

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

## What Can Machine Learning Do for the Public Procurement?

### **This is the author's manuscript**

*Original Citation:*

*Availability:*

This version is available <http://hdl.handle.net/2318/1888839> since 2023-01-31T10:07:14Z

*Terms of use:*

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

# What can Machine Learning do for public procurement?

I. Fatima<sup>PHD\_CS</sup>, F. Gorgerino<sup>PHD\_AL</sup>, **R. Meo**<sup>ASS\_CS</sup>, R. Nai<sup>PHD\_CS</sup>, P. Pasteris<sup>TEC\_CS</sup>,  
**G. M. Racca**<sup>FULL\_AL</sup>, E. Sulis<sup>RTDB\_CS</sup>



Computer Science Department - University of Turin  
Management Department – Administrative Law

# Agenda

- Introduction
- Case study
- Methodology to reach the goals
- Results
- Conclusion and future work



A decorative network diagram in the top-left corner, consisting of various sized circles (nodes) connected by thin lines (edges). Some nodes are solid grey, while others are hollow with a dashed border. The network is dense and irregular.

1.

# Introduction

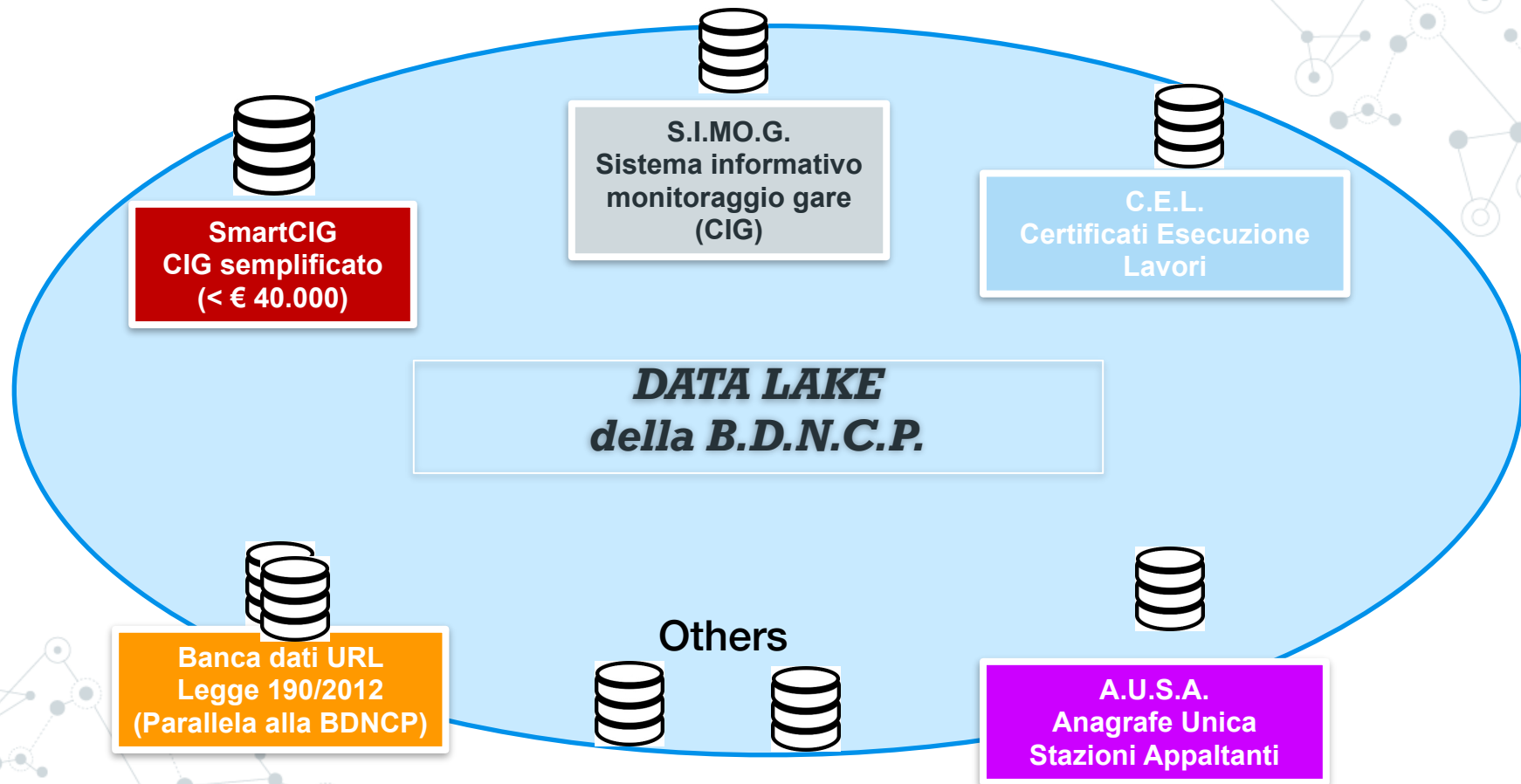
# Abstract

- We present the systematic work we conducted on the data about **public procurement** in Italy.
- The **goal** is to clean and integrate various public and open information sources and extract valuable information for the public sector and the companies interested in awarding a contract with the Public Administration.
- Included in the data analysis is the **Regional Administrative Justice** that receives recourses from the involved actors.
- This information coming from recourses is potentially useful for revealing some of the anomalies related to the incorrect behaviour of the partners.
- The obtained results can also make lighter the administrative judges' workload.

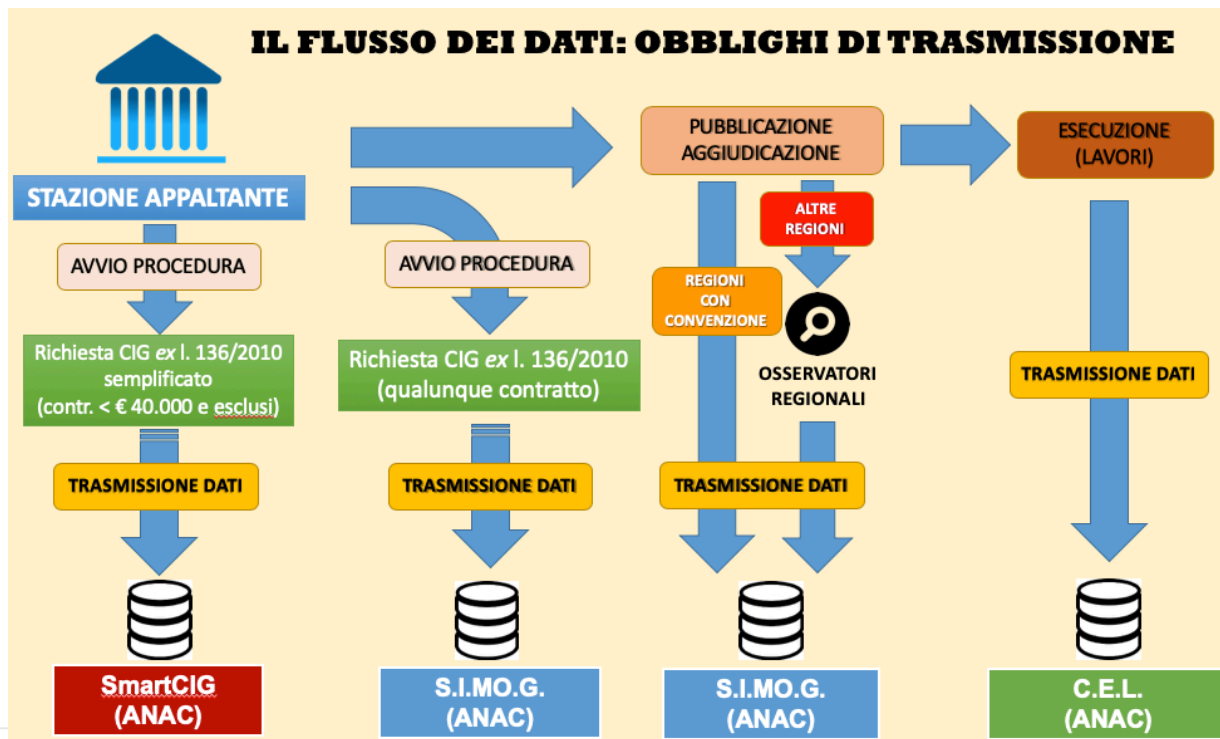
# The research project with ANAC

- ◎ *ANAC and University of Torino have stipulated a research project in 2019 for data analysis on public contracts*
- ◎ *BDNCP is a **data lake** with:*
  - *Sistema Informativo Monitoraggio Gare (SIMOG) with data on tenders*
  - *SMART CIG with data on contracts whose amount is under threshold (<40 K euro)*
  - *Anagrafe Unica delle Stazioni Appaltanti (AUSA)*
  - *CEL, with data on completion of works (*Certificati di Esecuzione Lavori – CEL*)*

