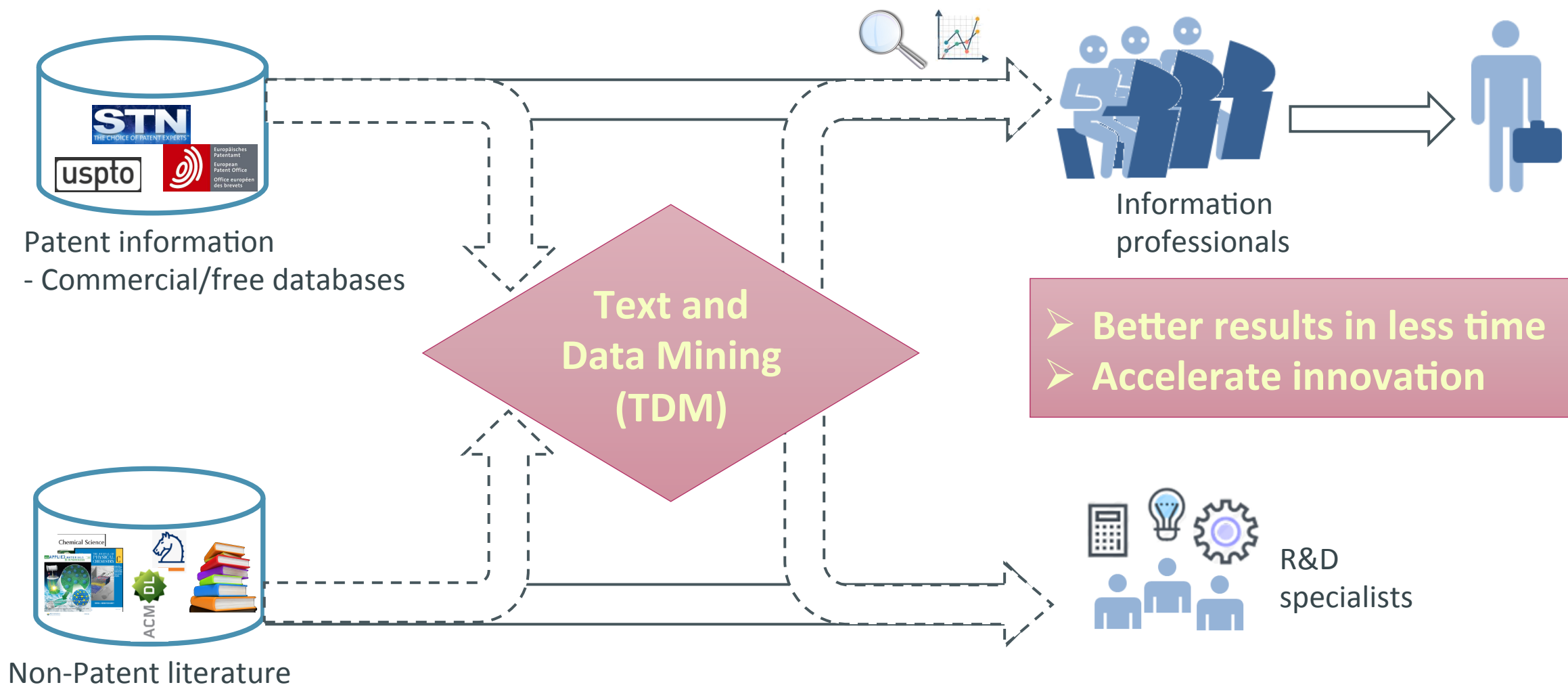*[source: wikipedia.org]*

A scalable framework for analyzing Big Data of Patent Information | Ahmad Alrifai | SEMANTiCS 2019

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

# Patent search and analysis – use cases

State of the Art

Novelty

Freedom to operate

Opposition

Others …

*Find all needles in the haystack!*

Before R&D investments

- Include non-patent literature
- Recall based

Any prior art related to the invention's core concept ?

Any Infringement of valid patents of others ?

- Portfolio analysis
- Patent landscaping
- etc.

