# NLP, AI and LLM WITH LEGALXML

Michele Corazza – Researcher at CIRSIFD-ALMA AI, Department of Legal Studies, University of Bologna, Italy

**ECAI 2025 –** University of Bologna
Day 1 - October 25, 2025

# Large Language Models and NLP

- The recent advancements in Natural Language Processing (NLP) have led to the creation of Large Language Models that can perform a multitude of tasks

- These are very powerful tools, however there is always a non-zero risk that the generated text contains errors (hallucinations)

- This is particularly severe for the legal domain, which has some characteristics that complicate any AI/NLP/LLM application.
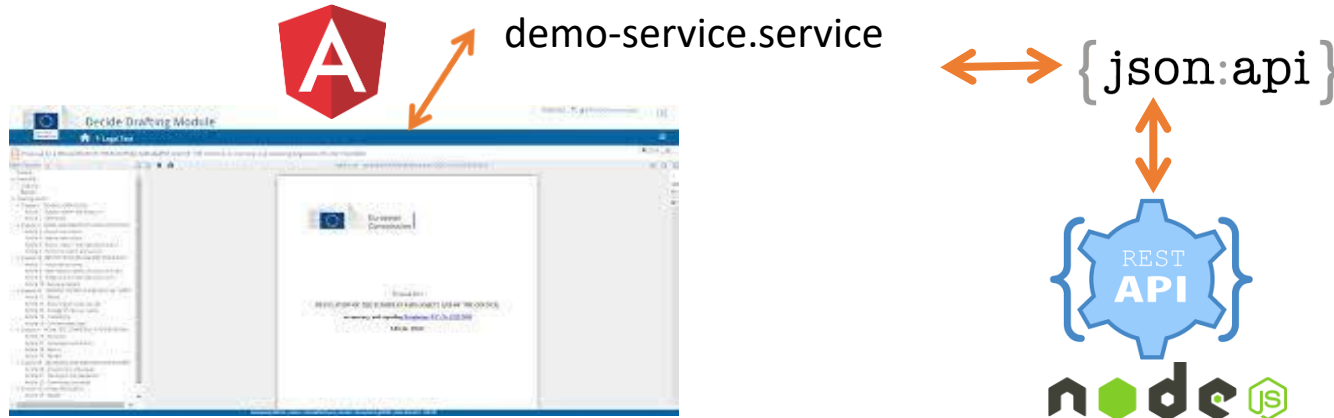
**USE-CASES**

# Problem

1. **Assist** the legislative drafting offices and the members of parliament in retrieving relevant legislative sources related to the bill that they want to draft (heterogeneous corpora) – **legal ex-ante analysis**

2. Monitor **policies** over time – **ex-post analysis**

3. Detect inconsistencies in legislative bills or **unconstitutionalities** with the Constitutional Court

4. **Generate** new legislative definitions according to the existing corpora with **Agentic AI**

# #1 Relevant documents with a lack of information in legislative heterogenous corpora
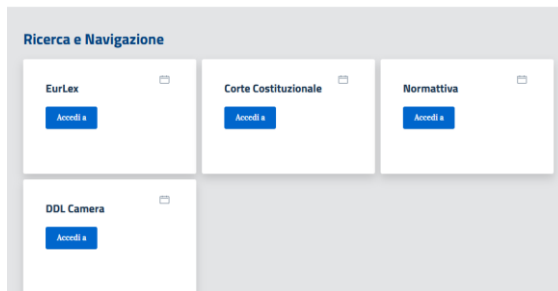
1. Complete drafting of normative references
2. Retrieve definitions
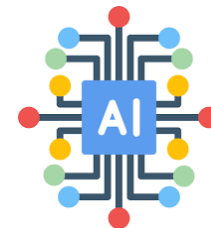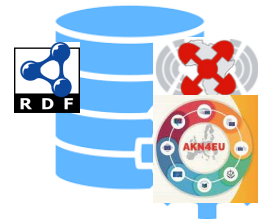3. Extract relevant legal documents

demo-service.service

$\{json:api\}$

REST API

node.js

Vue.js

Lucine, eXist: **AKN**
**CELLAR: RDF**
DBSQL

AKN4EU

# 1.1 Normative references suggestion

Given an **incorrect** or **partial text**, to get the list of relevant references, syntactically correct, in force, refer to documents that have a similar topic.