# Banca Dati Nazionale dei Contratti Pubblici (BDNCP)



# The data flow





# Goals of the research project



**BANCA DATI NAZIONALE  
DEI CONTRATTI  
PUBBLICI**

DATA ANALYSIS

ENFORCEMENT OF **KNOWLEGDE  
BASE**  
FOR THE  
PUBLIC DECISIONS

DATA ELABORATION

PROMOTION AND DISSEMINATION  
OF **BEST PRACTICES**

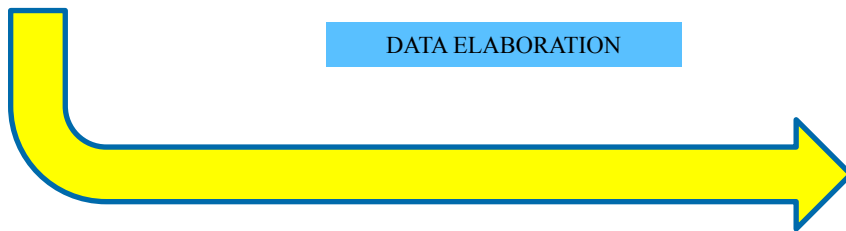
**EFFICIENCY AND INTEGRITY OF  
PUBBLIC ADMINISTRATION**

DATA COMPARISON

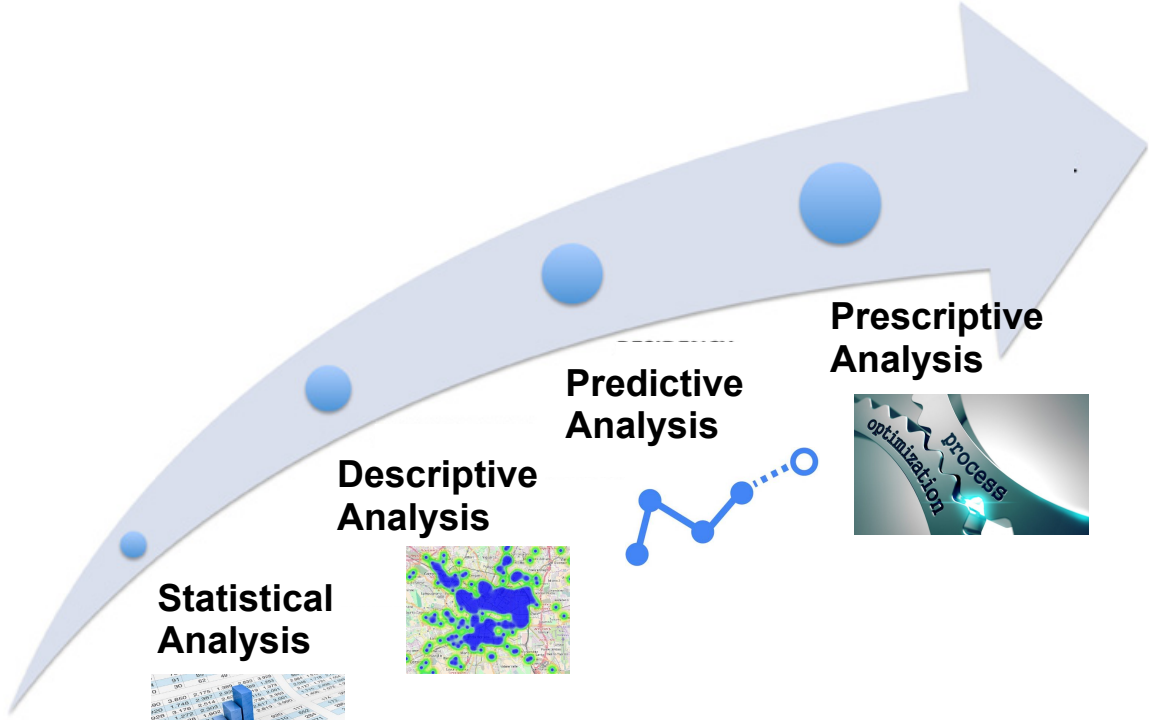
**QUALITY AND TRANSPARENCY  
IMPROVEMENT  
OF ADMINISTRATIVE SERVICES**

KPI DEFINITION AND  
USAGE


PREVENTION OF **CORRUPTION**





# The data analysis path



# Research questions

 **Research question 1 (RQ1):** how can we automatically **extract information from different legal archives**, in order to identify the entities involved in a public procurement?

 **Research question 2 (RQ2):** is it possible to set up an experiment to **predict the possible recourses** to administrative courts through the features of a public procurement?

 **Research question 3 (RQ3):** is it possible to provide services (**recommender systems**) to help the administrative operators (PA or economic)?

A decorative network diagram in the top-left corner, consisting of various sized circles (nodes) connected by thin lines (edges). Some nodes are solid grey, while others are hollow with a grey outline. The network is dense and irregular, extending from the top-left towards the center of the page.

# 2.

## Case study

# ANAC Open Data

- The **National Anti-Corruption Authority**, abbreviated to **ANAC**, is an independent Italian administrative authority whose task is to prevent corruption in the Italian public administration, implement transparency and supervise public contracts.
- **ANAC collects data on calls for procurement** from the public contract authority and provides a catalog of Open Data (in CSV, JSON and OCDS format) describing **public procurement**, **contract authority** (public administration), and **contractors** (economic operators winning the procurements).



<https://dati.anticorruzione.it/opendata>

# ANAC Open Data


2007 - 2022



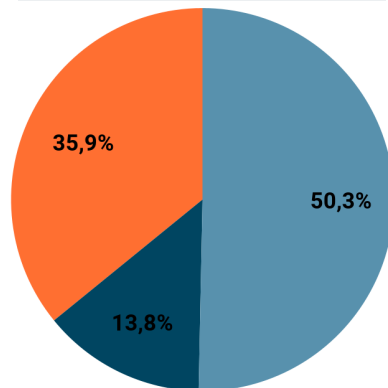
**ANAC  
Procurement**



7.2 million procurement  
~450.000 per year

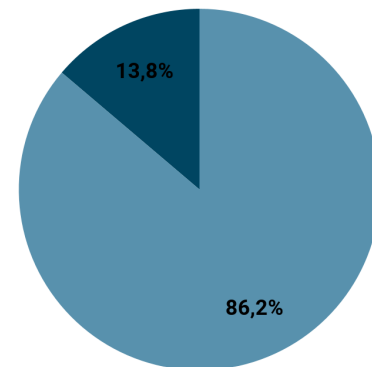
 **Identifier:** CIG (shared key)

Procurement object



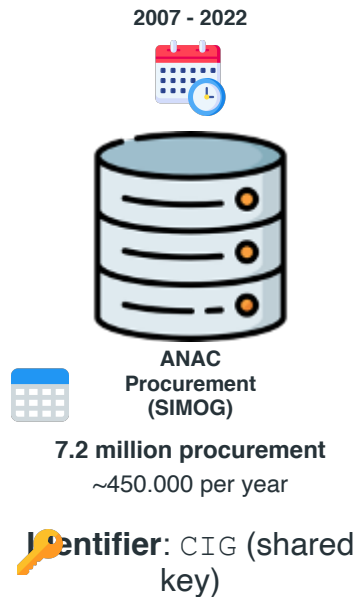
● Good/Supplies ● Public works ● Services

Procurement type



● Ordinary ● Special

# ANAC Open Data



**Contractors (PA): 42,393**



- Municipalities
- Regional governments
- Universities
- Hospitals
- etc.



**Awards:** 1,635,609  
(~22.71% of Procurement)



- **Private companies** specialized in goods/supplies, public works, services

**Economic operators: 265,039**

# Law 190/2012

- In compliance with Law 190/2012, **ANAC also collects the list of private companies participating in public procurement.**
- The data are transmitted by the individual administrations to ANAC via URL containing XML file.

 <https://dati.anticorruzione.it/#/190>



# Italian Administrative Justice

- On the other side, the **Italian Administrative Justice (IAJ)** collects the judges' sentences related to the public procurement appeals.
- Currently, about 80,360 judgments from the regional administrative courts (TAR) are available on the website.
- Every IAJ sentence is a textual file in html (60,284 files available), doc/docx (20,076 files available), or in pdf format (12 files available).



<https://www.giustizia-amministrativa.it>

A decorative network diagram in the top-left corner, consisting of various sized nodes (some solid grey, some hollow white) connected by thin grey lines, forming a complex web structure.