Prove of missing patentability of competitors' patents

A scalable framework for analyzing Big Data of Patent Information | Ahmad Alrifai | SEMANTiCS 2019

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

# Patent and scientific information



Patent information
- Commercial/free databases

Text and
Data Mining
(TDM)

Non-Patent literature

Information
professionals

> **Better results in less time**
> **Accelerate innovation**

R&D
specialists

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

A typical chemistry patent

➢ *Meta-data: dates, names, classification,…*

➢ *Title*

➢ *Abstract*

➢ *Detailed description*

➢ *Claims: legal scope of the protection*
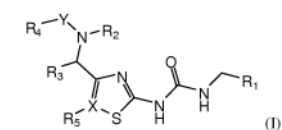
Multilingual
Patent Families

250 pages

System architecture for a powerful TDM

# Annotation and linking

➢ Patent chemical search is essential for chemical and pharmaceutical domains


Annotation, linking and indexing

➢ Trade names, chemical names and structures



**Diclofenac**

2-(2,6-Dichlorophenylamino)phenylacetic acid

Patent      Entity      Knowledge Base



- hydroxybenzoic acid
- chloride
- citrate
- sodium bisulfate
- EDTA
- ethyl alcohol
- propylene glycol
- sodium hydroxide
- hydrochloric acid
- citric acid
- lactic acid

**InchiKey:**
"KCXVZYZYPLLWCC-UHFFFAOYSA-N"

- IBM Patent System
- Atlas
- FDA SRS
- SureChEMBL
- PharmGKB
- Human Metabolome Database
- PubChem: Thomson Pharma
- PubChem
- Mcule
- NMRShiftDB
- ACToR

**PubChem**

| | |
|---|---|
| PubChem CID: | 6049 |
| Chemical Names: | EDTA; Edetic acid; Ethylenediaminetetraacetic acid; 60-00-4; Edathamil; Endrate |
| Molecular Formula: | $C_{10}H_{12}O_8CaN_2Na_2 \cdot 2H_2O$ or $C_{10}H_{16}N_2O_8$ or $((HOOCCH_2)_2NCH_2)_2$ |
| Molecular Weight: | 292.244 g/mol |
| InChI Key: | KCXVZYZYPLLWCC-UHFFFAOYSA-N |

**FIZ Karlsruhe**
Leibniz Institute for Information Infrastructure

# Structuring and analysis

"*The invention belongs to the field of* communication network*, in particular relates to a domain name of the* biometrics-based authentication system *and method.*"

Large-scale technical trends analysis



Segment the description part into predefined sections



- Abstract
- Claims
- Summary
- Examples Applicability
- Description
- Technical field
- Embodiments
- References
- Background art
- Appendex

➢ Targeted search

➢ Input to other tasks

A scalable framework for analyzing Big Data of Patent Information | Ahmad Alrifai | SEMANTiCS 2019

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

# Need for speed

- ➢ **Performance**

  - ▪ 1s processing time per document → ~ 2 months for a database

- ➢ **Special requirements**

  - ▪ High recall is required

  - ▪ Iterative evaluations and adaptations

- ➢ **Challenges**

  - ▪ NLP tools trained on different corpora → rule-based modifications

  - ▪ Lack of in-domain training data or golden corpora

  - ▪ Typographical errors from OCR and machine translation

A scalable framework for analyzing Big Data of Patent Information | Ahmad Alrifai | SEMANTiCS 2019

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

# It's all about scaling

Scale!

| | | |
|---|---|---|
| Keywords Extraction | | Scalable KE |
| Claim Structuring | | Scalable CS |
| Numerical Analyzer | | Scalable NA |
| Chemical Annotator | | Scalable CA |
| Desc. Segmentation | | Scalable DS |
| Analytics Service | | Scalable AS |
| ... | | |

hadoop YARN

APACHE Spark

kubernetes

APACHE PHOENIX

FIZ Karlsruhe
Leibniz Institute for Information Infrastructure

# Generic scalable framework

Keywords Extraction

Claim Structuring

Numerical Analyzer

Chemical Annotator

Desc. Segmentation

Analytics Service

...

