# Veracity of Big Data: Challenges of Cross-modal Truth Discovery

LAURE BERTI-EQUILLE, MOUHAMADOU LAMINE BA, Qatar Computing Research Institute, Hamad Bin Khalifa University, Doha, Qatar

### From Data Management to Truth Discovery Systems

As online user-generated content grows exponentially, the reliance on Web and social media data is increasing. Truth discovery from the Web has significant practical importance as online rumor and misinformation can have tremendous impacts on our society and everyday life. One of the fundamental difficulties is that data can be biased, noisy, outdated, incorrect, misleading and thus unreliable. Conflicting data from multiple sources amplifies this problem and veracity of data has to be estimated.

Beyond the emerging field of computational journalism and the success of online fact-checkers (e.g., FactCheck[1], ClaimBuster[2]), truth discovery is a long-standing and challenging problem studied by many research communities in artificial intelligence, databases, and complex systems, and under various names: fact-checking, data and knowledge fusion, information trustworthiness, credibility or information corroboration (see [1] for a survey). The ultimate goal is to predict the truth label of a set of assertions claimed by multiple information sources and to infer sources' reliability with no or few prior knowledge. One major line of previous work aimed at iteratively computing and updating the source's trustworthiness as a belief function in its claims, and then the belief score of each claim as a function of its sources' trustworthiness [12]. More complex probabilistic models have then incorporated various aspects beyond source trustworthiness and claim belief such as the dependence between sources, the correlation of claims [8], the notion of evolving truth. Recent contributions have further relaxed prior modeling assumptions to deal with truth existence [13], ap-

---

[1] www.factcheck.org/

[2] idir.uta.edu/claimbuster

proximate truth discovery [11; 5], truth evolution [6; 4], and applications in the context of social media and crowd-sourcing [2; 7]

However, some studies showed that most of prior work suffers from scalability issues, complex parameter setting and non repeatability of the results due to randomized initialization (see [9] for a thorough analysis). Moreover, it is unlikely that one method dominates all others across all application domains and most approaches, relying on majority voting, fall short for pessimistic scenarios where most of the sources are malicious and falsify information. As a consequence, it is currently unclear which techniques are the best suited as they are highly data-dependent and quality performance evaluation depends on the available samples of ground truth data.

In this challenge paper, we argue that the next generation of data management systems need to manage not only volume and variety of Big Data but most importantly veracity of data. Designing end-to-end truth discovery systems requires a fundamental paradigm shift largely driven by Machine Learning advances. It goes beyond adding new layers of data fusion heuristics or developing yet another probabilistic graphical truth discovery model. Actionable, cross-modal, and Web-scale truth discovery is needed in this perspective. It requires a transdisciplinary approach to analyze the dynamic and cross-modal dimensions of rapidly evolving networks of sources and multimedia contents. This paper highlights some of the challenges we deem as the most promising research perspectives towards an actionable, cross-modal, and Web-scale truth discovery.

**Cross-modal and Cross-lingual Truth Discovery.** The agility of a truth discovery system is of utmost importance to efficiently extract and map information: (i) from various languages; (ii) in various data formats, structures, and semantics (e.g., texts, Web tables, structured data, etc.); and (iii) conveyed by various media and technologies (e.g., tweets, Instagram images, Youtube videos, Web pages, RSS feeds, etc.). The main challenge is to address cross-modality and cross-language issues in the context of truth discovery where the linkage of various evidences from audio, image, video, text formats and languages has to be achieved as accurately as possible to corroborate events. This refers to cross-media entity and event linking where Deep Learning has just started to bring some interesting solutions and can be leveraged for cross-modal truth discovery.

**Timely and Actionable Truth Discovery.** In a humanitarian context for example, truth discovery from quasi real-time data could save lives. To be actionable, information extraction and truth discovery computation need to be streamlined, prioritized depending on the level of emergency and incompleteness of available information, and finally adjusted to the communities that will use the data (e.g., rescue team, NGOs). Time-dependent estimation and correction of observer bias, selection bias and long-tail phenomenon problem (e.g., where very few sources provide the first information after a disaster) are challenging tasks for quasi real-time truth discovery.