# 3.

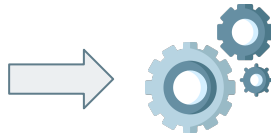
## Methodology

# Data gathering

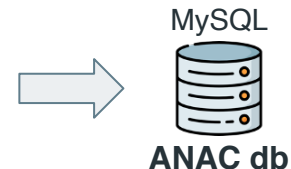
- Regarding **ANAC dataset**, we imported the Open Data of procurement from CSV format into a MySQL database table.
- Regarding **Law 190/2012**, we downloaded and imported the XMLs containing data of participants to procurement into a MySQL database table, with the fiscal code and name of each private company, the CIGs of procurement with its participation.
- We obtained the **IAJ judgments** via web scraping.
  - Since these are text files in html, doc/docx, and pdf format, they were indexed using tools specialized in Big Data processing.
  - According to November, 2022 figures on DB-Engines Ranking of Search Engines, **Elasticsearch** is the most popular search engine software used in Industry ([https://db-engines.com/en/ranking\\_trend/search+engine](https://db-engines.com/en/ranking_trend/search+engine)).

# Data gathering

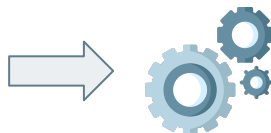
## ANAC Open Data



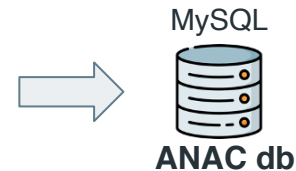
Procurement,  
Contractors,  
Economic operators,  
Awards



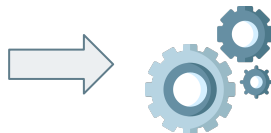
## Law 190/2012



Participants



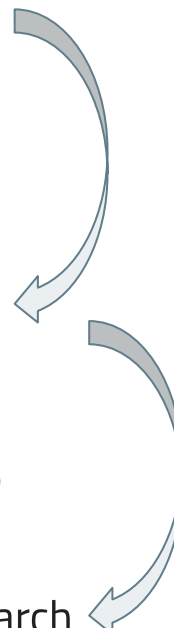
## IAJ files



NDJSON

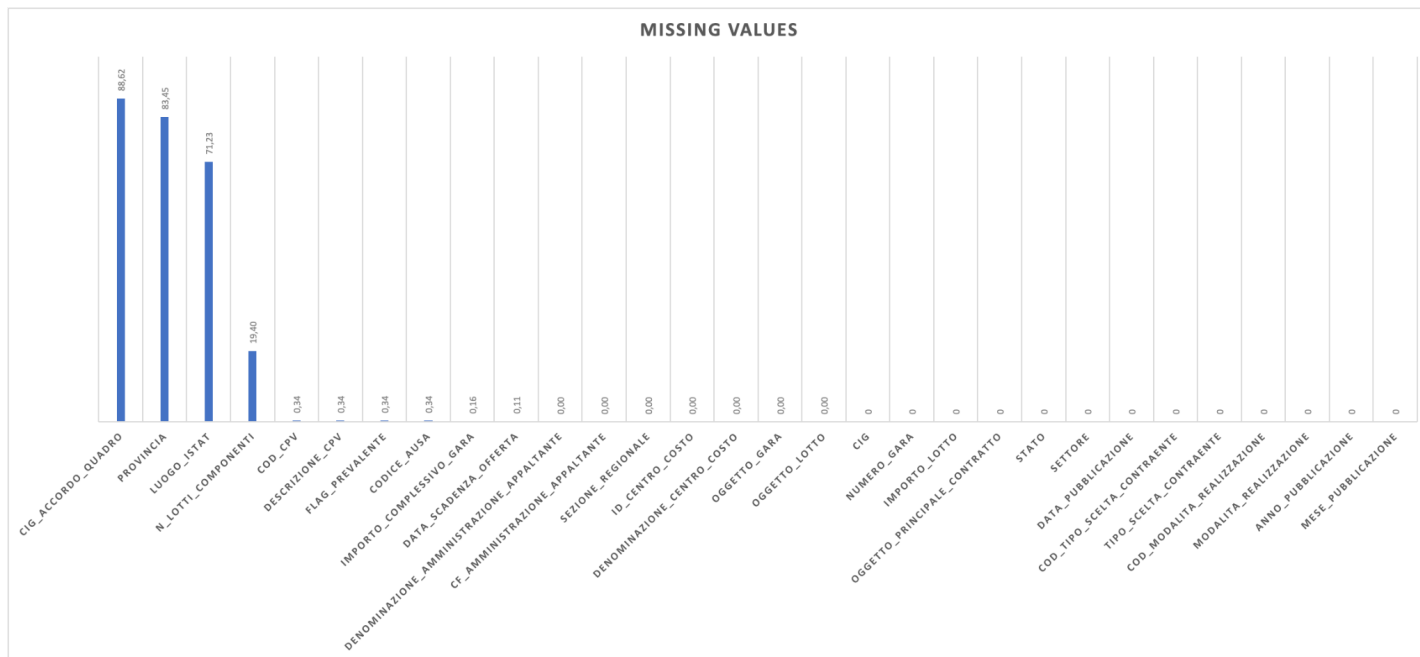


Judges'  
sentences



# Data quality problems

- ANAC Open Data: missing values



# Data quality problems

- Law 190/2012: missing XMLs

2015 2016 2017 2018 2019 2020 2021 2022

Udine

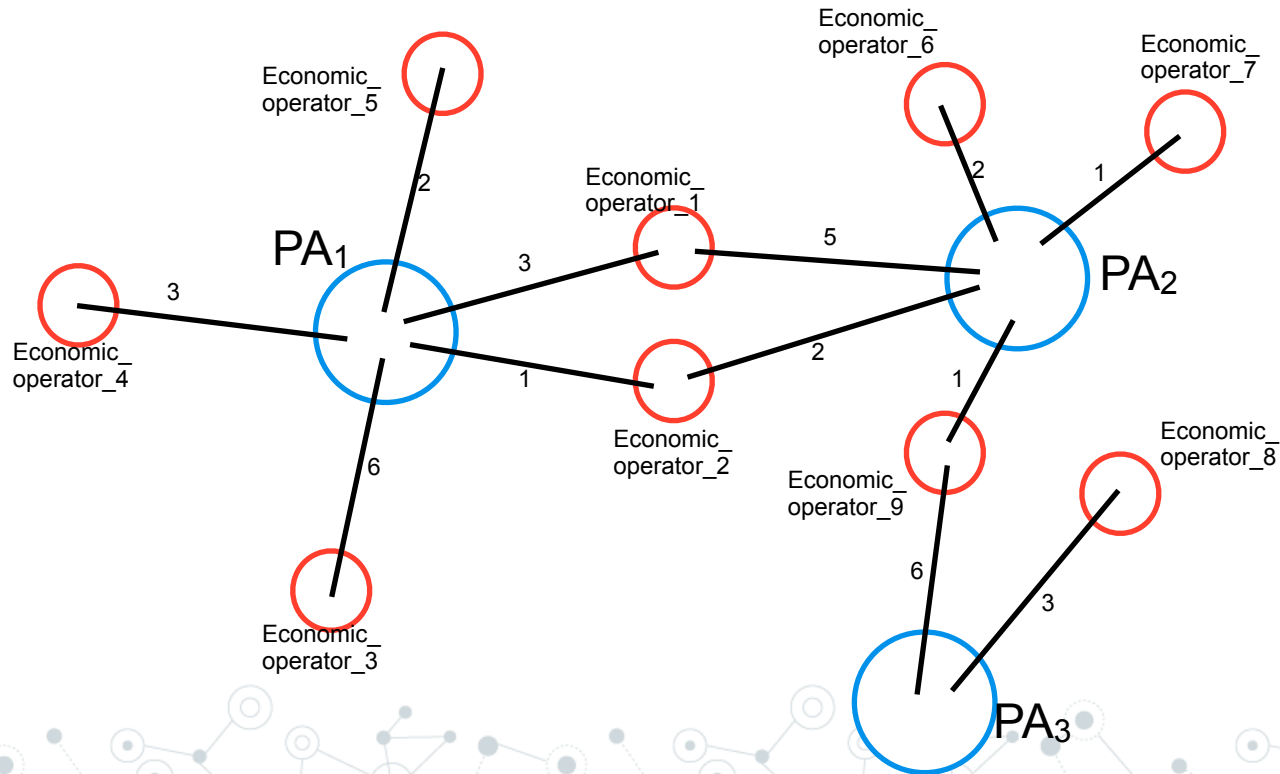
| CF Amministr... | Denominazione Amministrazione                              | Identificativo messaggio PEC                   | URL | Esito acce... | Data accesso |
|-----------------|--|--|-----|---------------|--------------|
| 94127260308     | ISTITUTO COMPRESIVO III DI UDINE                           | opcec296.20220127093747.444221.314.1.24@pe...  |     |               | 16/02/2022 1 |
| 80004840304     | ARCHIVIO DI STATO DI UDINE                                 | E3317E68-AFAC-E1F5-8A39-C2582BA2B9F8@te...     |     |               | 16/02/2022 1 |
| 94127270307     | ISTITUTO COMPRESIVO I DI UDINE                             | opcec296.20220124084121.16790.736.1.23@pec...  |     |               | 16/02/2022 1 |
| 00469890305     | COMUNE DI PAVIA DI UDINE                                   | 1BF8E7F4-213C-E847-5EBF-8017B8B2D1D3@t...      |     |               | 15/02/2022 1 |
| 80014550307     | UNIVERSITA' DEGLI STUDI DI UDINE                           | C98D4134.00F2B7FA.B42B3E9E.A2F69463.posta...   |     |               |              |
| 94150810300     | ENTE DI DECENTRAMENTO REGIONALE DI UDINE                   | 716EE483-B6A4-38FD-911C-48DC7D014296@t...      |     |               | 16/02/2022 1 |
| 80023240304     | LICEO GINNASIO STATALE JACOPO STELLINI UDINE               | opcec296.20220202121924.378025.452.2.23@pe...  |     |               | 15/02/2022 1 |
| 02162990309     | FONDAZIONE TEATRO NUOVO GIOVANNI DA UDINE                  | opcec296.20220119083257.01367.848.1.68@pec...  |     |               |              |
| 00164770307     | AZIENDA TERRITORIALE PER L'EDILIZIA RESIDENZIALE DI UDINE  | 25EF5973.027F352C.90916A4D.BDB32354.posta...   |     |               | 16/02/2022 1 |
| 02345670307     | SOCIETA' FERROVIE UDINE - CIVIDALE S. A R.L. UNIPERSONALE  | opcec296.20220128105202.30210.80.1.64@pec.a... |     |               | 15/02/2022 1 |
| 02935190302     | CAMERA DI COMMERCIO, INDUSTRIA, ARTIGIANATO E AGRICOLTU... | 25E9F05A.02939B28.9F9E6121.548A6726.posta...   |     |               | 15/02/2022 1 |
| 94106210308     | ORDINE DEI DOTTORI COMMERCIALISTI E DEGLI ESPERTI CONTA... | 25F5397B.029CAA6C.AF107A0C.E576A0A3.posta...   |     |               | 16/02/2022 1 |
| 94143540303     | CONS.COMUNI DEL BACINO IMBRIFERO MONTANO DELL'ISONZO ...   | 1FC02A94-9464-C4F9-F4BB-8A863C0250D2@te...     |     |               | 16/02/2022 1 |

