# Semantic Query Analysis from the Global Science Gateway

*Sara Goggi\*, Gabriella Pardelli\*, Roberto Bartolini\*, Monica Monachini\*, Stefania Biagioni\*\*, Carlo Carlesi\*\**

*Istituto di Linguistica Computazionale, "A. Zampolli", CNR-ILC, Italy \**
*Istituto di Scienza e Tecnologie dell'Informazione "A. Faedo", CNR-ISTI, Italy \*\**

## Scenario

❑ Web portals play an essential role in searching and retrieving information in the several fields of knowledge
❑ Web portals support the storage of a huge amount of information in NL originating from the queries launched by users worldwide
❑ A good example is given by *WorldWideScience.org* (The Global Science Gateway)

## Objective

❑ The aim is to retrieve information related to *social media* which as of today represent a considerable source of digital data more and more widely used for research ends

*the terms convey meaning*

## Focus

❑ The query logs registered by the *GreyGuide: Repository and Portal to Good Practices and Resources in Grey Literature* and received by the *WorldWideScience.org* portal
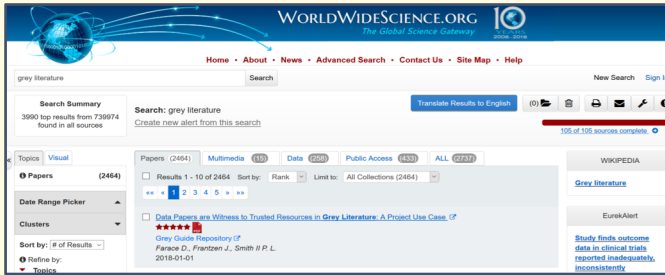

*https://worldwidescience.org/*


*http://greyguide.isti.cnr.it/*

## This project includes eight months of query logs registered between July 2017 and February 2018 for a total of 445,827 queries

## Methods and Tools

❑ A process of information retrieval from a rich digital catalogue of queries
   1) cleaning of the set of queries;
   2) filtering and ordering (alphabetically);
   3) using several trials for choosing the focus;
   4) processing the information and building the sample by NLP tools.

❑ NLP analysis
   a) free information extraction:
      ➢ measure the frequency of the words contained in the corpus;
      ➢ examine the lexical variety of the queries ;
      ➢ focus on a set of terms to build a micro-ontology.

   b) ontology-based extraction:
      ➢ enrich the domain;
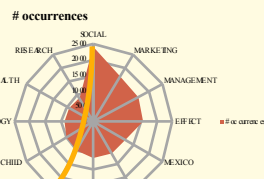      ➢ retrieve each occurrence of those terms contained in the ontology by using a search engine.
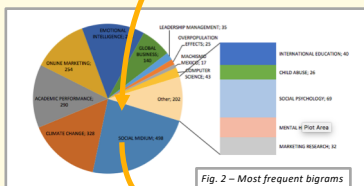

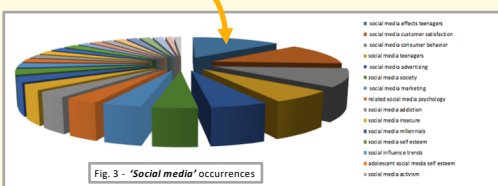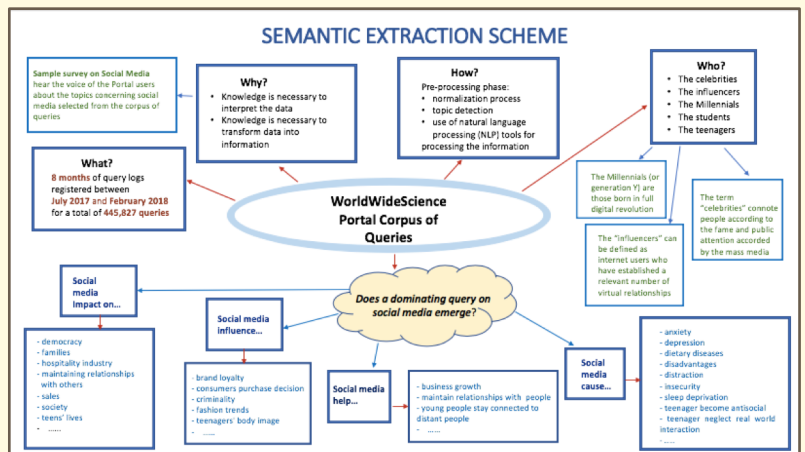*Fig.1 – Most frequent words*


*Fig. 2 – Most frequent bigrams*


*Fig. 3 - 'Social media' occurrences*


SEMANTIC EXTRACTION SCHEME

## WorldWideScience and 'Social Media'

### Why 'social media'

❑ Social Media are a very effective means of communication and vehiculate knowledge
❑ They are often quoted in bibliographical references amongst the more traditional categories
❑ The subject involves document types pertaining to Grey Literature

❑ A case study has been carried out involving **medicine, psychiatry and *'social media'*** Figures 1, 2, 3

❑ Some low-frequency terms (hapax) carry a **negative connotation** in relation to the use of *'social media'*:

   *<cyberbullying social media>; <depression social media>; <eating disorder social media>; <negative effects social media young adults>; <anxiety social media>; <social media compulsive buying>; <social media distraction>; <fake news social media>*

❑ An analysis of negative connotations in connection with *child/children*, is further investigated, as shown in Fig. 4
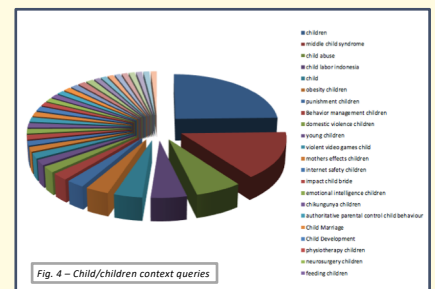

*Fig. 4 – Child/children context queries*

## Final Remarks

❑ The case study illustrates the main linguistic features of the *Global Science Gateway* by showing the lexical map which represents the most used/recurrent words
❑ Some critical issues: a diachronic analysis of the terms was not possible given the short temporal window taken into account; queries in different languages and many spelling/grammatical errors made our task more complicated by weighing the cleaning process down
❑ Terms extracted from the corpus of queries are largely referring to topics pertaining to the major problems of today's society, eg. *alcoholism, depression, obesity, pornography, drugs, violence....*
❑ NLP analysis allowed to browse the corpus through the most and less queried terms: once *social* has been identified as the most frequent one, the analysis was channeled into *'social media'* and the pertinent contexts.