

Authors: Alessandro Bombini<sup>1</sup>, Lisa Castelli<sup>1</sup>, Achille Felicetti<sup>2</sup>, Franco Niccolucci<sup>2</sup>, Anna Reccia<sup>2</sup>, and Francesco Taccetti<sup>1</sup>.



1. Istituto Nazionale di Fisica Nucleare (INFN), Sezione di Firenze, Florence, Italy; 2. Vast-Lab, PIN, Prato, Italy

## Abstract

When researchers and professionals try to retrieve stored analysis raw data, they usually face three major issues: data availability, missing (meta)data standardisation, and difficult query creation for database fetching.

To ease these issues, in the framework of the Eosc-Pillar project, the cloud platform *Tools for HERitage Science Processing, Integration, and ANalysis* (THESPIAN) was developed [1].

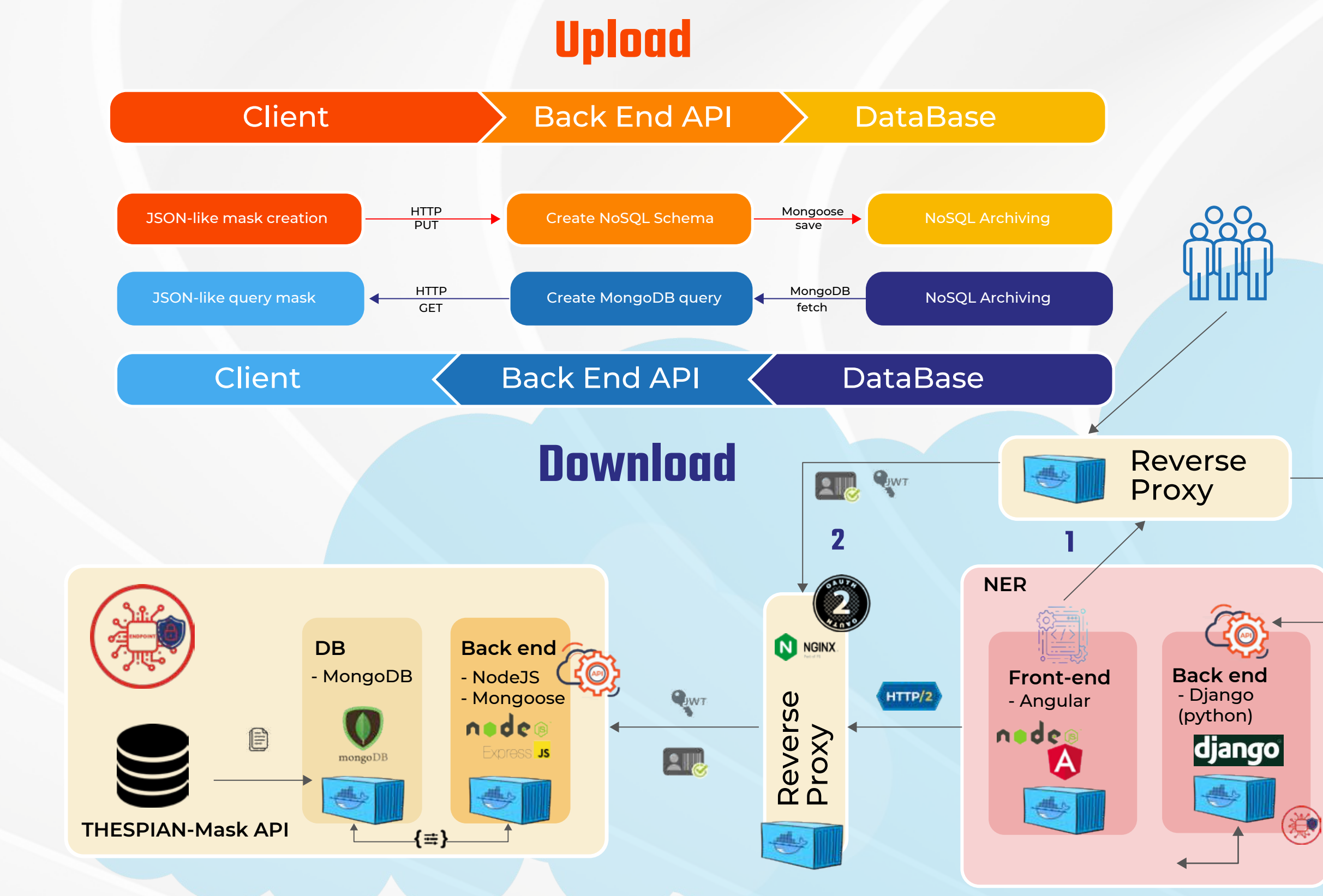
The first web service of the THESPIAN platform, THESPIAN-Mask, was developed to resolve the first two issues; it is a web-app for FAIR storage of scientific analysis on Cultural Heritage, and it is based on an hoc developed CIDOC-CRM-compatible ontology, CRMhs, describing the metadata of scientific files.