# Data quality problems

- **Law 190/2012:** missing XMLs

| Year        | URLs   | Downloaded | NOT downloaded | % NOT downloaded |
|-------------|--------|------------|----------------|------------------|
| <b>2022</b> | 20,170 | 18,554     | 1,616          | <b>8.01</b>      |
| <b>2021</b> | 19,760 | 17,362     | 2,398          | <b>12.14</b>     |
| <b>2020</b> | 18,986 | 15,541     | 3,445          | <b>18.14</b>     |
| <b>2019</b> | 19,277 | 13,624     | 5,653          | <b>29.33</b>     |
| <b>2018</b> | 18,950 | 11,658     | 7,292          | <b>38.48</b>     |
| <b>2017</b> | 18,038 | 10,494     | 7,544          | <b>41.82</b>     |
| <b>2016</b> | 17,469 | 9,222      | 8,247          | <b>47.21</b>     |
| <b>2015</b> | 17,962 | 8,546      | 9,416          | <b>52.42</b>     |

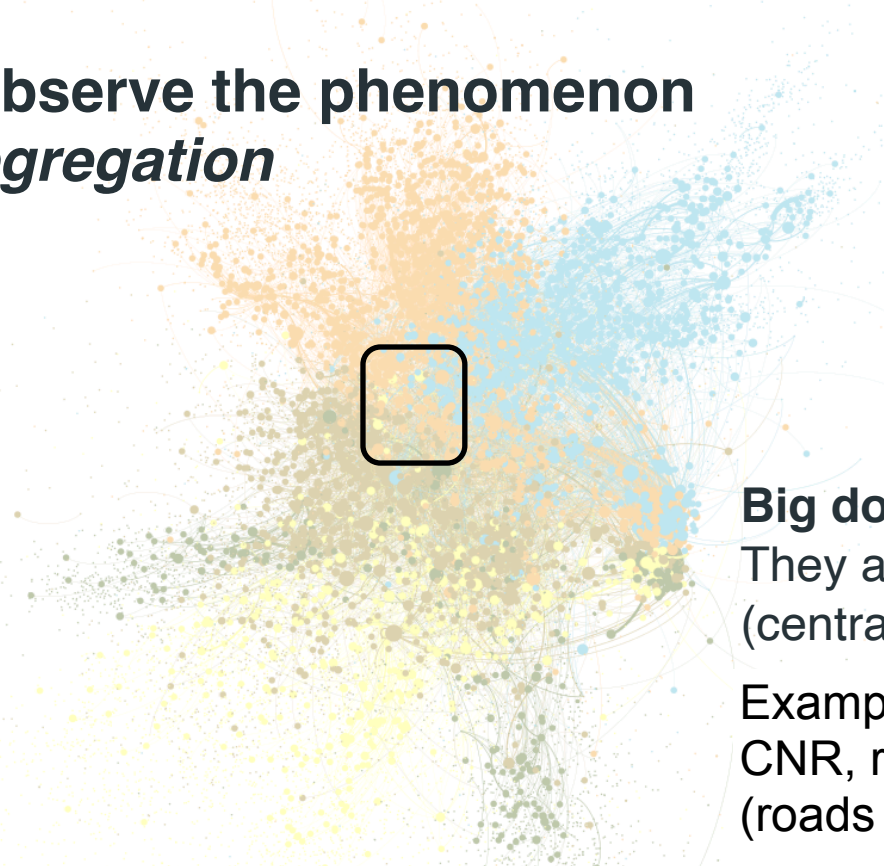
# Descriptive analysis: the graph of public contracts





# Analysis per macro-area

- © We observe the phenomenon of *segregation*



## Big dots:

They award contracts in entire Italy  
(central in graph)

## Examples:

CNR, roads signals, Tecnositaf  
(roads security)

A decorative network diagram in the top-left corner, consisting of various sized circles (nodes) connected by thin lines (edges). Some nodes are solid grey, while others are hollow with a grey outline. The connections form a complex, interconnected web.

4.