**Reference**: Regulation 122

**EuroVoc**:energy use

1)Regulation (EU) 2015/1222
Regulation (EU) 2011/1227
Regulation (EU) 2014/1227

**query**

**results**

References Suggestion Module

**VDB**

SQL

eXistdb

**AKN-XML documents temporal information**

# 1.1 Normative references suggestion

The module works as follows:

1. The year and number of the document are extracted from the query using heuristics
2. First, a query is used to find documents with the same number and/or year in the query (Directive 31 → Directive 2010/31/EU)
3. The model applies **Levensthein distance** to find documents that have a distance of at most 1 for both the year and document number
4. Averaged FastText embeddings are then used to encode the EuroVoc terms in the query and compare them with documents, obtaining a **ranking**

# 1.2 Legal definitions suggestion

Use prior existing definitions to suggest relevant ones, consolidated and updated, to the user using the topic of the bill (EuroVoc).

- Hydrogen definition in the Energy Bill updated today
- Hydrogen definition in the Food Bill in force in 2026

**query**

**Term**: "Hydrogen"

**EuroVoc**: "Energy"

**AKN-XML documents
definitions
EUROVOC
temporal information**

Definitions
Suggestion
Module

**Definition:**"Hydrogen Sensor" means...

**Long title**: Commission Regulation (EU) No 406/2010 of 26 April 2010 ...
**Partition**: art_1__list_1__point_1
**Version date:** 2010-04-26
Eurovoc: Technical standard, Pollution control, Motor vehicle, ...
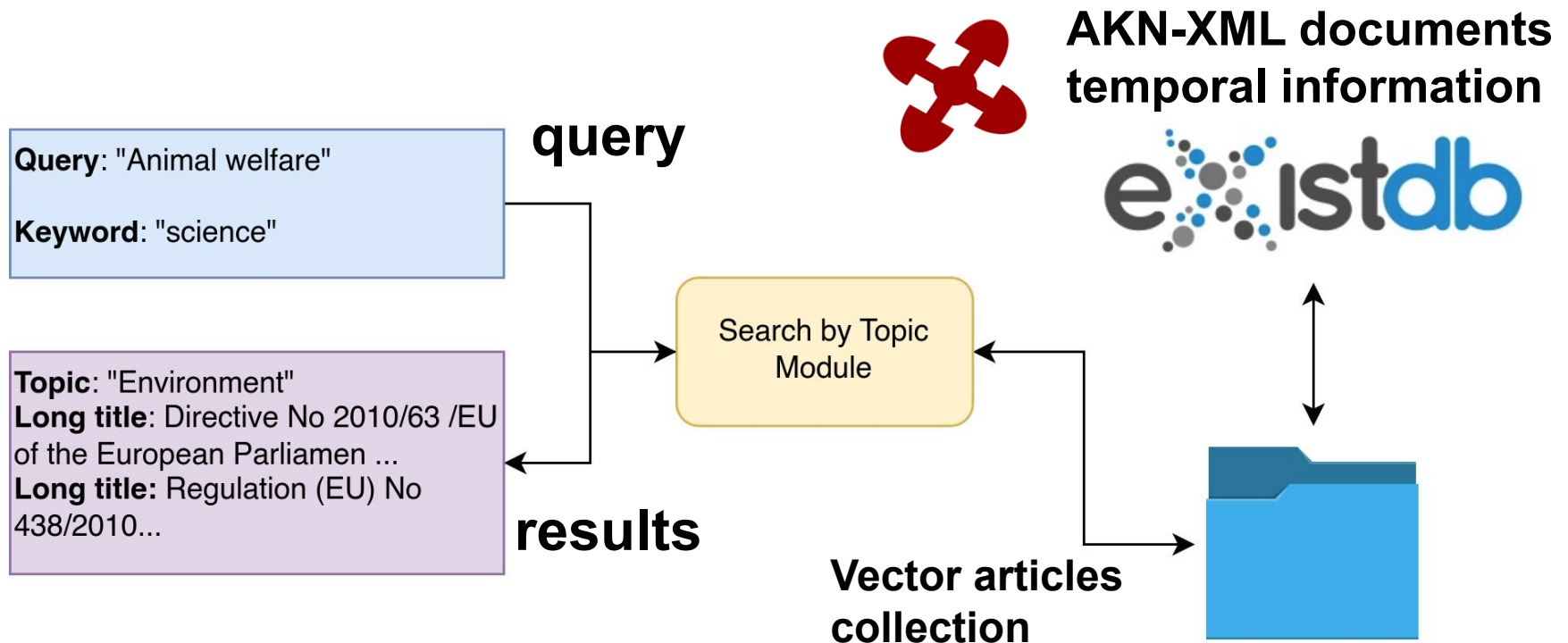
**results**