To ease the third issue, and also to help the metadata generation process, an additional cloud-native tool was developed: **THESPIAN-NER**.

It is a tool based on a deep neural network for Named Entity Recognition, enabling users to upload their Italian-written report files and obtain labelled named entities usable as keywords, either for (semi)automatically compose custom queries to the database, or fill (part of) the metadata form describing the file to be uploaded.

## Why THESPIAN-Mask?

- ★ Implements FAIR principles [2]
- ★ Offers assisted metadata generation
- ★ Cloud storage
- ★ Based on ad-hoc ontology, makes information fully interoperable with other data



## Why THESPIAN-NER?

- 👍 Offers assisted metadata generation
  - 👍 Based on ad-hoc ontology, makes information fully interoperable with other data
- To ease these two issues, the cloud-native web-app THESPIAN-NER was developed
- 👍 (semi)automatically generate custom THESPIAN-Mask queries, by similarity with a user-selected italian written documents or reports;
  - 👍 (semi)automatically fill-in a certain number of THESPIAN-Mask metadata fields, again by similarity with a set of Italian written documents or reports;

THESPIAN-NER - Named Entity Recogniton tool for Archeology

Collaborative semantic enrichment of text-based datasets - Italian language version

The screenshot shows the THESPIAN-NER web application interface. The main area displays a text document with various entities highlighted in different colors (ACT, TSP, SIT, ART, PRD, PER). Below the text, there is a sidebar with a list of entity types and their counts found in the text: Timespan (11), Artefacts (7), Site (4), Location (4), Person (2), Activity (2), Period (1), Organisation (0), and Biological Remains (0). At the bottom, it shows 'Number of selected entities: 0'.

## THESPIAN-NER web-app

### Front-End Web UI for

- ★ Upload .txt, .pdf documents
- ★ Easily compose custom queries
- ★ Easily fill in the metadata form

### Back-End RESTful-API(s) for

- ★ Perform NER analysis of uploaded text
- ★ Query the DataBase
- ★ Open THESPIAN-Mask



The Authentication/Authorisation process is handled by the reverse proxy

## DATASET(s)

### ArcheoNER:

- ★ 9 Entity labels
- ★ 92 Italian-written docs
- ★ 5230 total entities

### hsNER

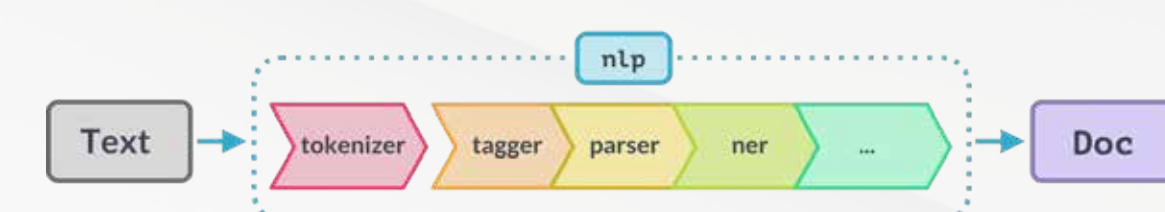
- ★ 14 Entity labels
- ★ 43 italian-written reports
- ★ 5676 entities



## Fine tune the spaCy NER model - I

- ★ open-source models used;
- ★ two models used:
  - ★ CNN-based
  - ★ Transformer-based

	ArcheoNER	hsNER	THESPIAN-Mask JSON keys
ACT	✓	✓	sample.preparation.method
ANL		✓	analysis.category.type
ART	✓	✓	studyObject.name
BIO	✓	✓	studyObject.name
DEV		✓	analysis.device
LOC	✓	✓	studyObject.locationLabel
MAT		✓	sample.material
MET		✓	sample.preparation.method
NAT		✓	studyObject.name
ORG	✓	✓	studyObject.owner
PER	✓	✓	studyObject.author
PRD	✓	✓	studyObject.periodLabel
SAM		✓	sample.type
SIT	✓	✓	studyObject.provenanceLabel
SOF		✓	analysis.device.software
RES		✓	analysis.result
TSP	✓	✓	analysis.startDate(endDate)



## Fine tune the spaCy NER model - II

- ★ Two training for each model:
  - ★ ArcheoNER, for Archaeological documents
  - ★ hsNER, for internal lab reports on heritage science.

Cons:

- 👎 Difficult to get huge, high quality training data;

Pros:

- 👍 Model still get the most relevant entities for the tasks (metadata generation and query creation)

## Acknowledgments

The present research has been partially funded by the European Commission within the Framework Programme Horizon 2020 with the projects EOSC-Pillar (Grant agreement number H2020-INFRAEOSC-05-2018-2019-857650)