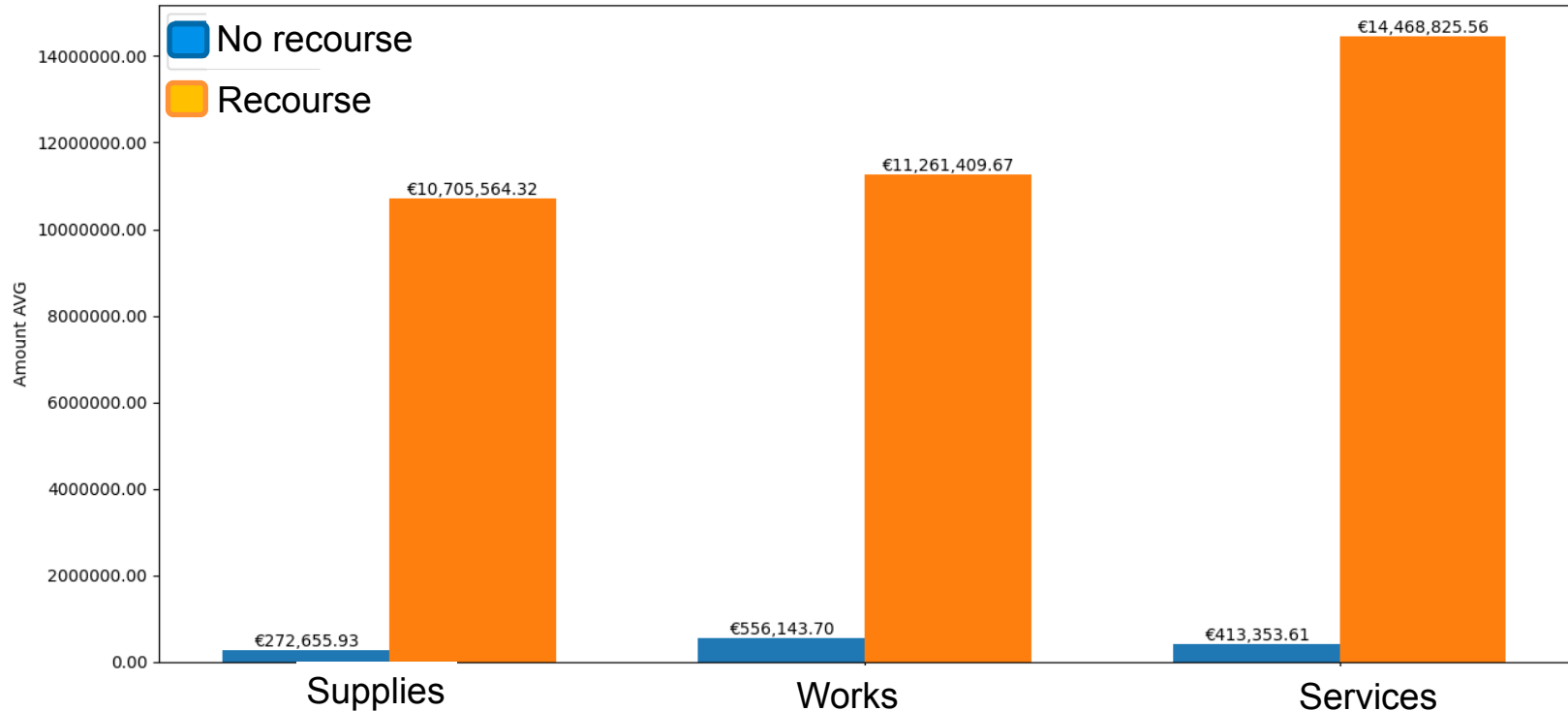
**Results**

# Search by procurement ID (CIG)

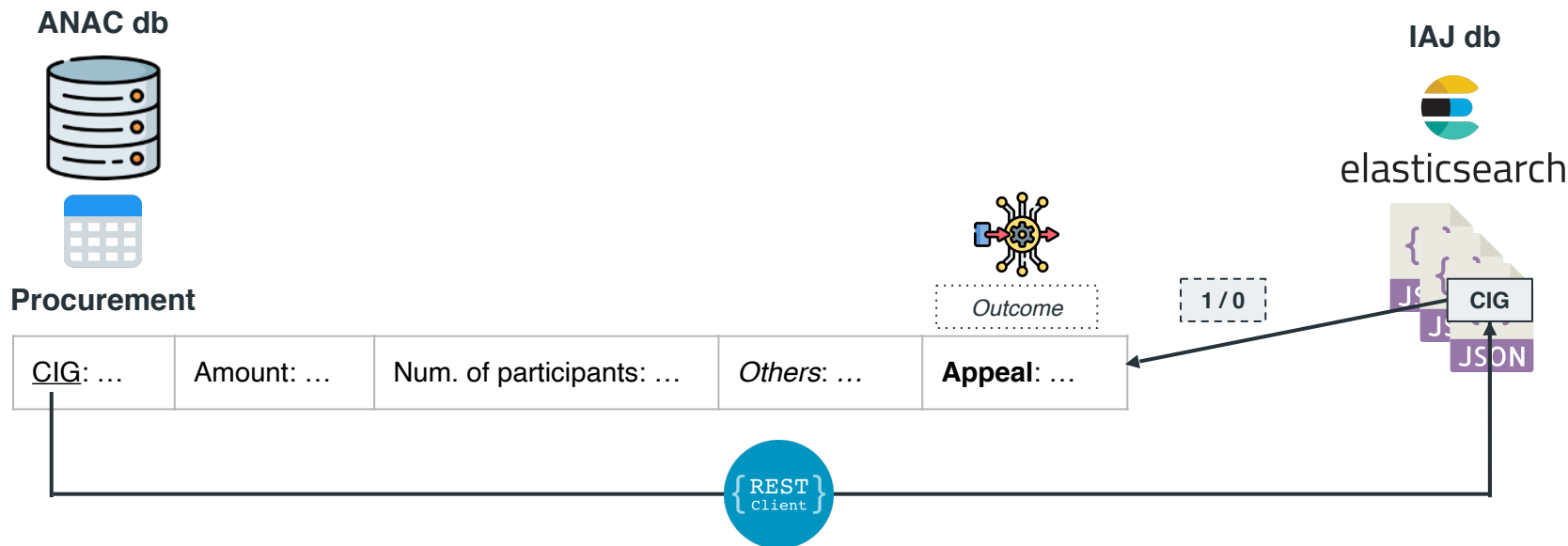
- Connection between the two information sources occurred by:
  - **searching for the procurement CIG in the IAJ sentences archive;**
  - the presence of the CIG in the IAJ dataset were used to label as positive the contract in ANAC procurement table
- 🎯 This is the supervised input to apply machine learning for the classification of “*irregular*” administrative practice in procurement
- However, the presence of a recourse could be due to a tendency of economic operators to stop the contract if they were not awarded a *substantious* contract

# Correlation with positive labels and contract amount

## Average Contract Amount per Procurement Category and Recourse

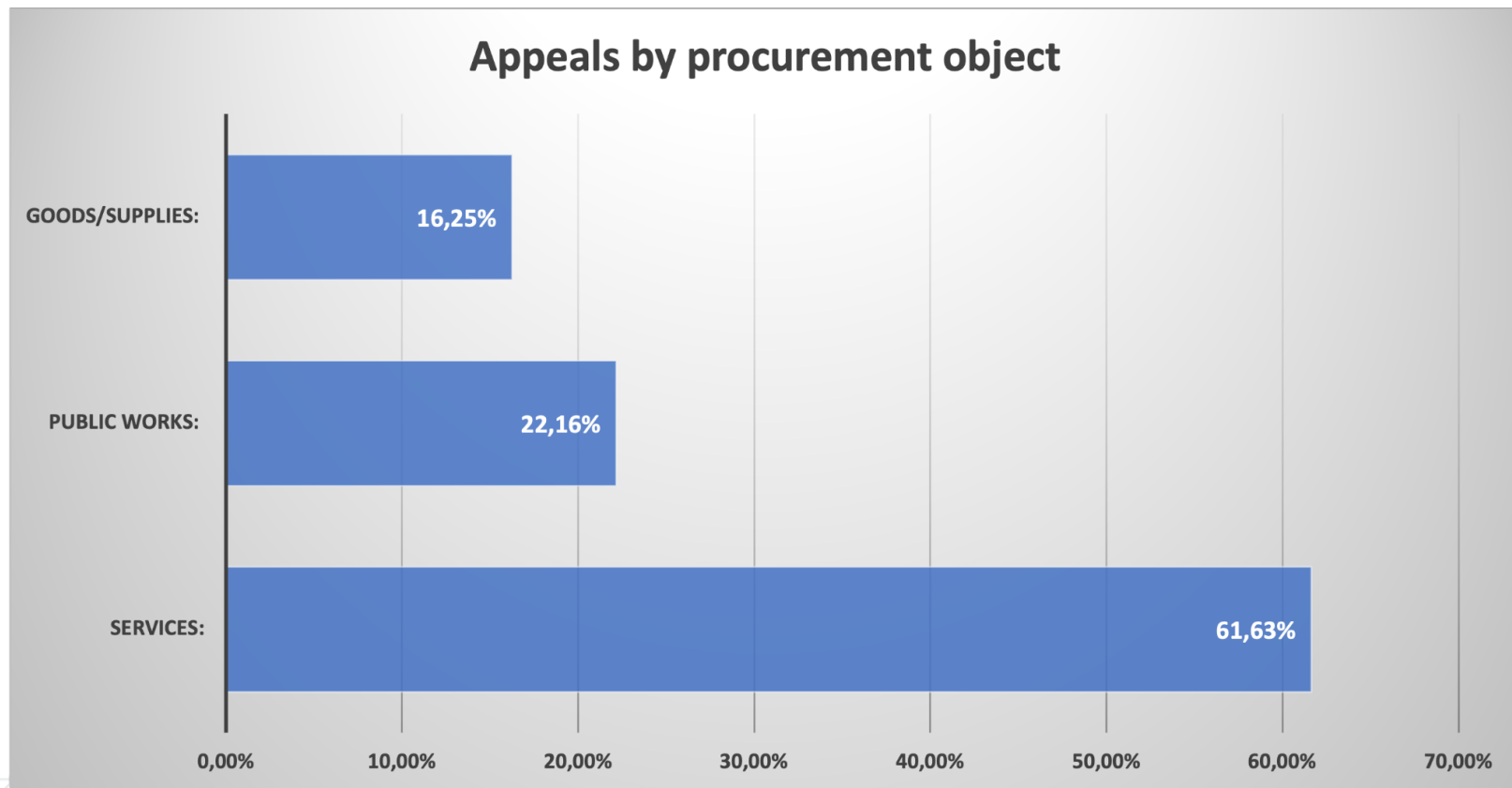


# Search by procurement ID (CIG)



The total number of CIG found is **8,062** over 80,360 judgments: this means that the probability that a sentence in our archive refers to a CIG is about 10%.

# Search by procurement ID (CIG)



# Search by participants and contractors

- The connection between the two information sources occurred by:
  - 🎯 **searching for the participants and contractors denominations in the IAJ sentences archive;**
  - the denominations found in the IAJ dataset were useful to find entities inside the legal archive and define a litigation measure between participants (economic operators) and contractors (PA).