Overhead

Dependencies

Timeout

Memory overflow

Utilization

Data locality

Scalable KE

Scalable CS

Scalable NA

Generic Scalable
Service

Scalable CA

Scalable DS

Scalable AS

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

# Setup, Execute, Aggregate, Cleanup

## Service.jar

```
if (numR != 0)
    return -1;

List<Map<String, Object>> allRecordIdsMap = new ArrayList<>();
List<Map<String, Object>> queryForList = tdmDAO.getJdbcTemplate().queryForList(SQLstate);

if (queryForList.isEmpty()) {
    String error = "Can't initialize ResultsetExtractor: SQL-Statement returned 0 answers!";
    throw new Exception(error);
}
allRecordIdsMap.addAll(queryForList);

primaryKeyListIterator = allRecordIdsMap.iterator();

if (primaryKeyListIterator == null) {
    String message ="Could not create key list for database iteration process with sql: ";
    LOG.error(message);
}

//write the text file
int numFiles = allRecordIdsMap.size()/lines+1;
FSDataOutputStream[] out = new FSDataOutputStream[numFiles];

int counter = 0;
int fileIndex = counter / lines;
out[fileIndex] = fsMaster.create(new Path(hdfsInput+"result" + Integer.toString(fileIndex) + ".txt"));

while(primaryKeyListIterator.hasNext()){

    if ((counter / lines)!=fileIndex){
        out[fileIndex].close();
        fileIndex = counter / lines;
        out[fileIndex] = fsMaster.create(new Path(hdfsInput+"result" + Integer.toString(fileIndex) + ".t)
    }
    Map<String, Object> row = primaryKeyListIterator.next();
    @SuppressWarnings("unchecked")
    Entry<String, Object> idEntry= (Entry<String, Object>) row.entrySet().toArray()[0];

    out[fileIndex].write(((String)idEntry.getValue().toString()+"\n").getBytes());

    row.clear();
    counter++;
}

primaryKeyListIterator = allRecordIdsMap.iterator();

if (primaryKeyListIterator == null) {
    String message ="Could not create key list for database iteration process with sql: ";
    LOG.error(message);
}

//write the text file
int numFiles = allRecordIdsMap.size()/lines+1;
FSDataOutputStream[] out = new FSDataOutputStream[numFiles];

int counter = 0;
int fileIndex = counter / lines;
out[fileIndex] = fsMaster.create(new Path(hdfsInput+"result" + Integer.toString(fileIndex) + ".txt"));

while(primaryKeyListIterator.hasNext()){

    if ((counter / lines)!=fileIndex){
        out[fileIndex].close();
        fileIndex = counter / lines;
        out[fileIndex] = fsMaster.create(new Path(hdfsInput+"result" + Integer.toString(fileIndex) + ".txt
    }
    Map<String, Object> row = primaryKeyListIterator.next();
    @SuppressWarnings("unchecked")
    Entry<String, Object> idEntry= (Entry<String, Object>) row.entrySet().toArray()[0];

    out[fileIndex].write(((String)idEntry.getValue().toString()+"\n").getBytes());

    row.clear();
    counter++;
}
out[fileIndex].close();
/**
 * Submit the job
 */
```

## Generic Scalable Service

**Setup** — Initialization code

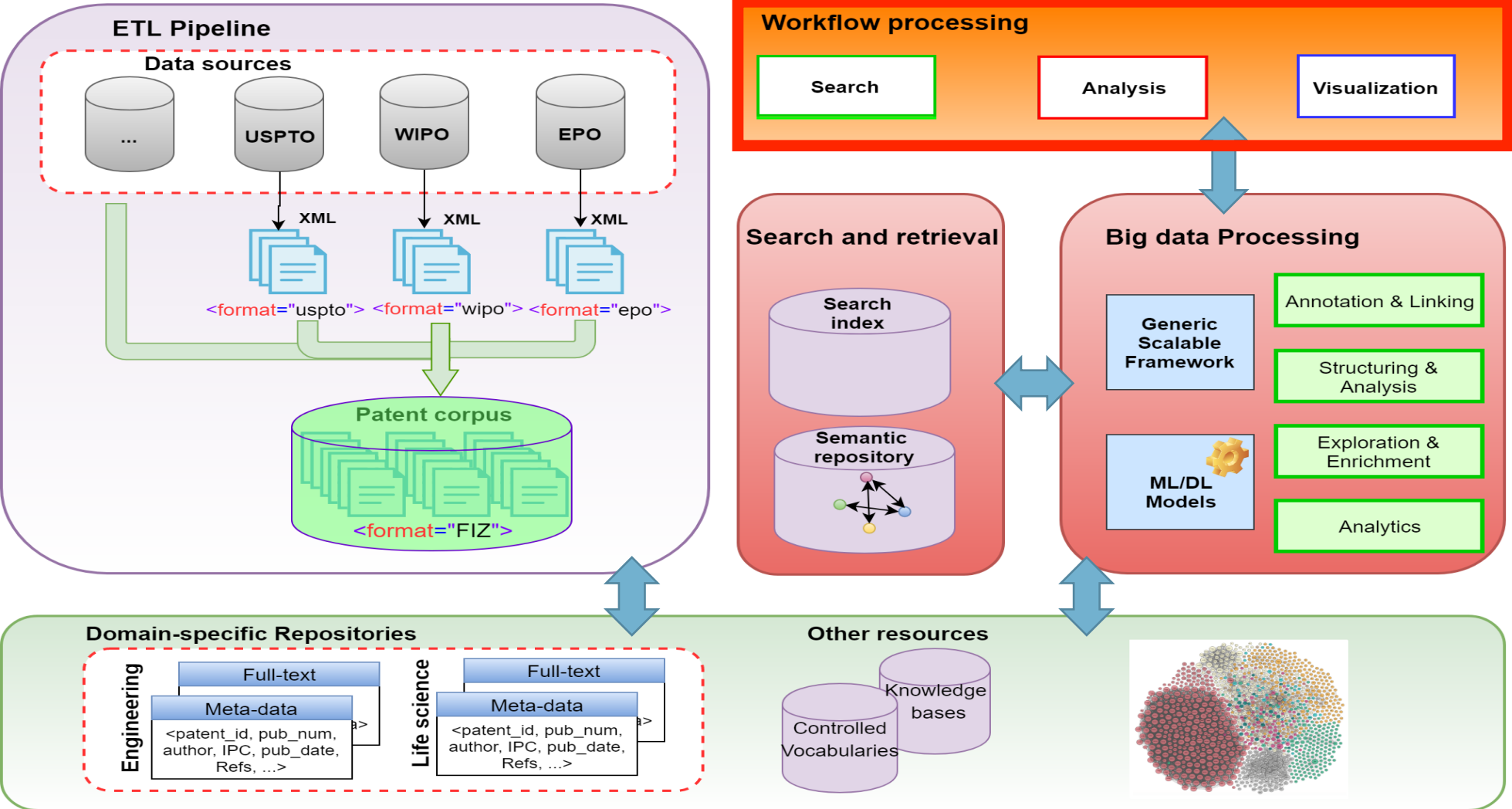**Execute** — Parallelization at document level

