# A brand scoring system for cryptocurrencies based on social media data

**Santomauro Giuseppe**    /    *DTE-ICT-HPC, ENEA – Italy*

**Alderuccio D., Ambrosino F., Fronzetti Colladon A., Migliori S.**

# Outline

- **The ENEA context**
  - ➢ ENEAGRID environment and CRESCO HPC clusters
  - ➢ Web Crawling in ENEAGRID
  - ➢ Semantic Brand Scoring in ENEAGRID

- **Proposal of current development**
  - ➢ Social Networks Crawling
  - ➢ Semantic Brand Score for Cryptocurrencies

- **Conclusions**

# ENEA

1 Headquarters;
9 Research Centers;
5 Laboratories.
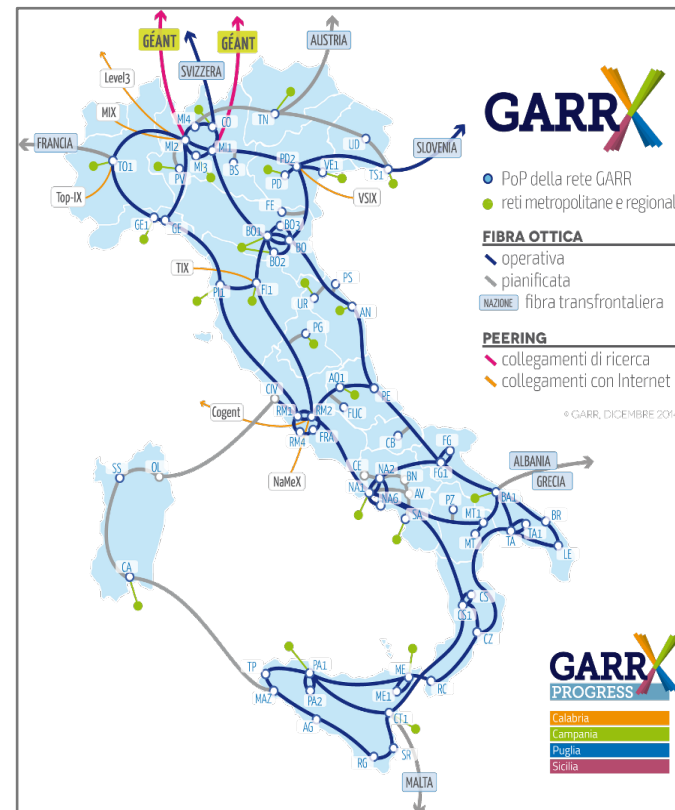


## Portici Research Center

## ENEAGRID

Computation & Storage **ENEA** distributed resources interconnected via **GARR** network.

### CRESCO HPC Clusters:

- 6 *Data Center*s in ENEA (**Portici** is the main site);

- More than 20000 cores;

- More than 400 computing nodes:
  - *Linux x86_64 + Special systems (GPU, PHI);*

- Storage resources:
  - *AFS (distributed);*
  - *GPFS (parallel high-speed) ~2.2PB;*

- Cloud computing facilities *(Openstack, VMWare);*

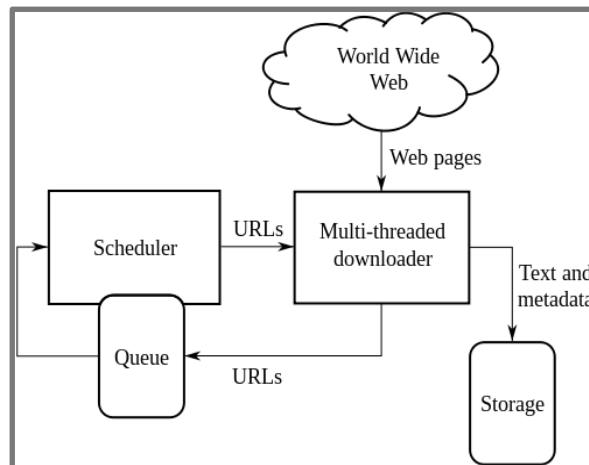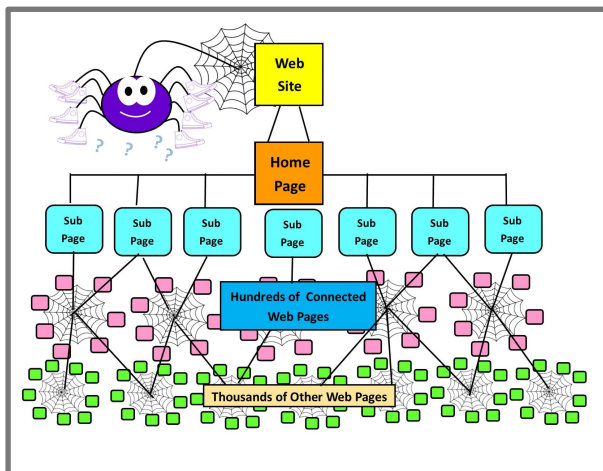- More than 1.4Pflops
  - *in **Top500** rank at Nov. 2018.*