**Estimation of Incompleteness, Biases and Errors in the Truth Discovery Process.** Information without context can be easily distorted and misinterpreted. When a piece of information is extracted from its original content, channel or thread, it may lose its context along with important "semantic markers" that explain *when, where, how, why,* and *for which* purpose or audience it has been produced. Observation may also be incomplete and biased for various reasons, e.g., security and privacy concerns, format limitations, *observer's bias* or *disclosure bias*. Estimating the biases and errors along the entire truth discovery pipeline is crucial and challenging, as well as estimating the credibility of user-generated content [3].

To overcome these challenges, we believe that an integrative framework is needed: (i) To define, in a principled way, a unified semantics of truth discovery; (ii) To proactively collect new evidences, contextual data, and external knowledge from multi-modal data; (iii) To support continuous inference and belief revision for computing and updating data veracity estimates; (iv) And finally, to monitor and estimate errors and biases in the truth discovery process.

To address these challenges, we have proposed DAFNA (*Data Forensics with Analytics* – dafna.qcri.org) at QCRI, an ambitious project for determining the veracity of cross-modal information from multiple Web sources. Beyond a first module demonstrated in [10], DAFNA's vision is to provide a platform for actionable and cross-modal truth discovery but still significant work is needed to address realistic and multi-modal truth discovery scenarios.

## REFERENCES

[1] L. Berti-Equille and J. Borge-Holthoefer. *Veracity of Data: From Truth Discovery Computation Algorithms to Models of Misinformation Dynamics*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2015.

[2] J. Gao, Q. Li, B. Zhao, W. Fan, and J. Han. Truth discovery and crowdsourcing aggregation: A unified perspective. *Proc. of the VLDB Endowment*, 8(12):2048–2059, 2015.

[3] G. Haralabopoulos, I. Anagnostopoulos, and S. Zeadally. The Challenge of Improving Credibility of User-Generated Content in Online Social Networks. *ACM JDIQ*, 2016.

[4] L. Jia, H. Wang, J. Li, and H. Gao. Incremental truth discovery for information from multiple data sources. In *Proc. of the Web-Age Information Management (WAIM) 2013 International Workshops*, pages 56–66, 2013.

[5] Q. Li, Y. Li, J. Gao, L. Su, B. Zhao, M. Demirbas, W. Fan, and J. Han. A confidence-aware approach for truth discovery on long-tail data. *Proc. of the VLDB Endowment*, 8(4):425–436, 2014.

[6] Y. Li, Q. Li, J. Gao, L. Su, B. Zhao, W. Fan, and J. Han. On the discovery of evolving truth. In *Proc. of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 675–684, 2015.

[7] F. Ma, Y. Li, Q. Li, M. Qiu, J. Gao, S. Zhi, L. Su, B. Zhao, H. Ji, and J. Han. Faitcrowd: Fine grained truth discovery for crowdsourced data aggregation. In *Proc. of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 745–754, 2015.

[8] R. Pochampally, A. Das Sarma, X. L. Dong, A. Meliou, and D. Srivastava. Fusing data with correlations. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data (SIGMOD'14)*, pages 433–444, 2014.

[9] D. A. Waguih and L. Berti-Equille. Truth Discovery Algorithms: An Experimental Evaluation. QCRI Technical Report, May 2014.

[10] D. A. Waguih, N. Goel, H. M. Hammady, and L. Berti-Equille. AllegatorTrack: combining and reporting results of truth discovery from multi-source data. In *Proc. of ICDE*, pages 1440–1443, Seoul, Korea, 2015. IEEE.

[11] X. Wang, Q. Z. Sheng, X. S. Fang, X. Li, X. Xu, and L. Yao. Approximate truth discovery via problem scale reduction. In *Proc. of the 24th ACM Conference on Information and Knowledge Management (CIKM)*, October 2015.

[12] X. Yin, J. Han, and P. S. Yu. Truth Discovery with Multiple Conflicting Information Providers on the Web. *IEEE Trans. Knowl. Data Eng.*, 20(6):796–808, 2008.

[13] S. Zhi, B. Zhao, W. Tong, J. Gao, D. Yu, H. Ji, and J. Han. Modeling truth existence in truth discovery. In *Proc. of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1543–1552, 2015.