Powered by *lucene*

# 1.2 Legal definitions suggestion

The module uses the Similarity class from Lucene (a refined version of the Vector Space Model (VSM) based on TF-IDF weights) to compare a query and existing definitions:
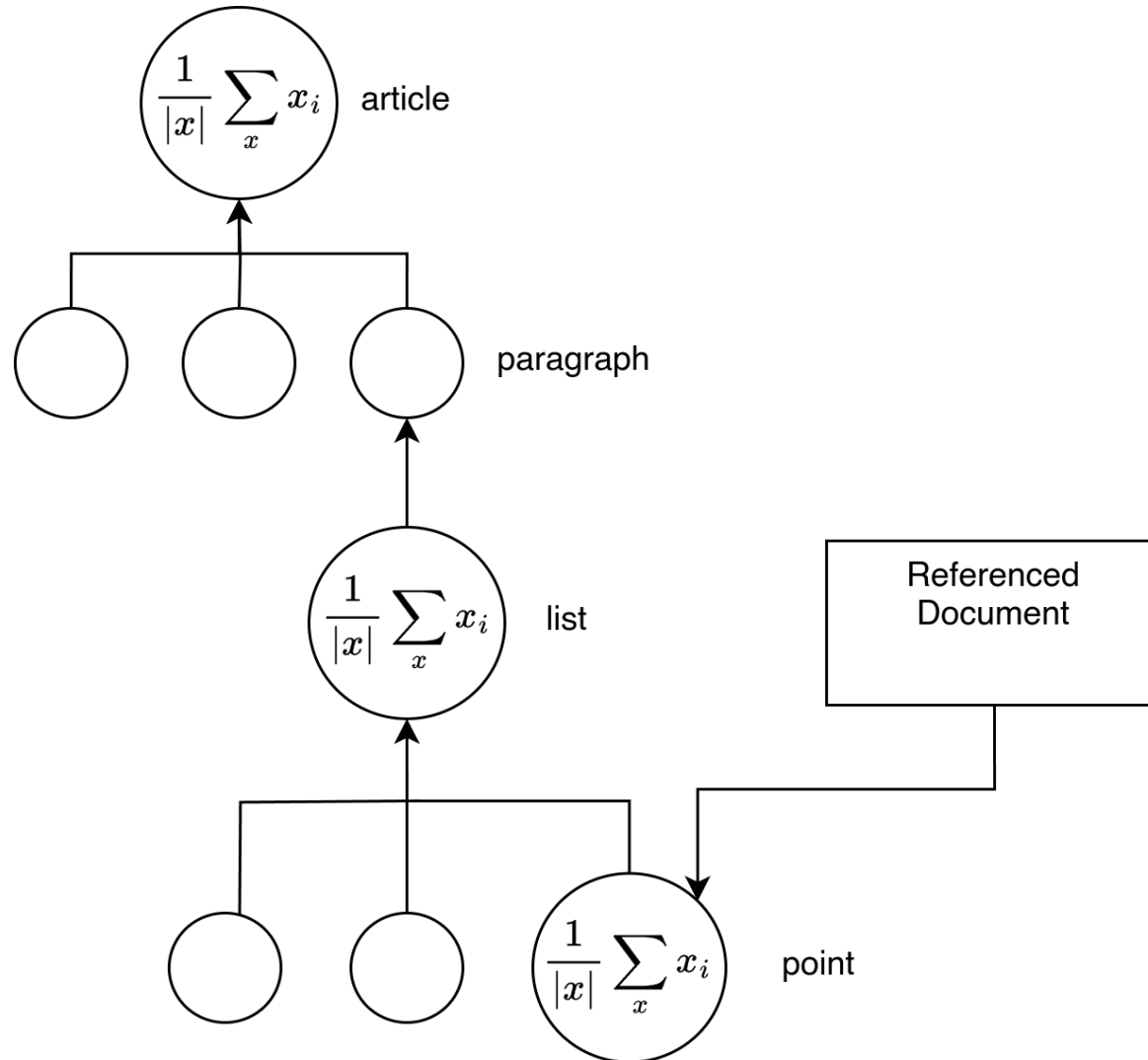
1.  The query EuroVoc and the EuroVoc from the document are compared
2.  The query is compared with both the definendum and definiens of the definition (<<"definendum" means "definiens")
3.  The results are filtered using user-specified manual constraints to obtain relevant results

