

معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

عضو في مؤسسة قطر
Member of Qatar Foundation

DISCOVERING THE TRUTH ON THE WEB One Facet of Data Forensics

Mouhamadou Lamine Ba, Laure Berti-Equille, Hossam M. Hammady

Qatar Computing Research Institute, HBKU, Doha, Qatar

{mlba, lberti, hhammady}@qf.org.qa

Truth Discovery Use Cases

Query Answering: Dealing with many possible answers

Example query: How many dead in Paris terrorist attacks?

Claim	Source	Value	Truthfulness
c_5	cnn.com	At least 128	Unknown
c_4	theguardian	120	Unknown
c_3	news.sky.com	130	Unknown
c_2	bbc.com	130	Unknown
c_1	@TBurgesWatson	35	Unknown

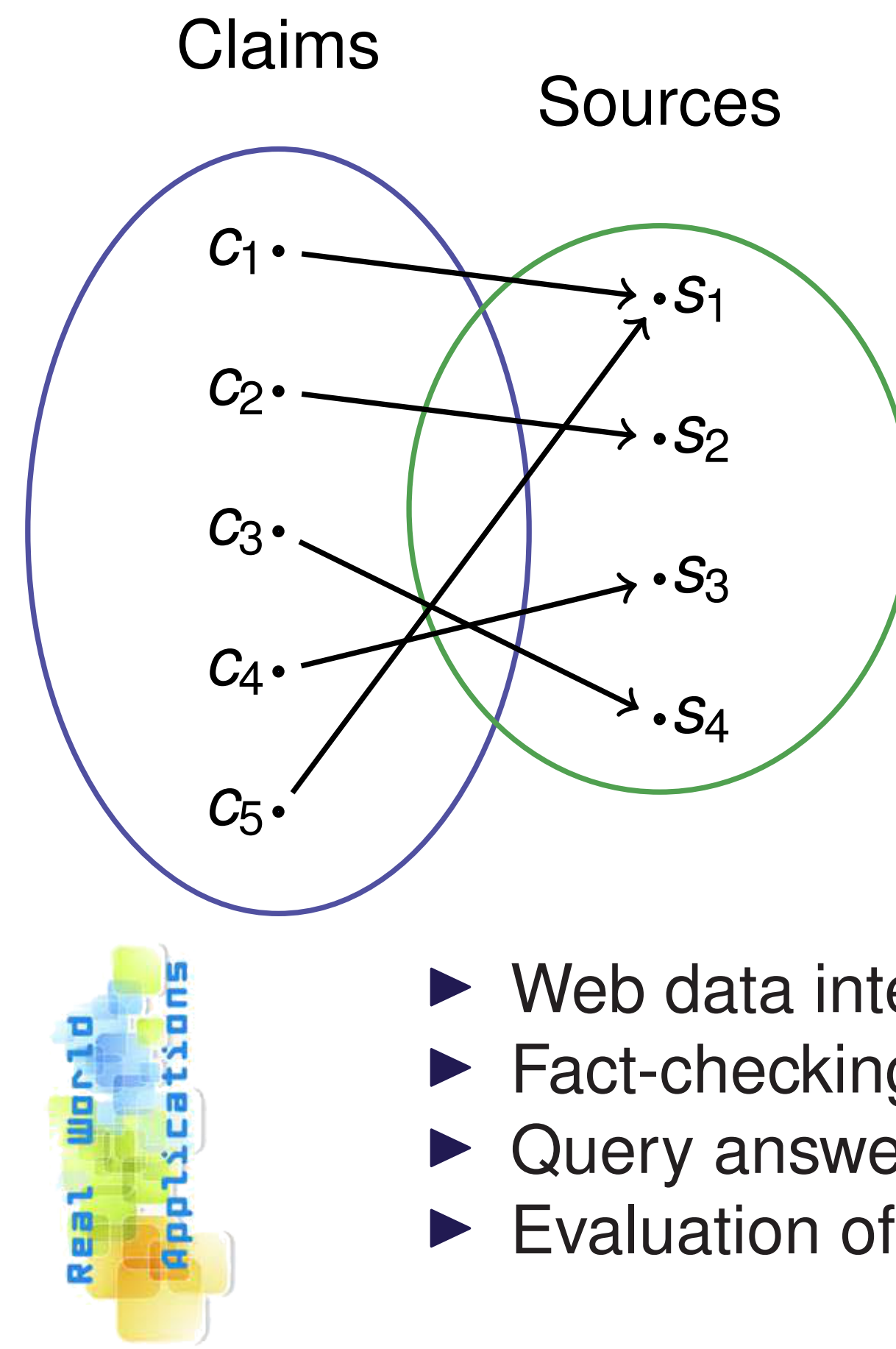
Multi-Source Data Fusion: Resolving conflicting values

Example: Percentage of expats in Qatar

S1: City	% expats	S2: City	% expats	S2: City	% expats
Qatar	75%	Qatar	80%	Qatar	60%
		City	% expats		
		Qatar	?		