# Search by participants' denomination

ANAC db



Participants

|                 |                   |                     |
|-----------------|-------------------|---------------------|
| <u>CE</u> : ... | Denomination: ... | <i>Others</i> : ... |
|-----------------|-------------------|---------------------|



{REST  
Client}

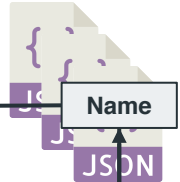
IAJ db



elasticsearch



Litigation  
measure



Name

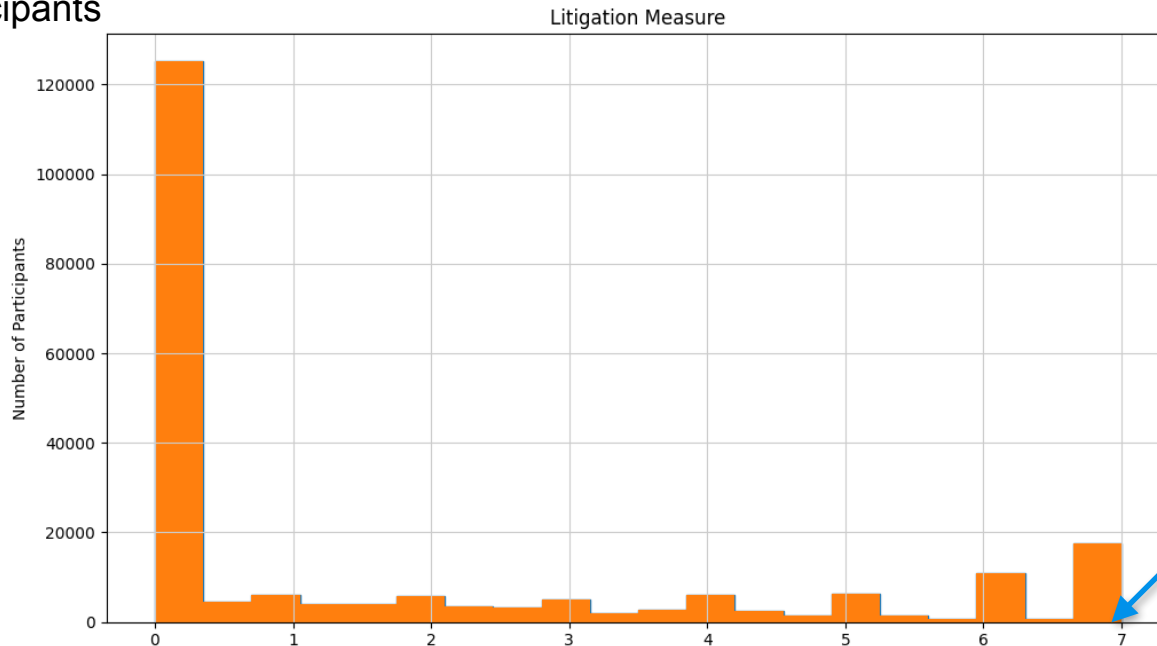
JSON



# Results of detection of participants' in IAJ

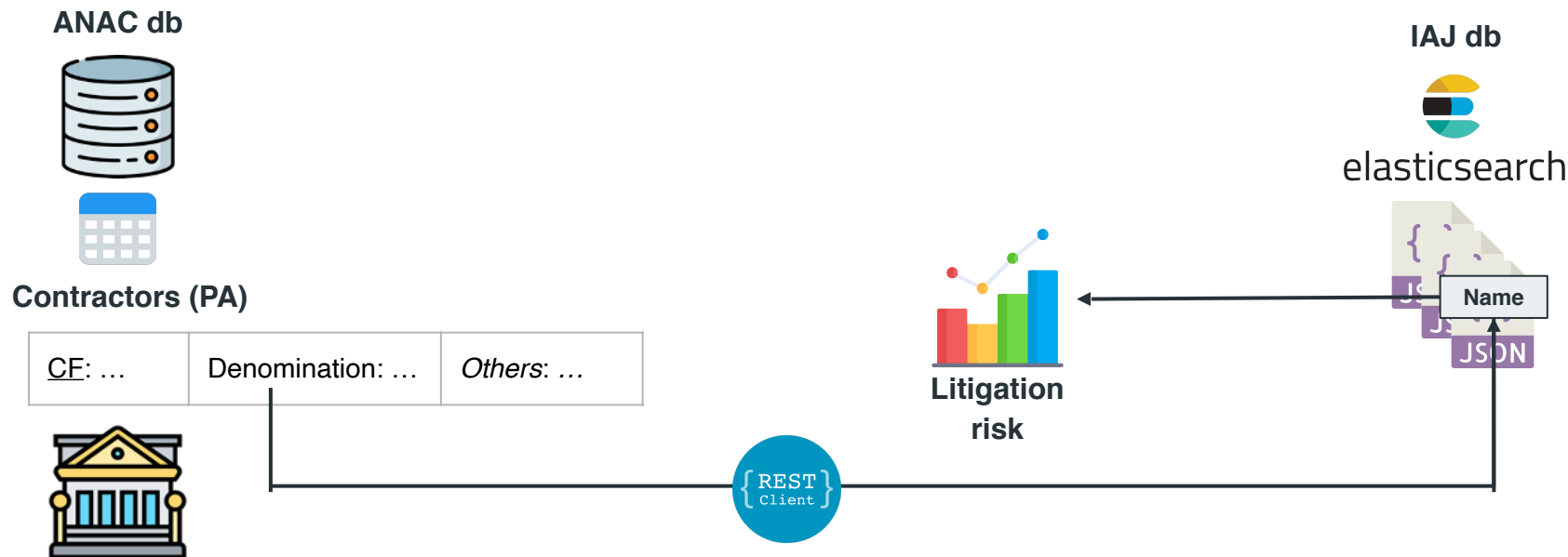
$$\text{Litigation measure} = \frac{\text{Number of recourses}}{\text{Number of participations}}$$

Number of participants



7 is the extreme of the confidence interval at alpha=5%

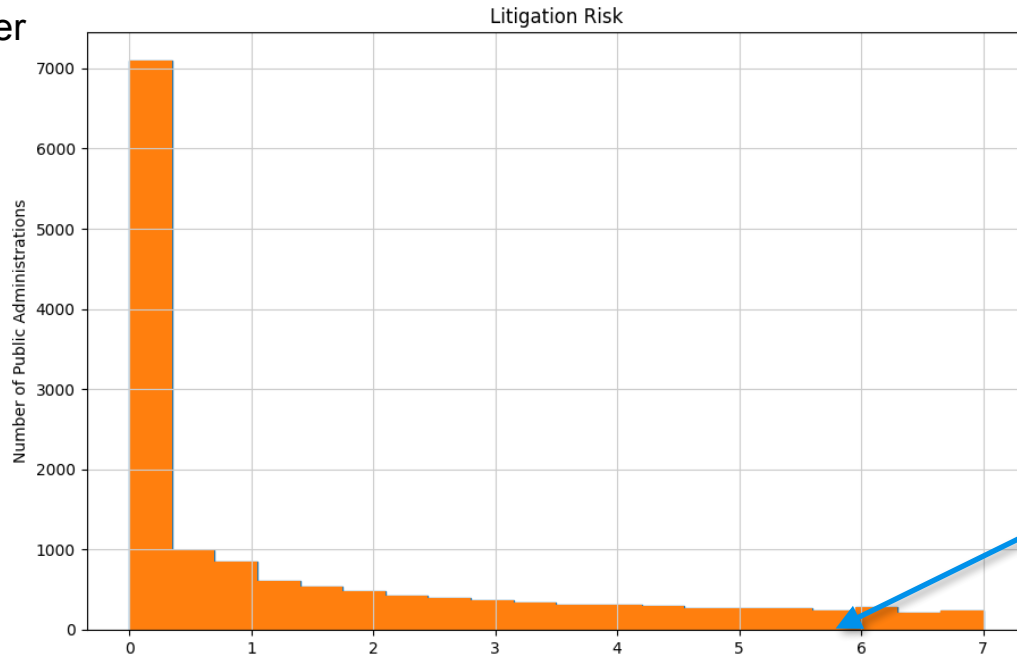
# Search by contractors' denomination



# Results of detection of contractors

$$\text{Litigation risk} = \frac{\text{Number of recourses}}{\text{Number of contracts}}$$