# 1.3 Search by topic

- Given a keyword of EuroVoc and the topic of bill, it returns a list of the 10 relevant documents divided in two relevant topics top level EuroVoc terms.

**Query**: "Animal welfare"

**Keyword**: "science"

**query**

**AKN-XML documents temporal information**

Search by Topic Module

**Topic**: "Environment"
**Long title**: Directive No 2010/63 /EU of the European Parliamen ...
**Long title:** Regulation (EU) No 438/2010...

**results**

**Vector articles collection**

# 1.4 Embedding method

# #1 - Evaluation

TABLE III

EVALUATION FOR THE NORMATIVE REFERENCES SUGGESTIONS.

| Number of Eurovoc | Top 1 Accuracy | Top 5 Accuracy | Top 10 Accuracy |
|---|---|---|---|
| 1 | 0.20 | 0.40 | 0.80 |
| 2 | 0.30 | 0.70 | 1.00 |
| 3 | 0.30 | 0.70 | 1.00 |

TABLE II

ACCURACY VALUES FOR DEFINITION SUGGESTION AND DOCUMENT CLUSTERING. FOR THE DEFINITION SUGGESTION TASK, WE SHOW BOTH THE RESULTS, INCLUDING NO OUTPUT, AS WELL AS THOSE WHERE THESE CASES HAVE BEEN EXCLUDED.
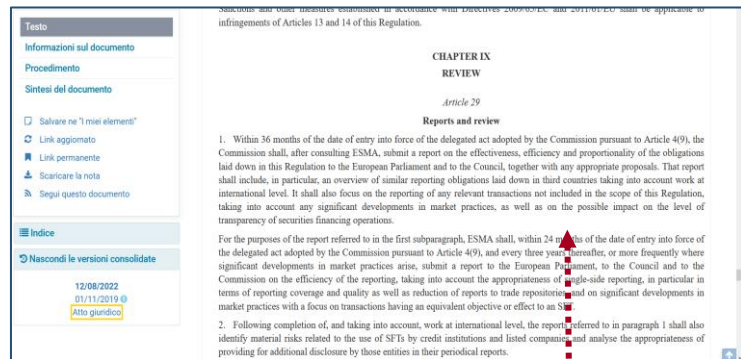
| Task | Accuracy (with "blank" outputs) | Accuracy (without "blank" outputs) |
|---|---|---|
| Definition suggestion | 0.62 | 0.72 |
| Document clustering | 0.52 | - |

