# AR-Science: From Raw Text Data to Structured Semantic Representation

Text Search Interface

© Artificial Researcher, 2022

# Graph Search Interface

3

# Requirements for Scientific Search

- Meta-data
  - <u>Bibliographic data</u>
  - References and citations
  - Resource sharing (code/data)
- Full-text
  - <u>Identification of domain-specific terminology and concepts</u>
  - <u>Identification category scientific discipline</u>
  - <u>Identification of related terms in discipline context</u>
  - Identification  text paragraph
  - The essence of the article in terms of claim, method, result, conclusion

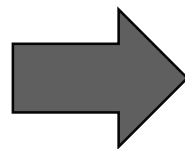# From Text to Structured Semantic Representation of Scientific Publications

International Patent Classification (IPC)

# Lexical Relations

Hyponym relations:
- part of
- kind of
- member of
- type of
- domain specific relation



Hypernym, **a broad term**,
  e.g. farm animals a hypernym of pig, cow, sheep, horse

Hyponym, **a narrow term**
  e.g. horse is a hyponym of farm animals

# Integrating Different Resources

# Assembly Line for Creating Indices and Ontologies



Documents are processed through several text analysis modules:

1. Text Normalization and harmonization (including PDF converter)
2. NLP and relation identification
3. Artificial Researcher NLP-Toolkit - Technical Terms
4. Packaged as enhanced XML/JSON
5. Software AR Ontology Services and AR Passage Retrieval
6. End user applications rest APIs, desktop script, GUI

# Visual display and direct access to original resource

Access to the Meta data such as taxonomy (IPC), and context semantic similarities

## Visualization of relations



| GENERAL | alkyl |
|---|---|
| SPECIFIC | hexyl |
| CLASSIFICATION | H04R, |
| SOURCE DATASET | EPO |
| SIMILARITY | PatBERT: 0.1342575, |
| | SciBERT: 0.83727044, |
| SOURCE PDF | https://data.epo.org/publication-server/pdf-document?cc=EP&pn=1701583&ki=A1&pd=2006-09-13 |
| LANGUAGE | en |
| VERSION | 1 |
| SOURCE TEXT | Examples of alkyl include hexyl , heptyl , octyl , isooctyl , nonyl , isononyl , decyl , isodecyl , undecyl , dodecyl , tridecyl , isotridecyl , tetradecyl , isotetradecyl , pentadecyl , isopentadecyl , hexadecyl , isohexadecyl , heptadecyl , isoheptadecyl , ocadecyl , nonadecyl , isononadecyl , eicocyl , isoeicocyl , henicosyl , isohenicosyl , docosyl , isodocosyl , tricosyl , isotricosyl , tetracosyl and isotetracosyl . |
| TEXT LOCATION | 1701583_DESCR_0 |
| DETECTOR | LSP_PATTERN_DETECTOR |

Direct access to the source and original sentence

## Query formulation Boolean

(ester OR hexyl OR acrylate OR methyl OR alkyl OR methacrylate OR alkyls OR diol OR radical) OR ("decyl group"~4 OR "alkyl methacrylate"~4 OR "alkyl acetate"~4 OR "alkenyl group"~4 OR "substituted alkyl"~4 OR "aliphatic alcohol"~4 OR "ethylene glycol"~4 OR "methacrylic acid"~4 OR "crotonic acid"~4 OR "alkyl radical"~4 OR "suitable alkyl acrylate"~5 OR "normal alkyl radical"~5 OR "lower alkyls"~4 OR "c alkyl"~4 OR "selenite ester"~4 OR "branched chain"~4 OR "alkoxyalkyl group"~4 OR "ester linkage"~4)

9

# Extracting semantic relation from natural language text

Input

```
1  {
2      "api_key":"4d39678c5cbd428db3312a66d3dae13f",
3      "text":"vinyl ethers methyl ketone such as vinyl hexyl ketone and methyl isopropenyl ketone ; N-vinyl.",
4      "classification": [
5          "Class 1",
6          "Class 2"
7      ],
8      "request_similarity": [
9          "SciBERT"
10     ],
11     "source_dataset": "some dataset",
12     "source_file": "some file",
13     "text_id": "in some file"
14 }
```

Output

```
"results": [
    {
        "CLASSIFICATION": [
            [
                "Class 1",
                "Class 2"
            ]
        ],
        "DETECTOR": "LSP_PATTERN_DETECTOR",
        "GENERAL": "vinyl_ether_methyl_ketone",
        "LANGUAGE": "en",
        "SEARCH_HELPER": [
            "hexyl",
            "vinyl",
            "methyl",
            "ether",
            "ketone"
        ],
        "SIMILARITY": {
            "SciBERT": 0.9645886
        },
        "SOURCE_DATASET": "some dataset",
        "SOURCE_FILE": "some file",
        "SOURCE_TEXT": "vinyl ethers methyl ketone such as vinyl hexyl ketone and methyl isopropenyl ketone ; N - vinyl .",
        "SPECIFIC": "vinyl_hexyl_ketone",
        "TERM_HELPER": {
            "term0": "hexyl",
            "term1": "vinyl",
            "term2": "methyl",
            "term3": "ether",
            "term4": "ketone"
        },
        "TEXT_LOCATION": "in some file",
        "VERSION": 14
    },
```