Finding the most relevant results in both use cases requires the estimation of the **truthfulness level** of each data claimed by the sources.

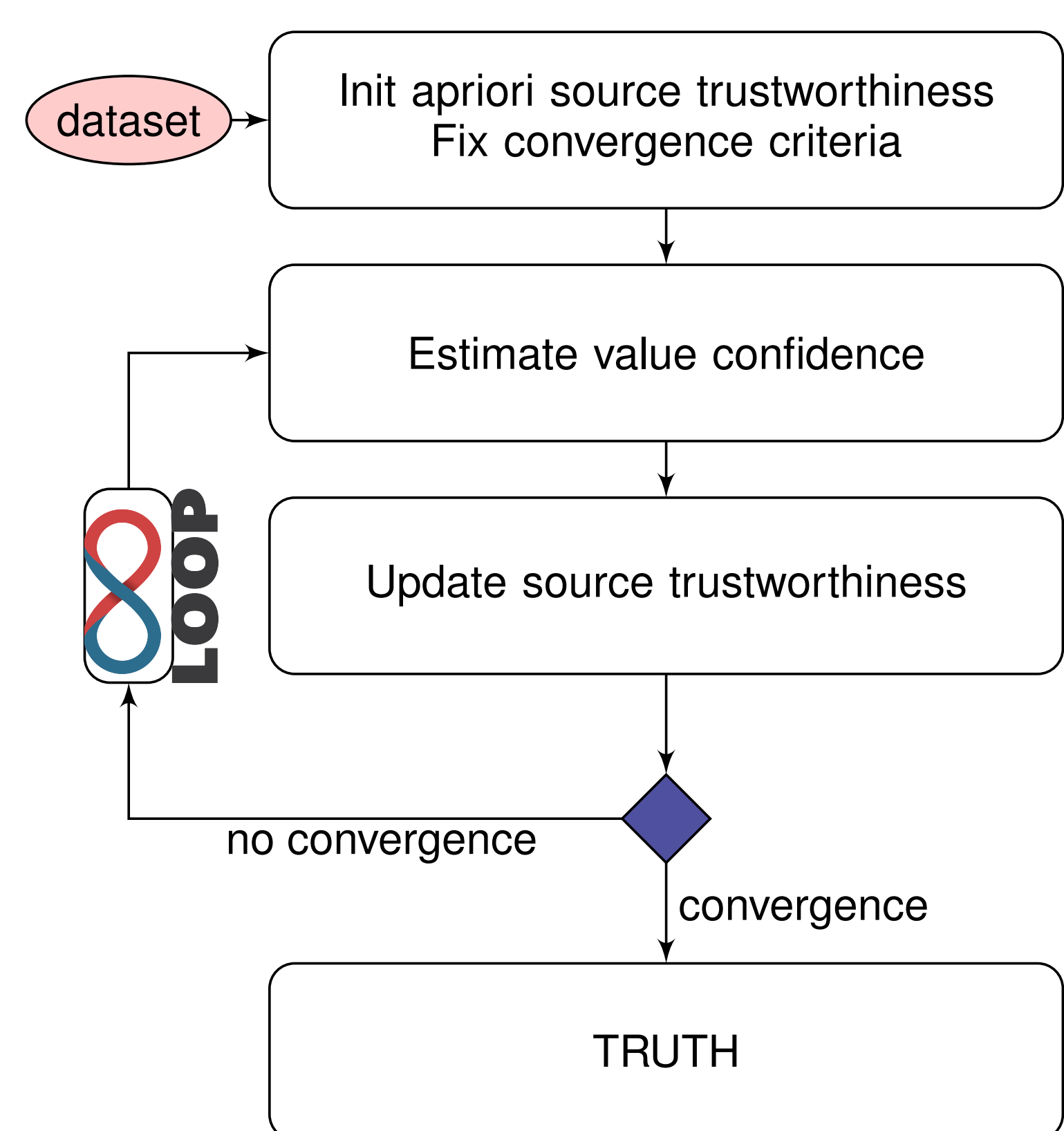
Problem statement: Given a set of claims $\{c_1, c_2, \dots, c_n\}$ about a real-world fact f claimed by n information sources, truth discovery computes a mapping $\mathcal{F} : \{c_1, c_2, \dots, c_n\} \mapsto \{\text{True}, \text{False}\}$



