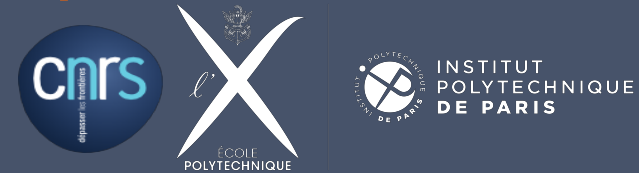# SourcesSay: Intelligent Analysis and Interconnexion of Heterogeneous Data in Digital Arenas

AI Chair project, ANR & DGA

**Ioana Manolescu**

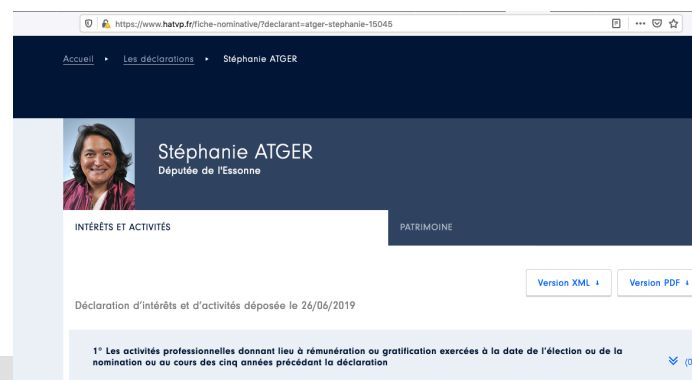**Inria and Institut Polytechnique de Paris**

# Motivation

**Data production has been democratized**: unprecedented data generation rates by humans, software, and (equipped) physical objects

Numerous opportunities to **add value by integrating data from several sources**. Examples from data journalism:

- Follow **official communication** by politicians together with their **social media presence**, **laws** they promote, and their **conflicts of interest**

SourcesSay project (ANR and DGA), Ioana Manolescu, Inria and Institut Polytechnique de Paris          AI Chair Kick-Off, September 2020          *Inría*

# Why data journalism?

Because I grew up in a dictatorship, and I value free press

Because journalists are threatened and killed still today in Europe



Daphne Galizia, 1964-2017



Jan Kuciak, 1990-2018

Because the press' economic model is threatened by IT giants

Because this industry is currently underserved by IT – and we could really make an impact!

# Data journalism problem: working with heterogeneous data

Digital data sources are **heterogeneous**



- For Open Data, W3C standard advocates RDF. Yet...

- **INSEE**: some RDF, lots of Excel and HTML**; NosDéputés.fr**: JSON, XML

- **HATVP** (Haute Autorité pour la Transparence de la Vie Publique): CSV, XML

- **EFSA** (European Food Safety Administration): PDF

Different format, organization, structure, value representation convention...

# Application: analyzing a fake news ecosystem

Fact-checking: verification of public statements in the (social) media

Collaboration since 2014 with:

**Le Monde**

**LES DÉCODEURS**
VENONS-EN AUX FAITS

Les Décodeurs publish as Open Data their classification of 1300 web sites in:

{ **rather reliable**; **satirical**; **has published fakes**; **agregateur (re-check)** }

https://www.lemonde.fr/webservice/decodex/updates
https://toolbox.google.com/factcheck/

**Google** Fact Check Tools

# Application: analyzing a fake news **arena**

An **arena** consists of a set of **entities** (users, organizations etc.) and **contents** they author, share, or are mentioned in

**Fake news arena**:

- content: HTML, JSON, text or PDF (pages, articles, posts, tweets)

- publishers and distributors, e.g., in relational or JSON;

- fact-checks (usually semi-structured, XML or JSON)