http://www.cresco.enea.it

https://www.garr.it

# Web Crawling in ENEAGRID

➢ Activity to browse *www* systematically and download web content;

➢ *Google*, *Bing* and *Yahoo* periodically download the content of a wide web space;

➢ Data are stored locally and processed to build indexes, statistics and to structure them;

**Application contexts:**

- Web Searching;
- Intelligence & security;
- Blog analysis;
- User behaviors;
- Marketing.

# Web Crawling: Software and Virtual Lab

## Problems, rules and requirements

- Best practices to avoid improperly using the network (e.g. DDOS);

- Laws to comply with privacy and/or copyright (e.g. GDPR).

- Open source solution (BUbiNG*)
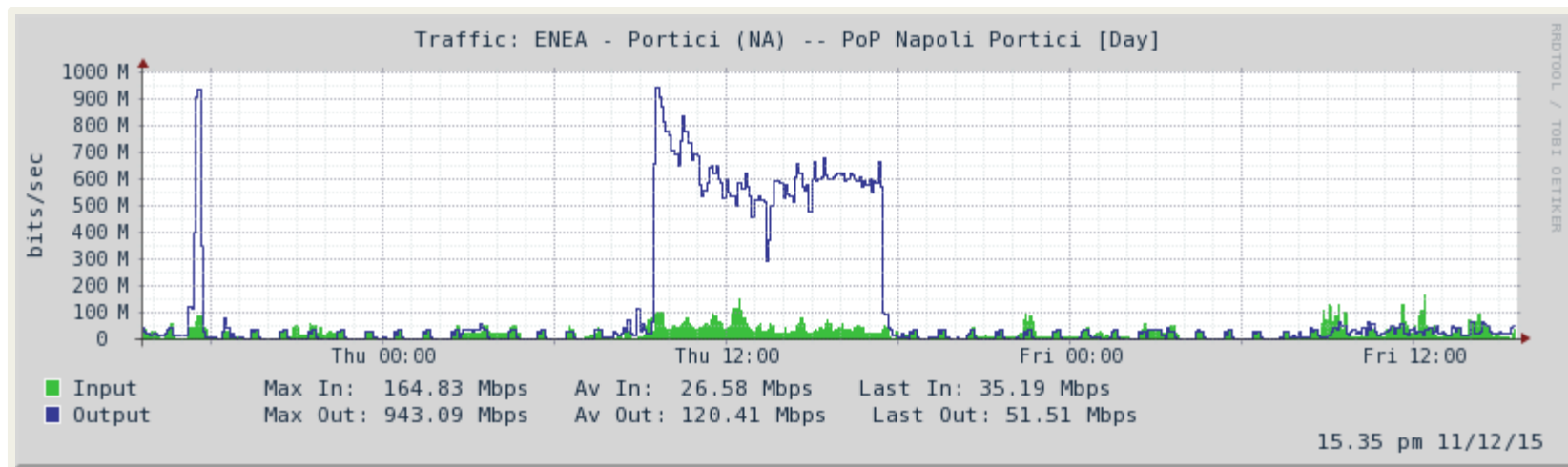
## Virtual Laboratory