$c_1 \mapsto \underline{FALSE}$
 $c_2 \mapsto \underline{FALSE}$
 $c_3 \mapsto \underline{TRUE}$
 $c_4 \mapsto \underline{FALSE}$
 $c_5 \mapsto \underline{FALSE}$

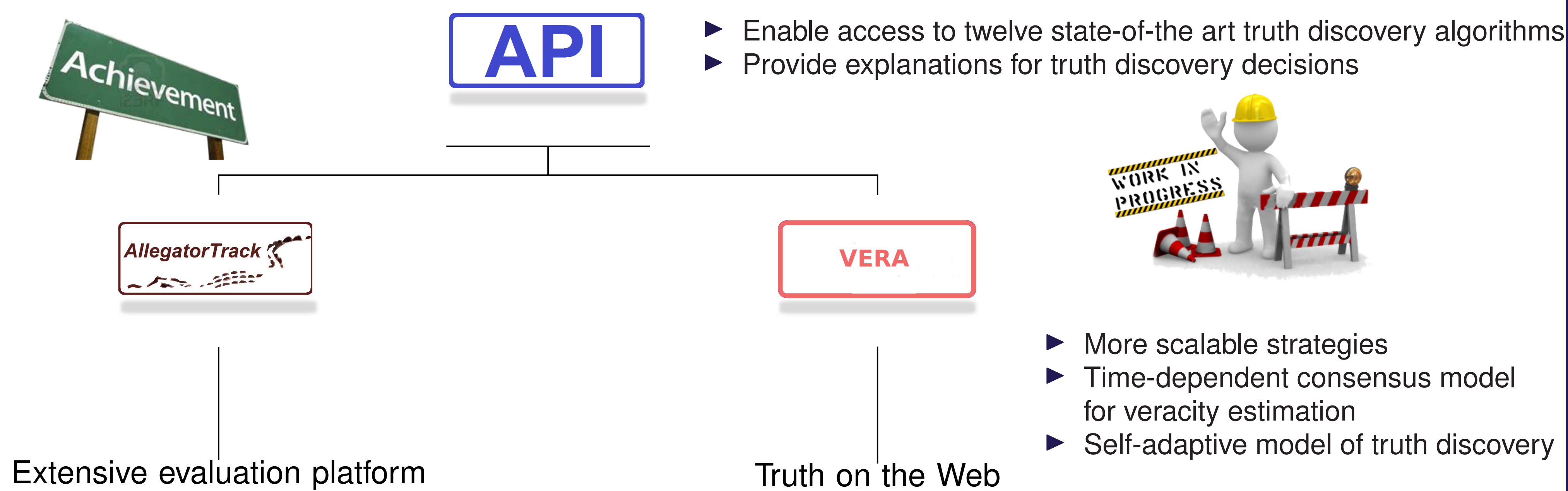
- Web data integration
- Fact-checking in computational journalism
- Query answering in Web open information extraction
- Evaluation of tasks and workers in crowd-sourcing apps

Truth Discovery Pipeline



DAFNA – Data Forensics with Analytics

Data Forensics with Analytics (DAFNA) is a project initiated by the Data Analytics group at QCRI. DAFNA aims at providing a suite of tools for estimating data veracity for Data Forensics and solving some of the limitations of current truth discovery methods.



- More scalable strategies
- Time-dependent consensus model for veracity estimation
- Self-adaptive model of truth discovery

Three Classes of Algorithms

One can classify existing truth discovery approaches in three classes.