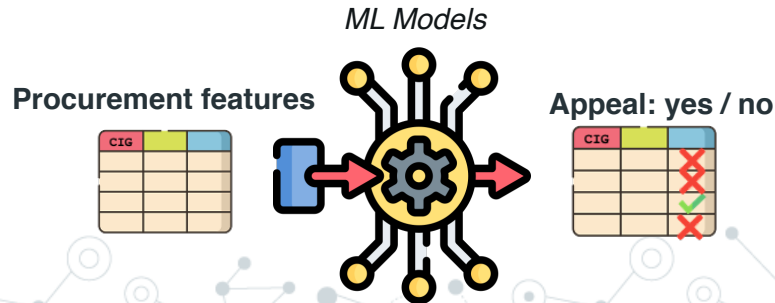
Number of PAs



5.9 is the extreme of the confidence interval at alpha=5%

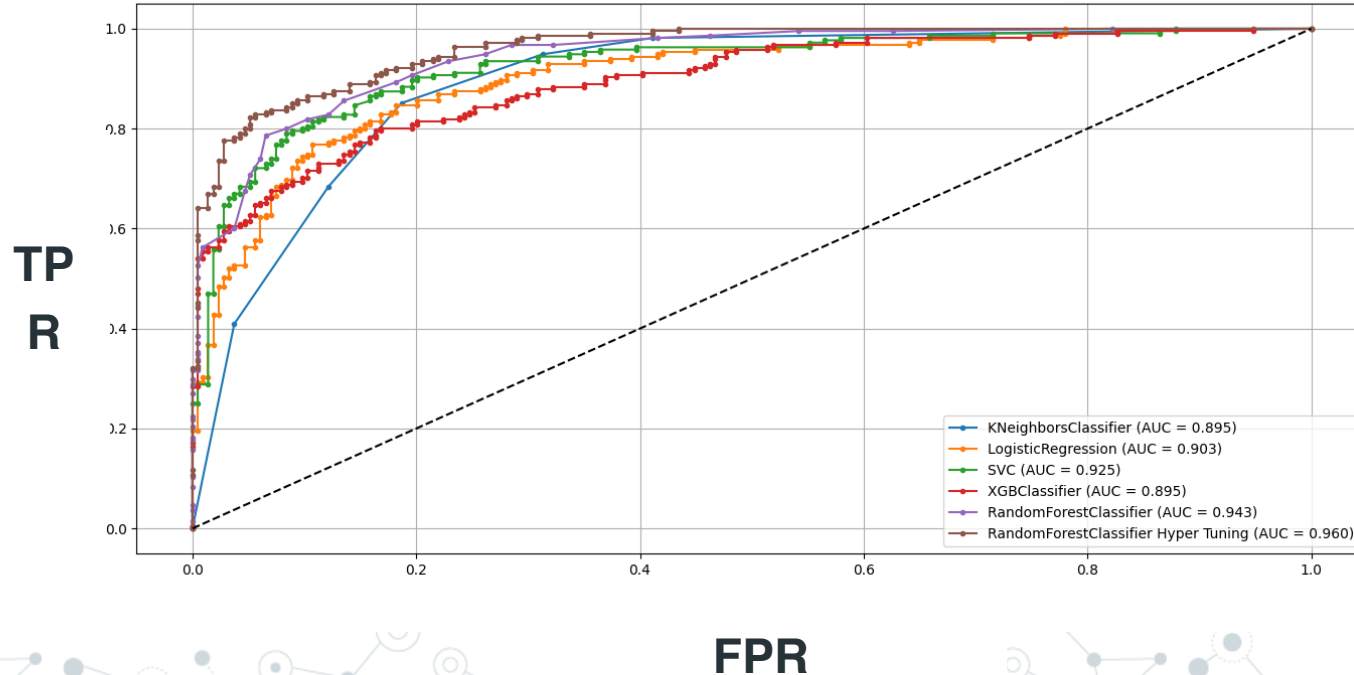
# Machine Learning Models

- Following the RQ2, the problem statement is given: find an algorithm such that, given the description of a procurement, it predicts with the highest expected accuracy, the presence of a possible recourse related to that procurement to the Administrative Justice courts.
- We explore consolidated MLAs: it can help to dig into a large amount of data and to discover the predictive patterns based on the procurements features.