**Aggregate** — Intra-document analysis

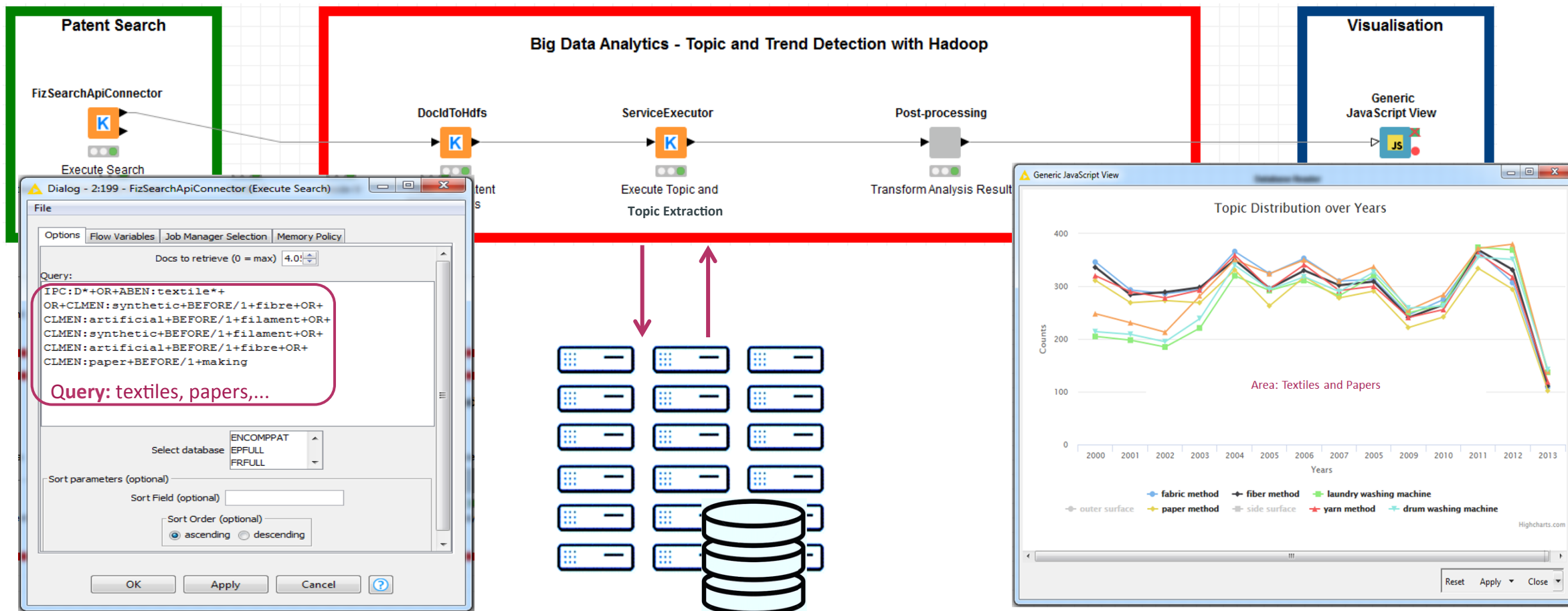**Cleanup** — Post processing, filtering, ranking, sorting, etc.