- Agreement-based algorithms
- MAP Estimation-based algorithms
- MLE Estimation-based algorithms

MAP (Max. a posteriori)
 Setup: • Given data $D = (x_1, \dots, x_n)$
 • Assume a joint dist $p(D, \theta) = p(D|\theta)p(\theta)$
 • Goal: Choose a good value of θ for D
 • Check $\theta_{MAP} = \arg\max_{\theta} p(\theta|D)$
 $\theta_{MLE} = \arg\max_{\theta} p(D|\theta)$
 Pros: Easy & Interpretable

Known Limitations

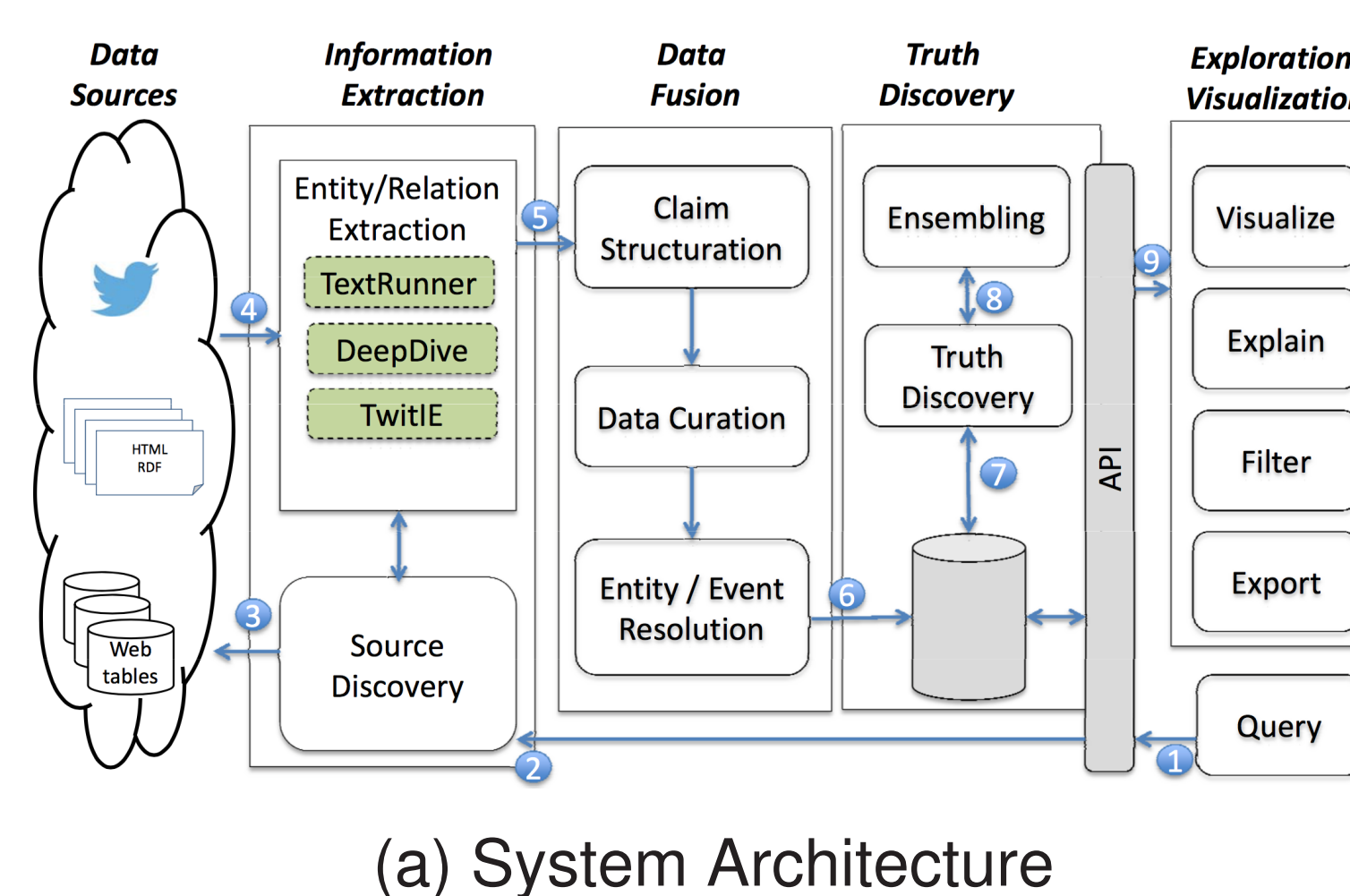
Applicability of truth discovery algorithms on the Web faces the following challenges.