*[Boldi, P., Marino, A., Santini, M., Vigna, S.: BUbiNG: Massive Crawling for the Masses. (2016)]

# First Test: **Single Snapshot**
## (Efficiency and Robustness)

**Number of agents:** 16;
**Running time:** 8 h;
**Amount of downloaded data:** 2.94 TB;
**Amount of downloaded resources:** 66.806.790 Pages;
**Data downloading speed:** 850 Mbps;
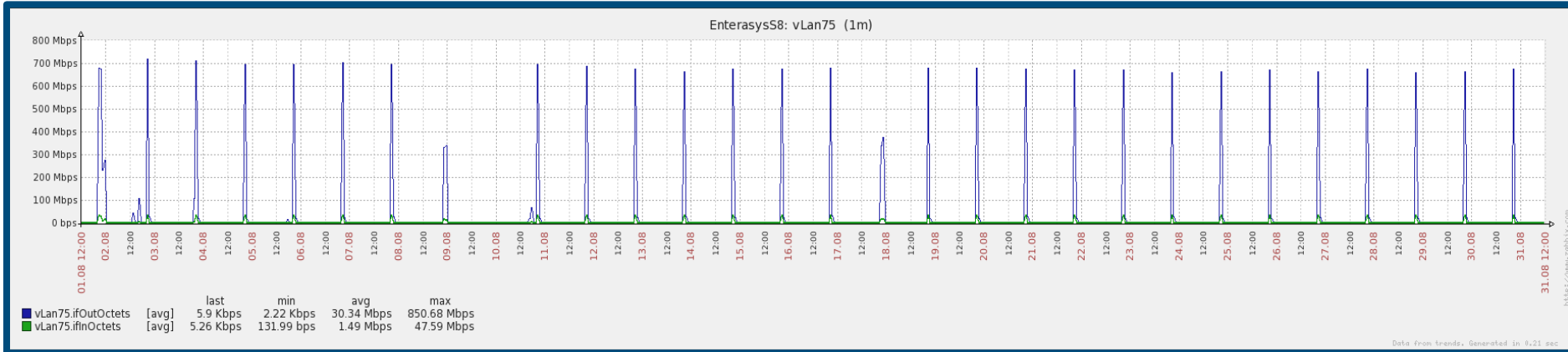**Resources downloadeding speed:** 2305 Pages/Sec.



Network traffic measured by GARR on Napoli-Portici PoP.

[Santomauro, G. et Al.: A collaborative environment for web crawling and web data analysis in ENEAGRID. In: DATA 2017, (2017)]

# Second Test: **Periodic Snapshots** (Reliability)

**Daily** web crawling sessions, during **one month**, each of them kept alive for **one hour** (from 21:00 to 22:00), by considering only web pages from *.it* domain.



Network traffic in the snapshot period.



Average of download speed for each snapshot.

- 15 TB of downloaded data (484 GB/snapshot);

- 3,3 TB saved on the storage (111 GB/snapshot);

- **Downloading Speeds**: $\mu$: 1,00 Gbps, $\sigma$: 0,0005 Gbps.

# *Semantic Brand Scoring* in ENEAGRID

➤ The Semantic Brand Score* (SBS) is a novel metric designed to assess the importance of one or more brands, in different contexts and whenever it is possible to analyze textual data, even big data.