**Given** a new content with their authors, re-distributors, links etc.

**What can we say** about the trustworthiness of the content and its environment?

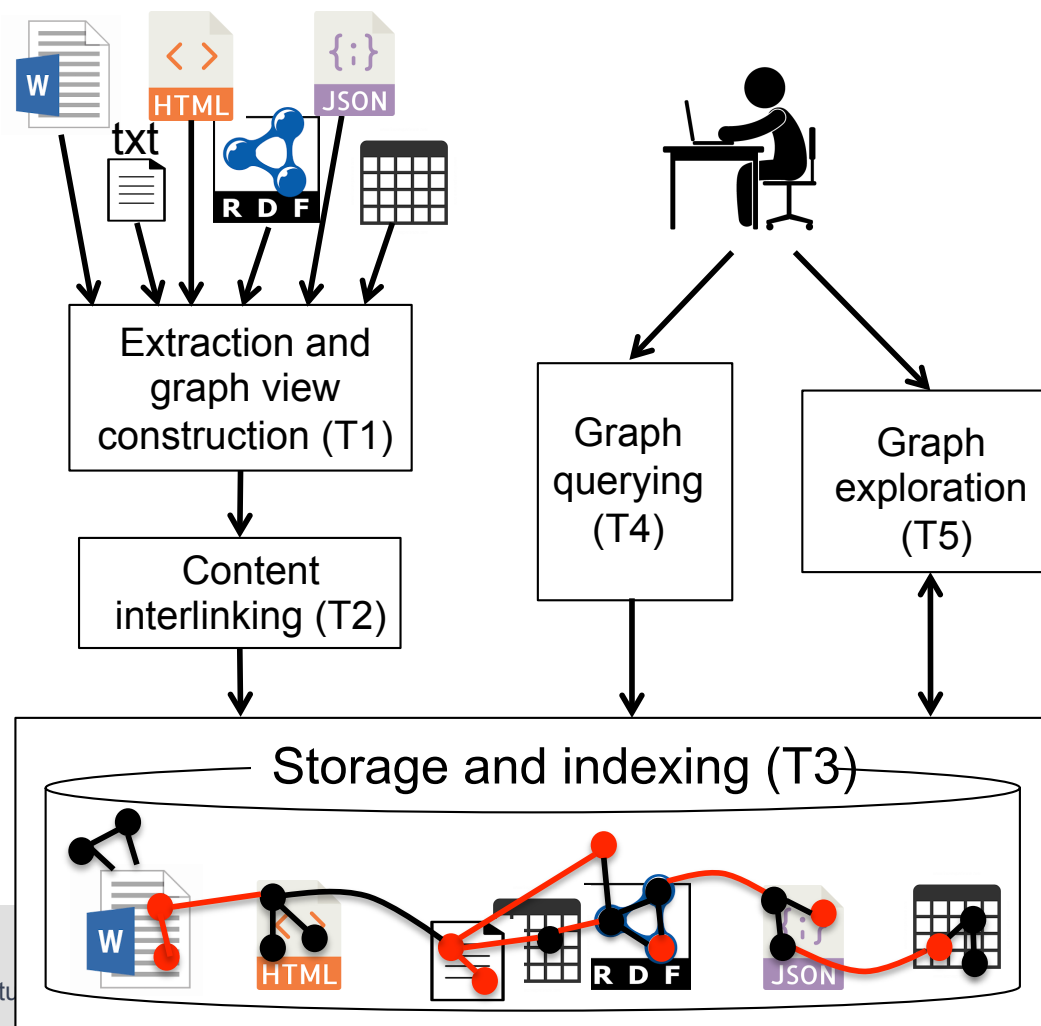     AI Chair Kick-Off, September 2020     *Inria*

# Other digital arenas

- **Scientific and general-audience publications** on a topic

  - Particle air pollution, controversial drugs, a company's products...
- **Journalistic investigations**

  - Tax evasion (Panama Papers):
    relational database + PDF documents
  - Mongering doubt on tobacco effects
    or global warming
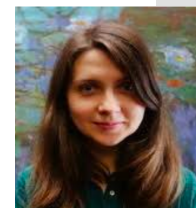
# SourcesSay Architecture

Unifying technical hypothesis:
integrate data in a **graph**



SourcesSay project (ANR and DGA), Ioana Manolescu, Inria and Institu...

# Challenges

How to **enrich and interconnect sources**?        Collab. O. Balalau (CEDAR)
                                                                                    H. Galhardas (U. Lisbon)

- Entity and relationship extraction through NLP and AI

- Node/entity matching, disambiguation w/r knowledge base...

How to **efficiently store large volumes of heterogeneous content, and the connections** extracted from them?                    Collab. A. Anadiotis (CEDAR)
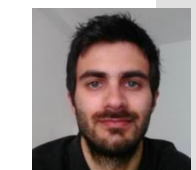
- Learn usage patterns, devise novel data processing engines

How to **efficiently and flexibly query the data** using keywords or NL?
How to know when an answer is **interesting**?        Collab. O. Balalau (CEDAR)

- Learn to rank

How to **explore and interact with** the graphs?        Collab E. Pietriga (ILDA)

*Inria*

# Context and collaborations



- Non-funded partners bring applications:
  Le Monde, WeDoData

- Inria AI engineer (2019-2022): Tayeb Merabti

- DIM RSFI PhD (2020-2022)  with WeDoData

- Work started in seven 2020 internships: I. Burger, F. Chimienti,
  J. Feitz, Y. Haddad, J. You, Y. Youssef, X. Zhang

SourcesSay project (ANR and DGA), Ioana Manolescu, Inria and Institut Polytechnique de Paris          AI Chair Kick-Off, September 2020          *Inria*

# The vision

- **Integrate** and **interpret** heterogeneous data from digital arenas

- **Sourcing** precisely every info

- Novel **query** and **exploration**

- Precision, efficiency, friendliness to non-technical users

     AI Chair Kick-Off, September 2020     *Inria*