

# VERA: A Platform for Veracity Estimation over Web Data

Mouhamadou Lamine Ba, Laure Berti-Equille, Kushal Shah, Hossam M. Hammady  
Qatar Computing Research Institute  
Hamad Bin Khalifa University  
Doha, Qatar  
{mlba,lberty,kshah,hhammady}@qf.org.qa

## ABSTRACT

Social networks and the Web in general are characterized by multiple information sources often claiming conflicting data values. Data veracity is hard to estimate, especially when there is no prior knowledge about the sources or the claims in time-dependent scenarios (e.g., crisis situation) where initially very few observers can report first information. Despite the wide set of recently proposed truth discovery approaches, “no-one-fits-all” solution emerges for estimating the veracity of on-line information in open contexts. However, analyzing the space of conflicting information and disagreeing sources might be relevant, as well as ensembling multiple truth discovery methods. This demonstration presents VERA, a Web-based platform that supports information extraction from Web textual data and micro-texts from Twitter and estimates data veracity. Given a user query, VERA systematically extracts entities and relations from Web content, structures them as claims relevant to the query and gathers more conflicting/corroborating information. VERA combines multiple truth discovery algorithms through ensembling and returns the veracity label and score of each data value as well as the trustworthiness scores of the sources. VERA will be demonstrated through several real-world scenarios to show its potential value for fact-checking from Web data.

## 1. INTRODUCTION

With the recent development of computational journalism [3, 5], on-line fact-checkers such as FactCheck<sup>1</sup>, Snopes<sup>2</sup>, PolitiFact<sup>3</sup>, TruthorFiction<sup>4</sup> or OpenSecrets<sup>5</sup>, and ClaimBuster<sup>6</sup> have lately gained unprecedented attention as their goal is to verify on-line information for public opinion and automate Web-scale fact-checking. But estimating the veracity of data still remains a challenging problem: extracting structured information from large, heterogeneous

<sup>1</sup><http://www.factcheck.org>

<sup>2</sup><http://www.snopes.com>

<sup>3</sup><http://www.politifact.com>

<sup>4</sup><http://www.truthorfiction.com>

<sup>5</sup><http://www.opensecrets.org>

<sup>6</sup><http://idir-server2.uta.edu/claimbuster>

corpora of textual and multimedia documents, and integrating these multi-source data are difficult tasks. Web data and micro-texts from social media can be noisy, outdated, incorrect, conflicting, and thus unreliable, often due to information extraction errors, disagreements, biased observations, disparate or low quality of the sources.

Many truth discovery methods have been proposed to deal with data veracity estimation (see [2] for a survey). They are mostly applied to structured data and compute iteratively the accuracy of the sources claiming some data as a function of the veracity scores of their data and the veracity scores are computed as a function of the accuracy of their sources. Recent approaches have been developed to discover true values extracted from textual content in a large corpus of Web sources using various information extractors [4, 7]. These solutions extend previous probabilistic models based on iterative vote counting and integrate the extraction systems’ error in truth discovery computation.

Nevertheless, most approaches operate on a static set of structured claims from a fixed corpus of information sources. They usually do not expand dynamically the search space to gather additional evidences and controversial or corroborating claims. Moreover, several studies have proven that a “one-fits-all” solution does not seem to be achievable for a wide range of truth discovery scenarios [6] and we argue that ensembling truth discovery methods can significantly improve the quality performance of current results [1].

In this demo, we present VERA, a Web-based platform that supports the pipeline of truth discovery from Web unstructured corpus and tweets: ranging from information extraction from raw texts and micro-texts and data fusion to truth discovery and visualization. VERA offers several advantages over previous work as it includes:

- Extraction and fusion of multi-source information to answer a factual query defined by the user;
- Combination of multiple truth discovery algorithms using ensembling in order to effectively discover true values from conflicting ones;
- Explanation of the truth discovery results;
- Visualization artifacts to better understand the information space with disagreeing vs. agreeing sources and corroborating vs. conflicting claims.

To the best of our knowledge, this work is the first attempt to demonstrate truth discovery in action from Web data and Twitter data, overcoming limitations of single truth discovery methods with ensembling to estimate data veracity. VERA platform, RESTful API, and additional material including real-world datasets and a synthetic dataset generator are available at: <http://da.qcri.org/dafna/>.

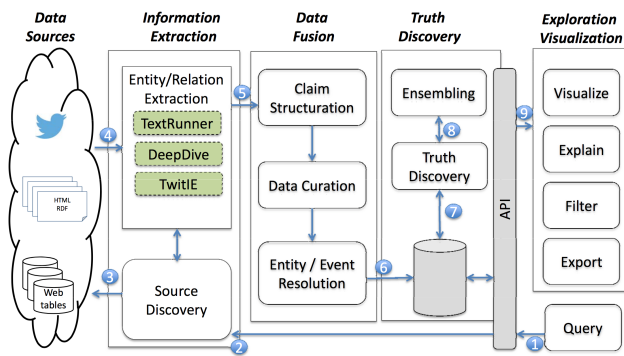


Figure 1: VERA Architecture

## 2. REFERENCES

- [1] L. Berti-Equille. Data Veracity Estimation with Ensembling Truth Discovery Methods. In *IEEE Big Data Workshop on Data Quality Issues in Big Data*, 2015.
- [2] L. Berti-Equille and J. Borge-Holthoefer. *Veracity of Big Data: From Truth Discovery Computation Algorithms to Models of Misinformation Dynamics*. Morgan & Claypool, 2015.
- [3] S. Cohen, J. T. Hamilton, and F. Turner. Computational Journalism. *CACM*, 54(10):66–71, 2011.
- [4] X. L. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmman, S. Sun, and W. Zhang. Knowledge Vault: A Web-scale Approach to Probabilistic Knowledge Fusion. In *KDD'14*, pages 601–610, 2014.
- [5] N. Hassan, C. Li, and M. Tremayne. Detecting Check-worthy Factual Claims in Presidential Debates. In *CIKM'15*, pages 1835–1838, 2015.
- [6] D. A. Waguih and L. Berti-Equille. Truth Discovery Algorithms: An Experimental Evaluation. *CoRR*, 1409.6428, 2014.
- [7] D. Yu, H. Huang, T. Cassidy, H. Ji, C. Wang, S. Zhi, J. Han, C. R. Voss, and M. Magdon-Ismail. The Wisdom of Minority: Unsupervised Slot Filling Validation based on Multi-dimensional Truth-Finding. In *COLING 2014*, pages 1567–1578, 2014.