➤ The advantage with respect to some traditional measures is that the SBS do not relies on surveys administered to small samples of consumers.
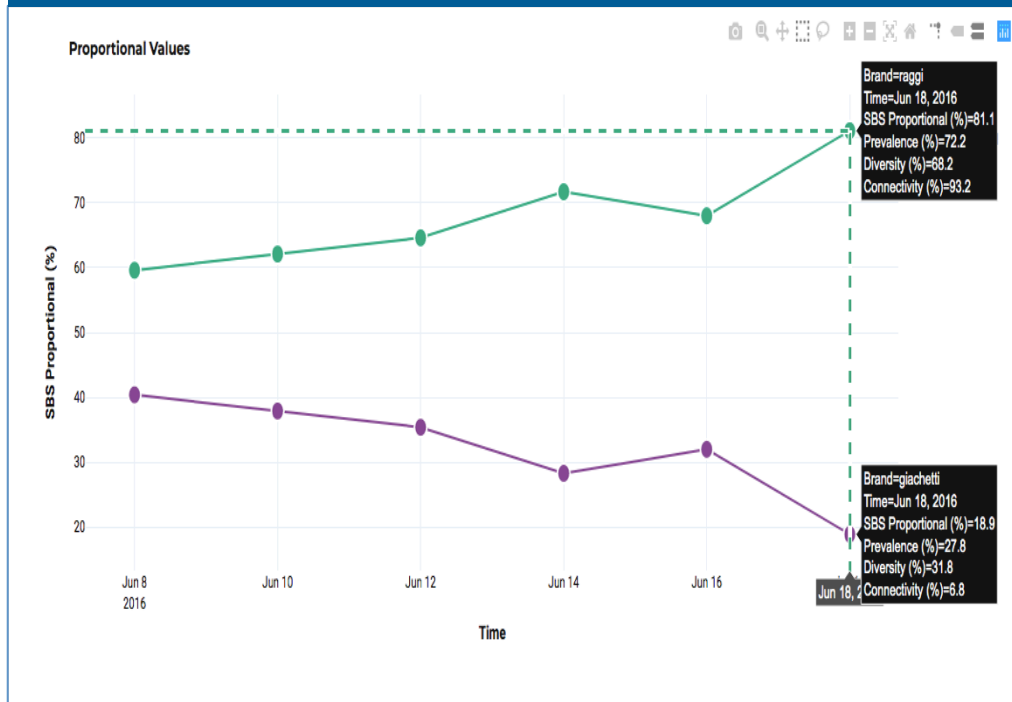
## The metric is measured along the 3 dimensions:

- **Prevalence** measures the frequency of use of the brand name, i.e. the number of times a brand is directly mentioned;

- **Diversity** measures the diversity of the words associated with the brand;

- **Connectivity** represents the brand ability to bridge connections between other words or groups of words (sometimes seen as discourse topics).

*[Fronzetti Colladon, A.: The semantic brand score. Journal of Business Research 88, 150 – 160 (2018)]
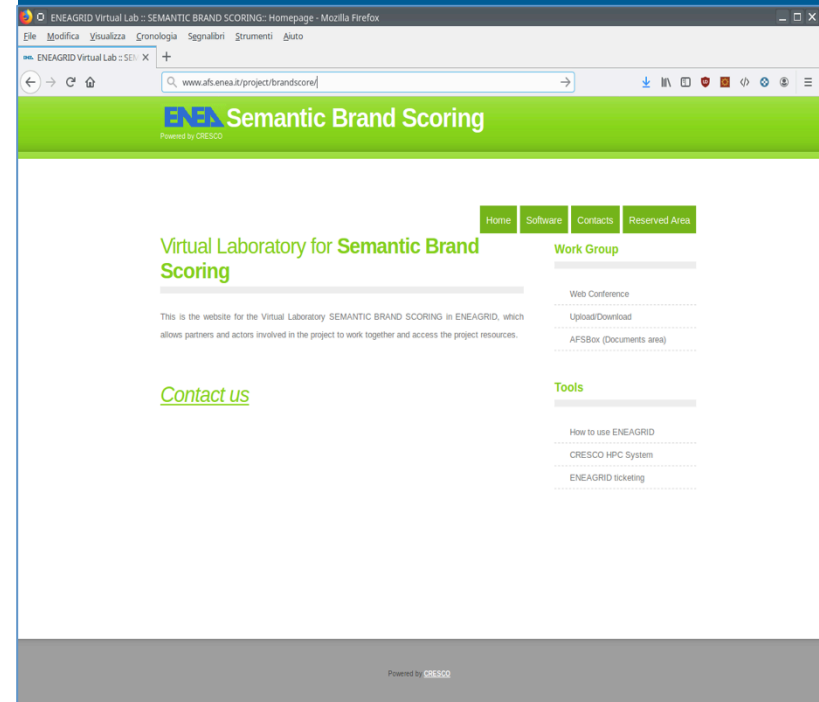
# Semantic Brand Score: Test and Virtual Lab