# #2 SORTIS: Monitoring and Measuring the Policy



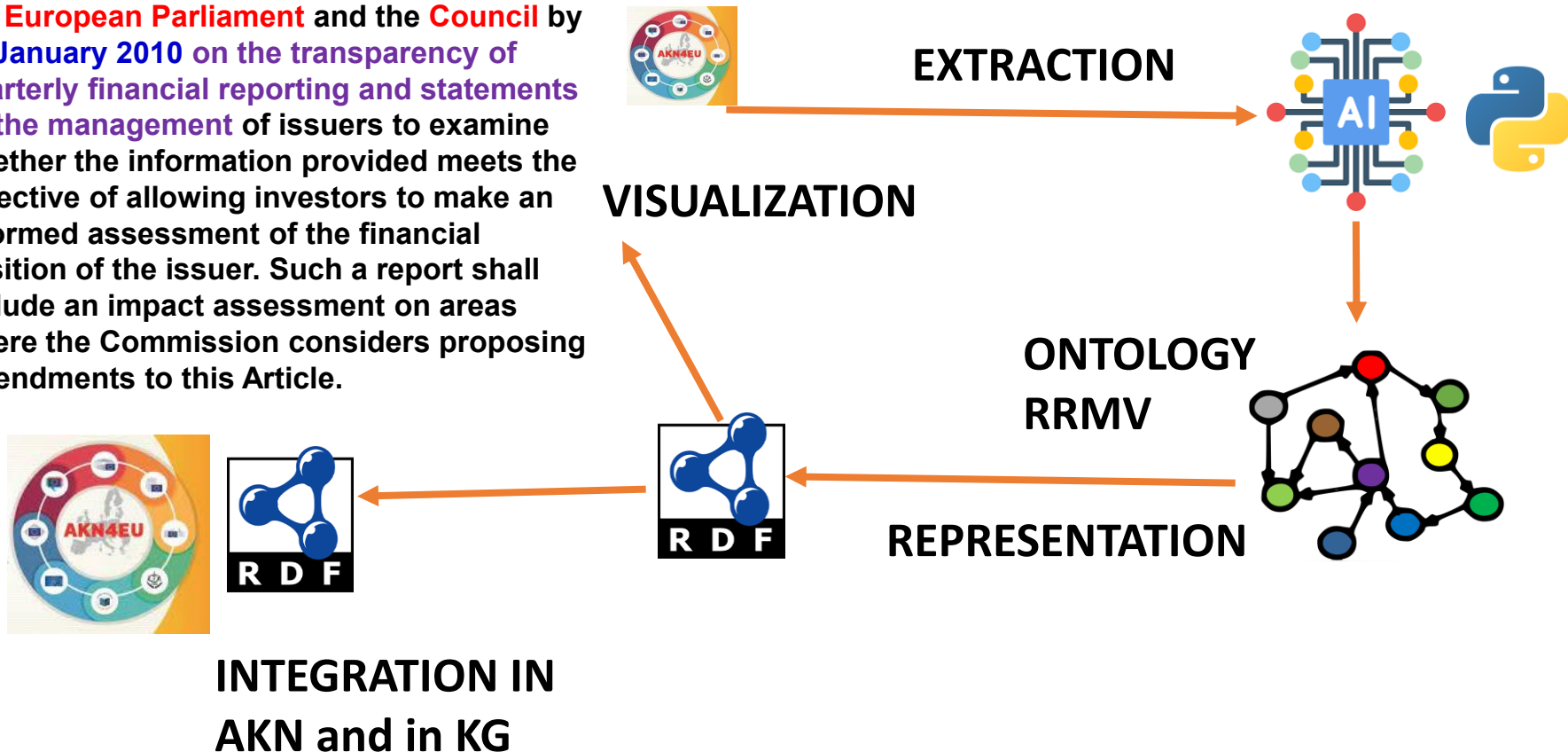