- No one-fits-all solution
- No common evaluation platform
- No query-specific claim extraction
- Limited support of rapidly evolving truth

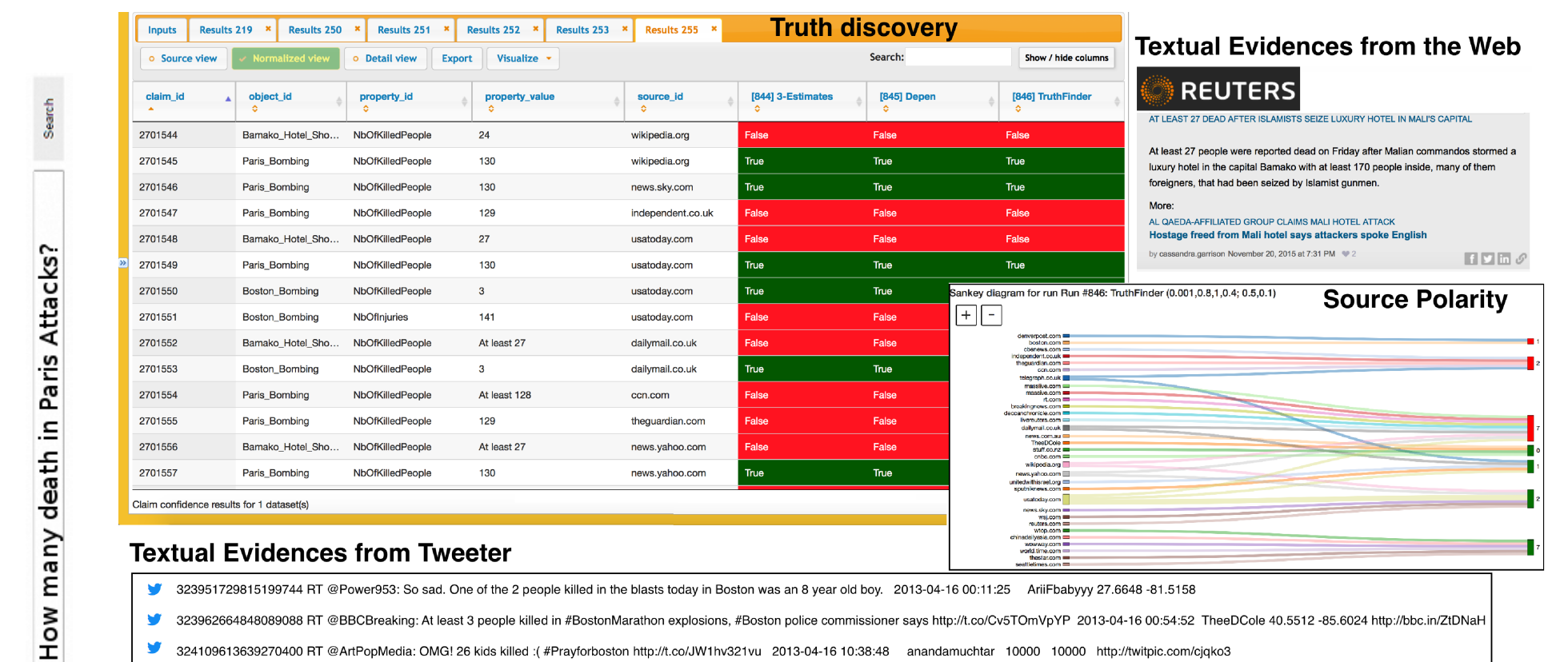
VERA – Veracity Estimation on Web Data

VERA is a Web-based platform supporting 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;
- Use of time-dependent model in order to capture evolving truth, in particular from social media;
- Ensembling multiple truth discovery models to effectively find true values from conflicting ones;
- Visualization artifacts to better understand the information space with disagreeing vs. agreeing sources and corroborating vs. conflicting claims.



(a) System Architecture



(b) VERA Back-End and Visualization Artifacts

Publications

- Mouhamadou Lamine Ba, Laure Berti-Equille, Kushal Shah, Hossam M. Hammady: VERA – A Platform for Veracity Estimation over Web Data. In World Wide Web (WWW), 25th International World Wide Web Conference, April 2016.
- Laure Berti-Equille and J. Borge-Holthoefer: Veracity of Big Data – From Truth Discovery Computation Algorithms to Models of Misinformation Dynamics. Lectures on Data Management, Morgan & Claypool Publishers, 2015.
- 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 Data Engineering (ICDE), 2015 IEEE 31st International Conference on Data Engineering, pages 1440 –1443, April 2015.
- Dalia Attia Waguih and Laure Berti-Equille. Truth discovery algorithms: An experimental evaluation. CoRR, abs/1409.6428, 2014.