Preliminary tests on the configuration and on the performance demonstrate a correct integration

## Test: Rome Mayor election



https://semanticbrandscore.com

## Virtual Laboratory
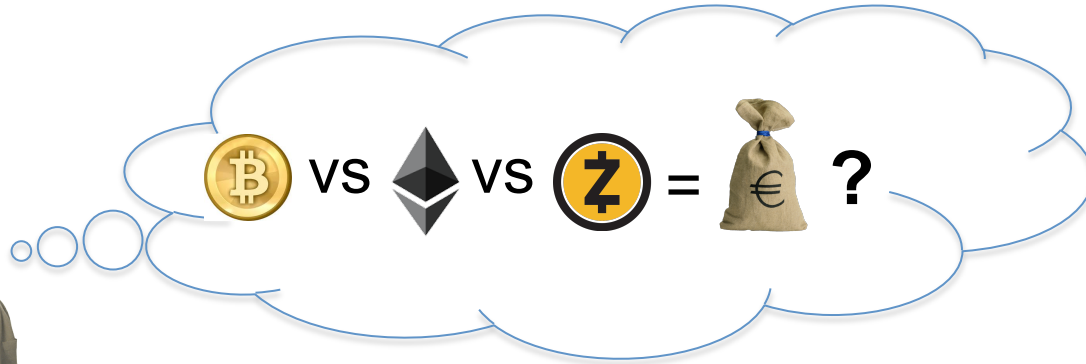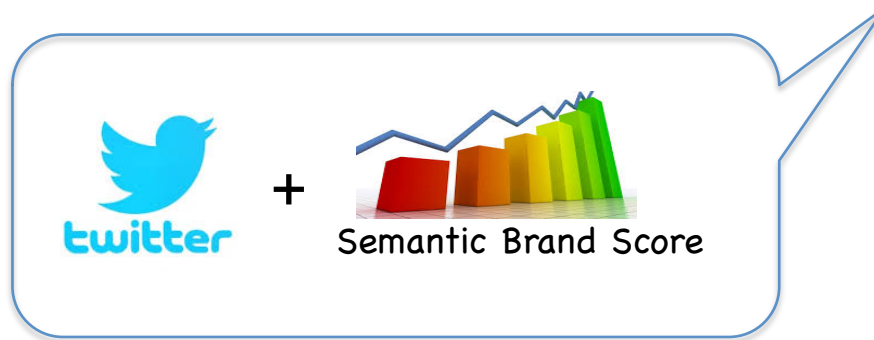


http://www.afs.enea.it/project/brandscore/

Trader

Our idea
of solution

# Current developments

## Social Networks Crawling

- **Aim:**
  - Creating dataset on a specific topic;
  - Downloading texts, messages, news from Social Networks.

- **Methodology:**
  - Installing and configuring a social crawler for *Twitter;*
  - Use of parallel developer accounts to avoid the limitation on the number of tweets downloaded per user.

## Semantic Brand Score for Cryptocurrencies

- **Aim:**
  - Creating a rank among most popular cryptocurrencies (e.g. *Bitcoin*, *Ethereum*, *Zcash*);

- **Methodology:**
  - Running periodic sessions of crawling in order to create a database of tweets that concern news and discussions about digital coins;
  - Applying the Semantic Brand Score to rank cryptocurrency importances.

# Conclusion

✓ We provided an overview of activity about the implementation of web crawler integrated in our HPC ENEAGRID/CRESCO infrastructure;

✓ Currently we are also equipping our framework with a social media crawler that downloads contents from *Twitter* and *a* Semantic Brand Scoring (SBS) tool which uses ENEA computational and storage power;

✓ First tests on the social crawler and on the SBS software demonstrate good results.

## *Future work…*

✐ *Performing experiments to tune our framework;*

✐ *Refining our semantic filter to obtain a more accurate dataset.*

# Thanks for the attention

Dr. Giuseppe Santomauro
giuseppe.santomauro@enea.it