**Measuring**
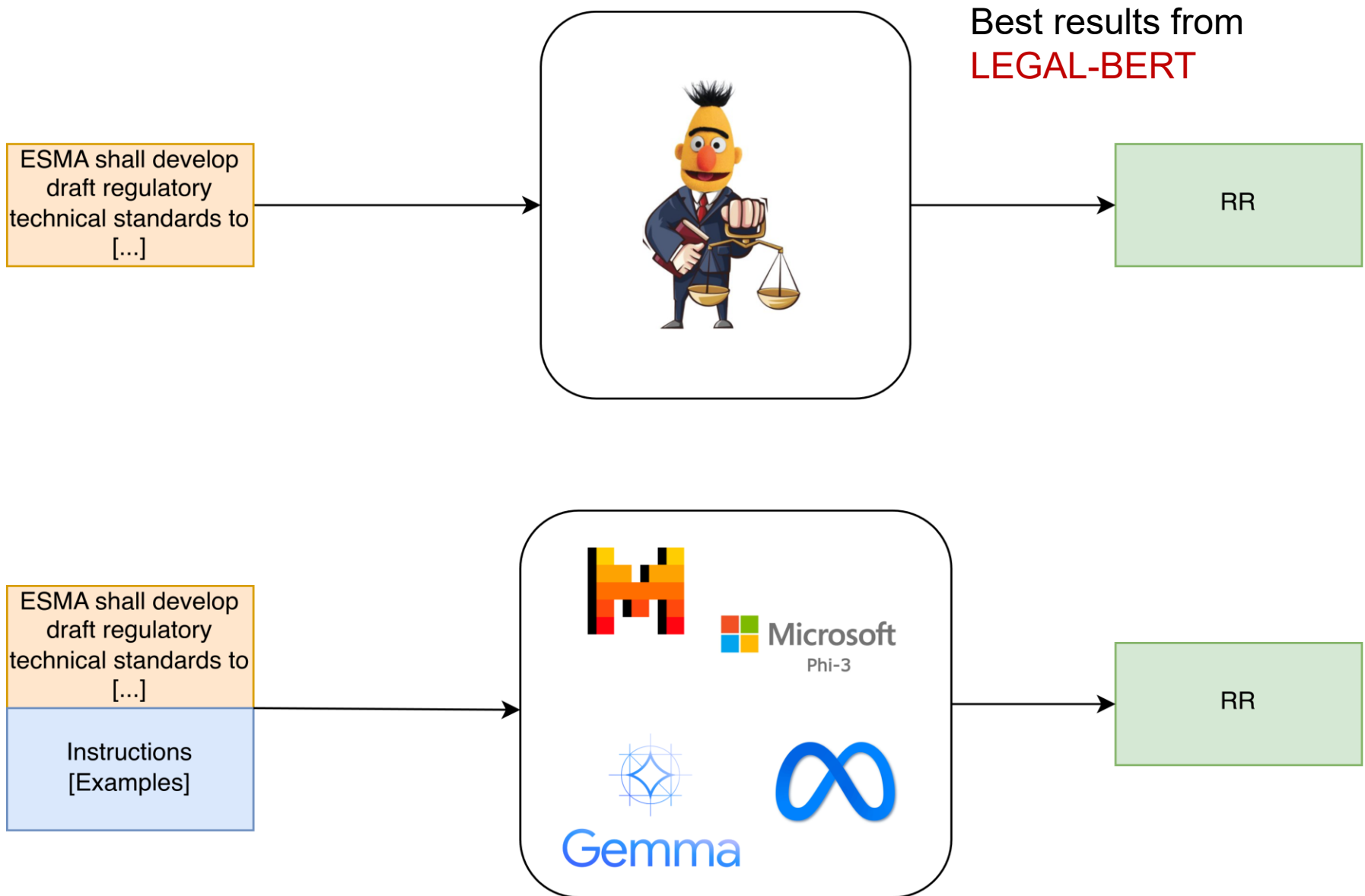
**Reporting Requirement**