Automatic classification according to taxonomies such as IPC/CPC, MeSH, PLoS, WIPO Scientific Subject fields

# Better than Distributional Semantics

Using semantically enriched data can retrieve up to 60% more related concepts than traditional pre-trained contextual models (a.k.a. Neural Network-based methods, BERT-family).

### Artificial Researcher's text mining technology

**Technical terms**

Related concept

**brake pedal:**
vehicle operating pedal,
conventional hydraulic brake system
pedal devices
position brake actuating member
brake actuating member
hydraulically-assisted rack pinion steering
gear
brake operating member
conventional braking system
pair pedals

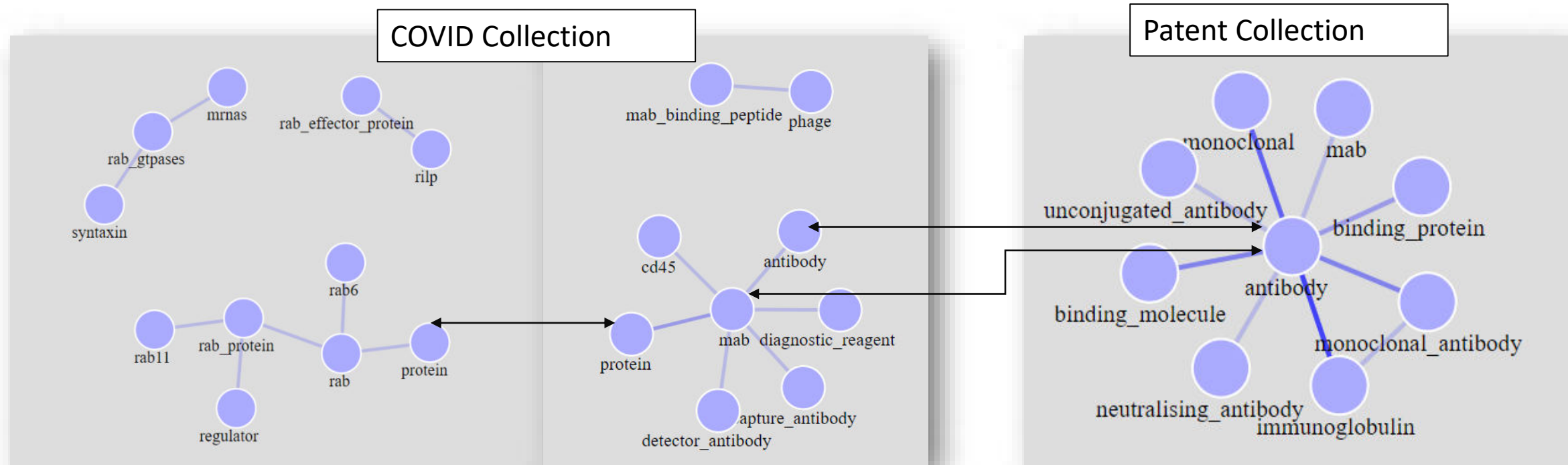**accelerator pedal**
case pedal device
pedal device

### Extraction using pre-trained contextual models only

| brake pedal | PatBERT | SciBERT |
|---|---|---|
| conventional hydraulic brake system | 0.69 | 0.89 |
| hydraulically-assisted rack pinion steering gear | 0.49 | 0.80 |
| conventional braking system | 0.66 | 0.84 |

# Connecting Different Collections and Domains

# Verified Use Cases

Artificial Researcher in Open Access

**Intellectual property Industry**

Cloud-based architecture
- Scientific publication
- Patent data

*Pre-Seed 2019-2020*

Artificial Researcher in Science

**Next generation scientific search tool**

Onsite architecture
- Open access
- Internal resources

*FemPower IKT 2019-2022*

EUROPEAN OPEN SCIENCE CLOUD

**iFAIR**
Identifying datasets in scientific publications. (2021)

**AR-ONTO-COVID**
A Knowledge-Based Resource for Covid-19 (2020)

# Pain points

- Data quality
- Data quality
- Data quality

– descriptions missing
– no statistics about the data
– wrong content (OCRs, classification symbols)
– …

# Thank you for your attention

artificialresearcher.com

florina.piroi@artificialresearcher.com

linda.andersson@artificialresearcher.com