



## TETYS: towards the next-generation open-source Web topic explorer



Anna Bernasconi, Francesco Invernici, and Stefano Ceri

Contact: anna.bernasconi@polimi.it

Department of Electronics, Information and Bioengineering, Politecnico di Milano, Milan, Italy

Web users who search specialized content typically lack knowledge of the precise topology of the dataset.

 $\rightarrow$  We propose a *next-generation open-source* Web topic explorer, an architecture with:

- a **pipeline** for ingesting huge data corpora, extracting highly relevant topics 1)
- an interactive dashboard, supporting topic visualization, exploration of temporal series, and statistical testing 2)

Our prototype, CORToViz, explores the CORD-19 dataset (COVID-19 research abstracts). Other domains will be explored in TETYS (e.g., climate change).

## Topic modeling

Ο



- Based on up-to-date technologies (including LLMs) and self-optimizing models (see BERTopic framework [1]);
- Unsupervised topic modeling via clustering of articles along the orthogonal dimensions of the document embeddings (latent representations generated by language

## models) – computed with SPECTER [2];

- Visualization of topics through wordclouds (generated via gensim Python package)
- Temporal topic modeling via time series mining and statistical testing.



## CORToViz dashboard (http://gmql.eu/cortoviz) [3]



- Line plot of the intensities (i.e., the relative frequencies of appearance) of the two selected topics – with configurable bin resolution (1-4 weeks).