**RRMV-RDF**

**Monitoring**

# #2 - Extraction and Representation of «Reporting Requirements»

3.   The **Commission** shall provide a **report** to the **European Parliament** and the **Council** by **20 January 2010** on the transparency of quarterly financial reporting and statements by the management of issuers to examine whether the information provided meets the objective of allowing investors to make an informed assessment of the financial position of the issuer. Such a report shall include an impact assessment on areas where the Commission considers proposing amendments to this Article.

**EXTRACTION**

**VISUALIZATION**

**ONTOLOGY RRMV**

**REPRESENTATION**

**INTEGRATION IN AKN and in KG**

# #2 - RR binary classification

ESMA shall develop draft regulatory technical standards to [...]

Best results from LEGAL-BERT

RR

ESMA shall develop draft regulatory technical standards to [...]

Instructions [Examples]

Microsoft Phi-3

Gemma

RR
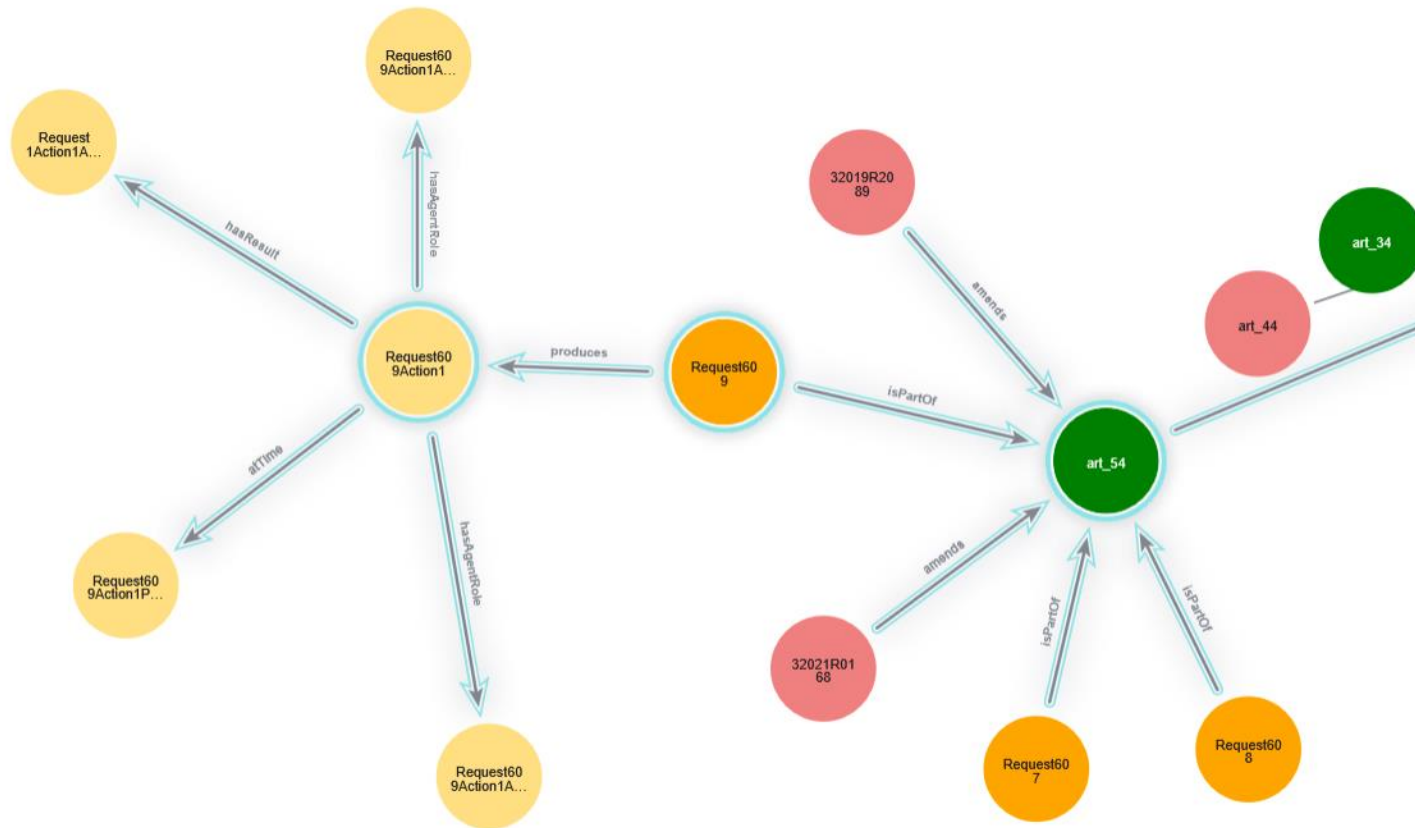
# #2 - Entity Extraction

Best result from LEGAL-BERT, followed closely by GEMMA 10-shot and Llama3.3 10-shot RDF

| ESMA | | Addresser |
| shall | | - |
| develop | | Action |
| [...] | | ... |

PeriodOfTime better with GEMMA

| ESMA shall develop draft [...] | | Microsoft Phi-3 |
| Instructions [Examples] | | Gemma |

Addresser: ESMA
Action: develop
[...]

# #2 - Art. 54 is modified, it contains RR – Neo4J



Reporting Requirements Ontology for European Legislation
Monica Palmirani[1][0000-0002-8557-8084], Andrea Giovanni Nuzzolese[2][0000-0003-2928-9496], and Generoso Longo[1][https://orcid.org/0009-0003-2687-5884]

EGOVIS 2025

# **Conclusions**

- In order to mitigate the risks of the application of LLMs to the legal domain, we can apply a hybrid approach

- The Akoma Ntoso standard allows us to leverage its rich information and to use them with sophisticated NLP/AI models (structure of the documents, definitions, references, jurisdiction, etc)

- By enriching statistical models with symbolic information, we reduce the risk of hallucinations and ensure that the models consider the relevant (i.e. in force) documents while also considering the temporal aspects of the law.