**FIZ** Karlsruhe
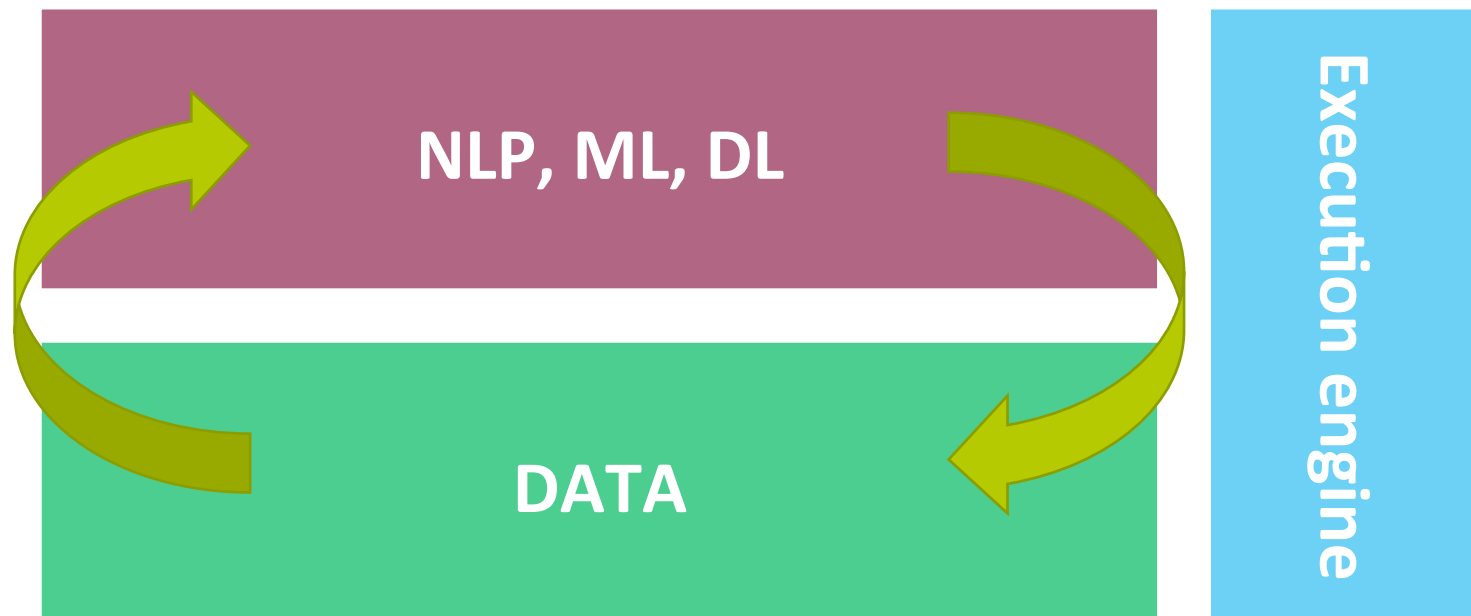Leibniz Institute for Information Infrastructure

Analysis and recognition of technology trends

# Whats's next  - next generic Generation

Interoperability of NLP and machine learning steps in interactive workflows

➢ Data level: standardized formats

➢ Execution engine: abstractions and optimizations – Ontology for TDM

A scalable framework for analyzing Big Data of Patent Information | Ahmad Alrifai | SEMANTiCS 2019

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

# Conclusion

- Patents form a unique and valuable knowledge source (also beyond IP domain)

- Need for the most advanced techniques to generate synergies and added values

- Patent analytics domain is catching up with adapting Machine learning and semantic technologies

- Scalable infrastructure and generic frameworks advance semantic technologies and boost the performance of TDM applications

- Considering and linking with domain-specific KBs to maximize the potentials of patent data mining

- Integrity, compatibility and interoperability of NLP and Deep learning for better comparability and reusability

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

# THANK YOU!
## Questions ?

## Visit us!



**Contact**

Ahmad Alrifai
Senior Researcher

+49-7247 808-545
Ahmad.Alrifai@fiz-karlsruhe.de

**Text and Data mining**
https://www.fiz-karlsruhe.de/de/forschung/text-und-data-mining-tdm

**FIZ** Karlsruhe
Leibniz Institute for Information Infrastructure

Zertifikat seit 2016
audit berufundfamilie

Leibniz
Association