# Machine Learning Models Results

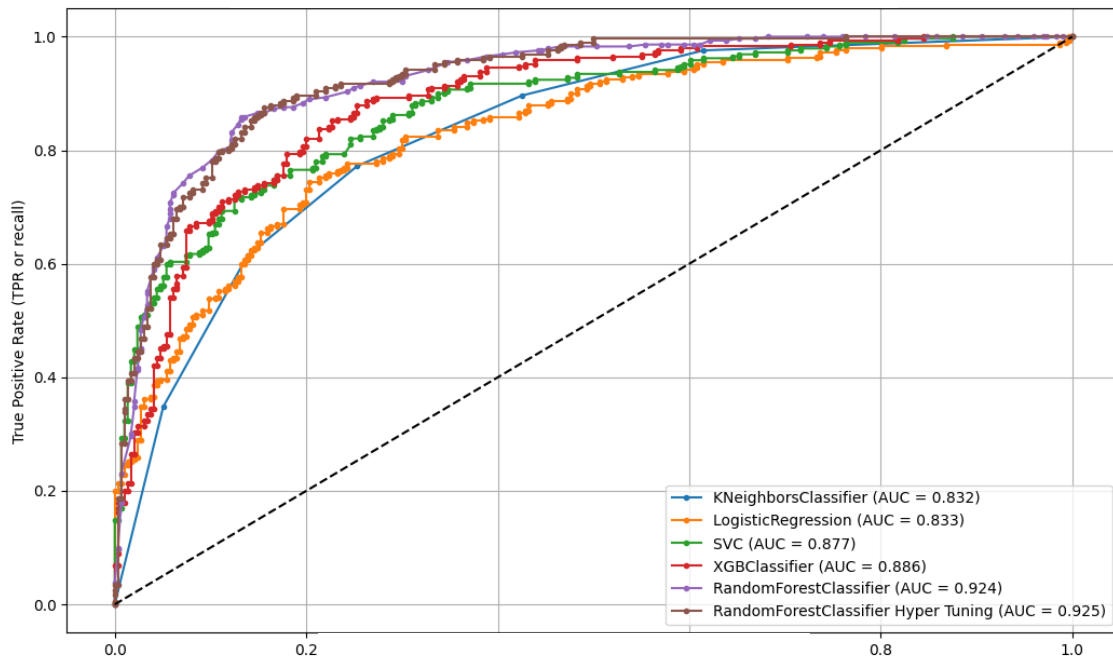
Procurement object: good/  
supplies



# Machine Learning Models

Procurement object:  
public works

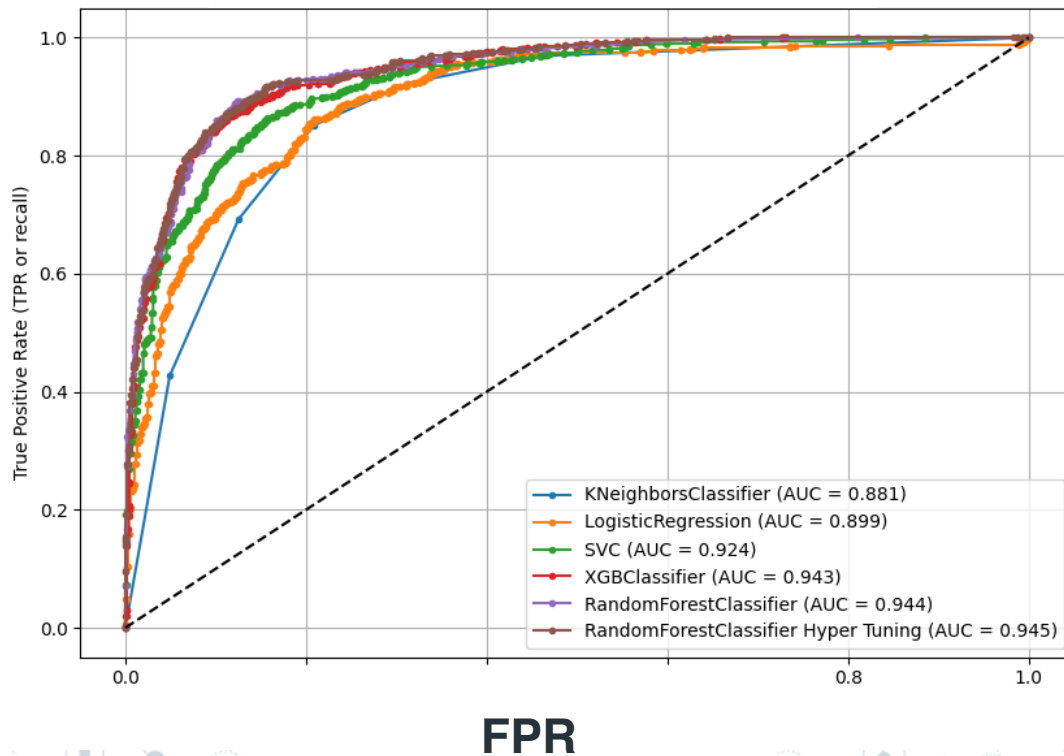
TP  
R




# Machine Learning Models

Procurement object:  
services

TP  
R



# Recommender Systems on Contracts (RQ3)

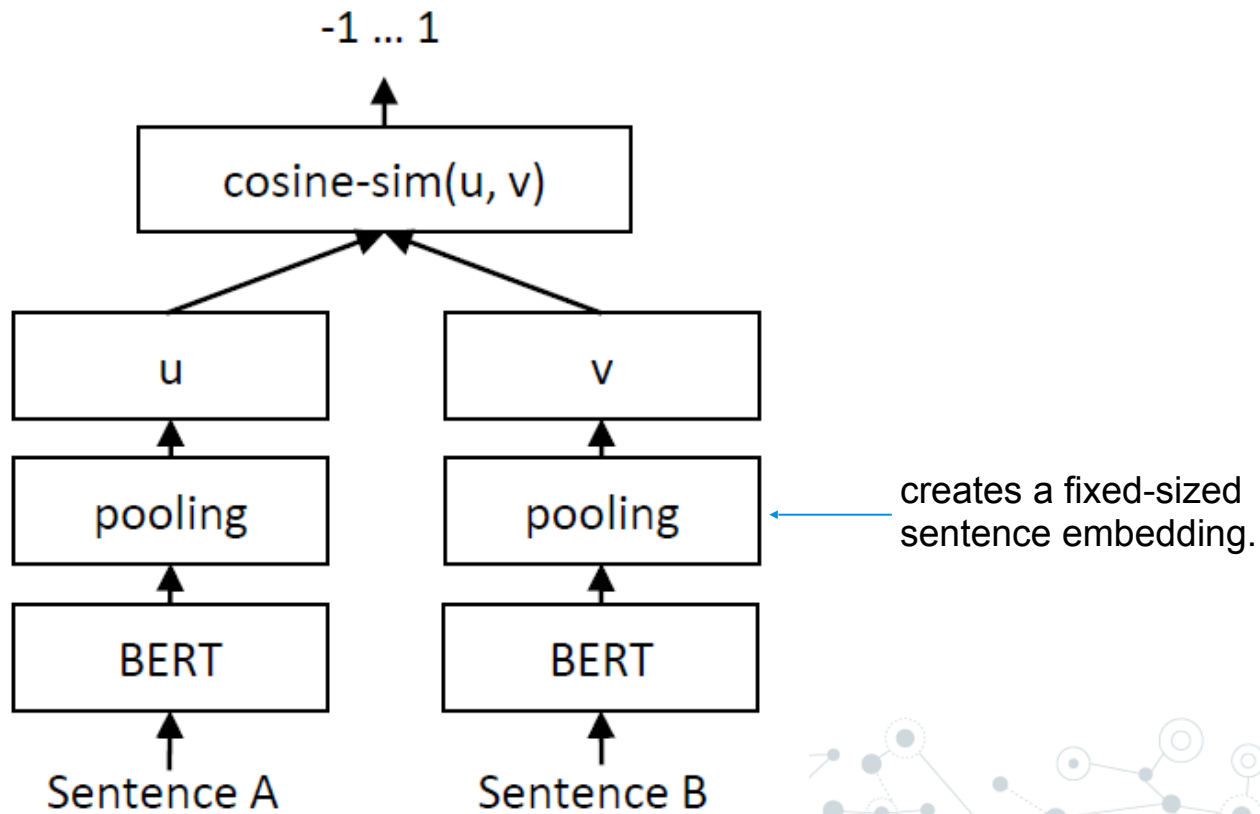
 We want to provide systems to help people (employees from PA or from economic operators searching for the right tenders) working and finding contracts in a large database

- We trained recommender systems on *oggetto del contratto*, a short description that the responsible person of the procurement included in the tender.
- A problem we have in the ANAC data is that the textual descriptions are often ambiguous and carelessly filled by the responsible people of the tenders



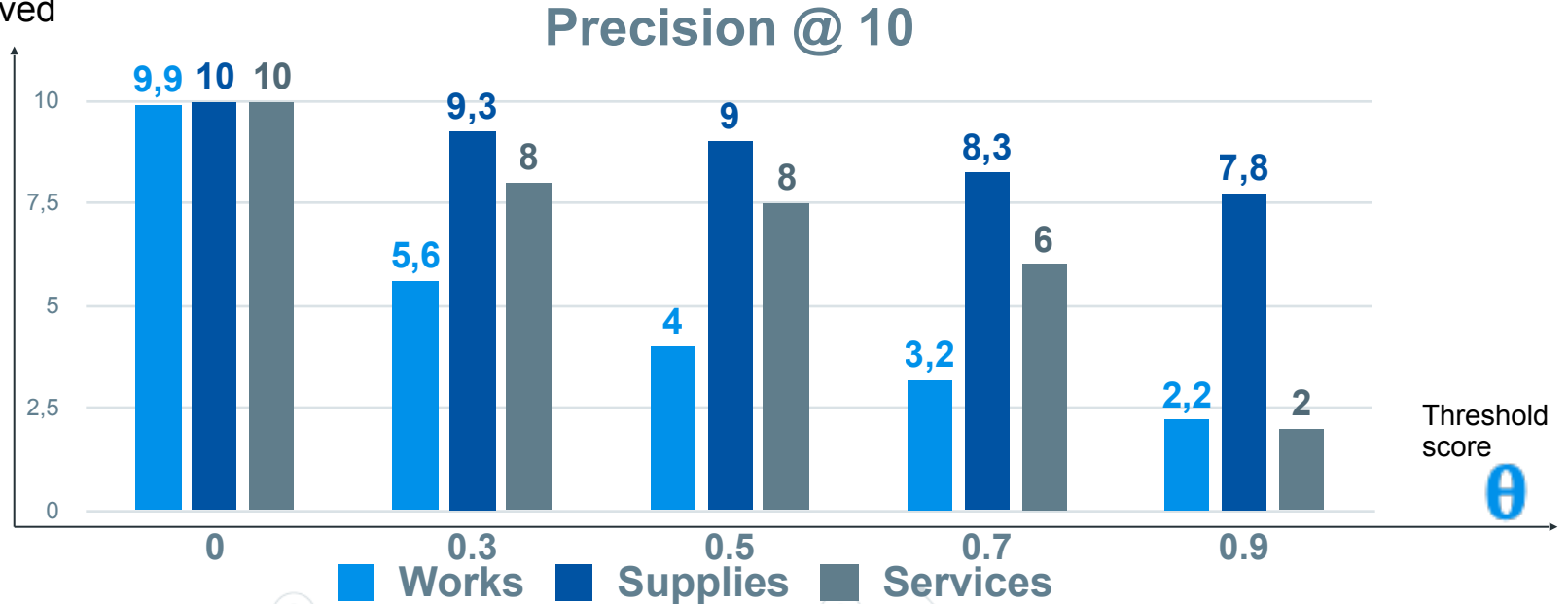
# Recommender Systems on Contracts

- We applied S-Bert Deep Neural networks:  
<https://ai.googleblog.com/2020/08/language-agnostic-bert-sentence.html>  
on the short textual description
- Transformers outperform both RNN and CNN in translation benchmarks and are trained by ranking on **semantic similarity** a large set of sentences.
- Translation occurs first by encoding the meaning of each word in the sentence and then decoding.
  1. We use the encoding vectorial representation of the textual *object of the contract*
  2. use *cosine similarity* to retrieve the most similar contracts in the database.
- Transformers, on top of higher translation quality, require less computation, speeding up training by up to an order of magnitude.



# Evaluation by a panel of 3 independent experts on a small random training set

Similar contracts retrieved



## Problems with *oggetto del contratto*

- A problem we have in the ANAC data is that the textual descriptions are often ambiguous and carelessly filled by the responsible people of the tenders
- **Examples:**
- Open awarding procedure for assigning a public work contract:  
*LAVORI DI 'RESTAURO ARCHITETTONICO DEGLI INFISSI ESTERNI E LA MANUTENZIONE DELLA COPERTURA DEL PALAZZO PRINCIPE DI ARAGONA SEDE DELL'OPERA PIA 'ISTITUTO PRINCIPE DI ARAGONA'*
- It is described in a very similar way of the following service whose description is ambiguous:
- LAVORI DI RISANAMENTO CONSERVATIVO PER LA MANUTENZIONE EDILE ED IMPIANTISTICA PRESSO IL PALAZZO GADALETA, SEDE DEGLI UFFICI GIUDIZIARI DI TRANI - SERVIZI DI ARCHITETTURA ED INGEGNERIA

A decorative network diagram in the top-left corner, consisting of various sized circles (nodes) connected by thin lines (edges). Some nodes are solid grey, while others are hollow white with a grey outline. The connections form a complex, interconnected web.

# 5.

## Conclusion and future work

# Conclusion

- This work demonstrates the ability to manage a huge juridical dataset from a national public entity to automatically extract meaningful knowledge to address machine learning experiments (RQ1).

# Conclusion

- In addition, we explored the results of a predictive experiment to test recourse prediction of the administrative courts on the basis of the features of a public procurement (RQ2).

## Future work

- As a future work, we plan to investigate furthermore the explainable AI techniques,
- We need to answer the *causal* questions from the domain experts (law people) on the reasons of a recourse or of an award to an economic operator.



## Future work

- As a future work, we plan to expand the usage of sentence embedding and cosine similarity to match the short description of a recourse and of contracts
- We hope this will help us reduce the number of candidate recourses without a corresponding match in the contract database

## Future work

- Another direction concerns the adoption of *process mining* techniques for conformance checking and predictive process monitoring applied to a log of temporal events obtained by this database.

# Thanks!

## Any questions?