# From Fact-Checking to Journalistic Data Integration

Ioana Manolescu
Inria Saclay-Île-de-France and Institut Polytechnique de Paris

# Outline

1. Group presentation

2. Motivation: why journalism?

3. Towards automatic checking of statistic claims

4. Integrating (very!) heterogeneous journalistic data
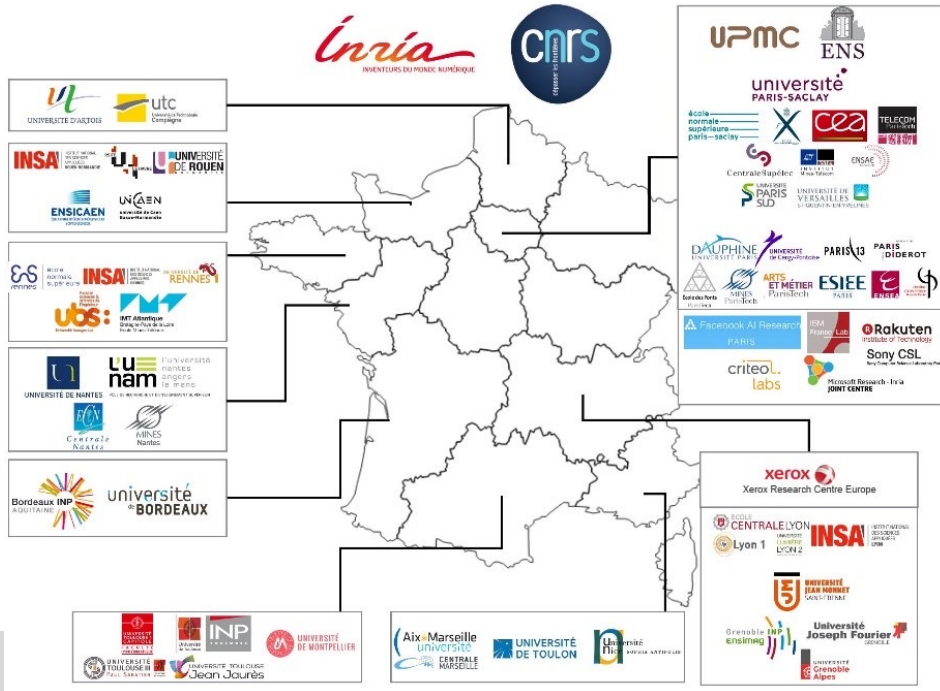
5. Related work and perspective

Ioana Manolescu, Inria and Institut Polytechnique de Paris        Trustworthy Data Science and AI, SFU   April 2021

*Inria*

# CEDAR team presentation

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

# CEDAR: Inria and Ecole Polytechnique

**Inria**: French national research institute in Computer Science and Applied Mathematics, since 1976

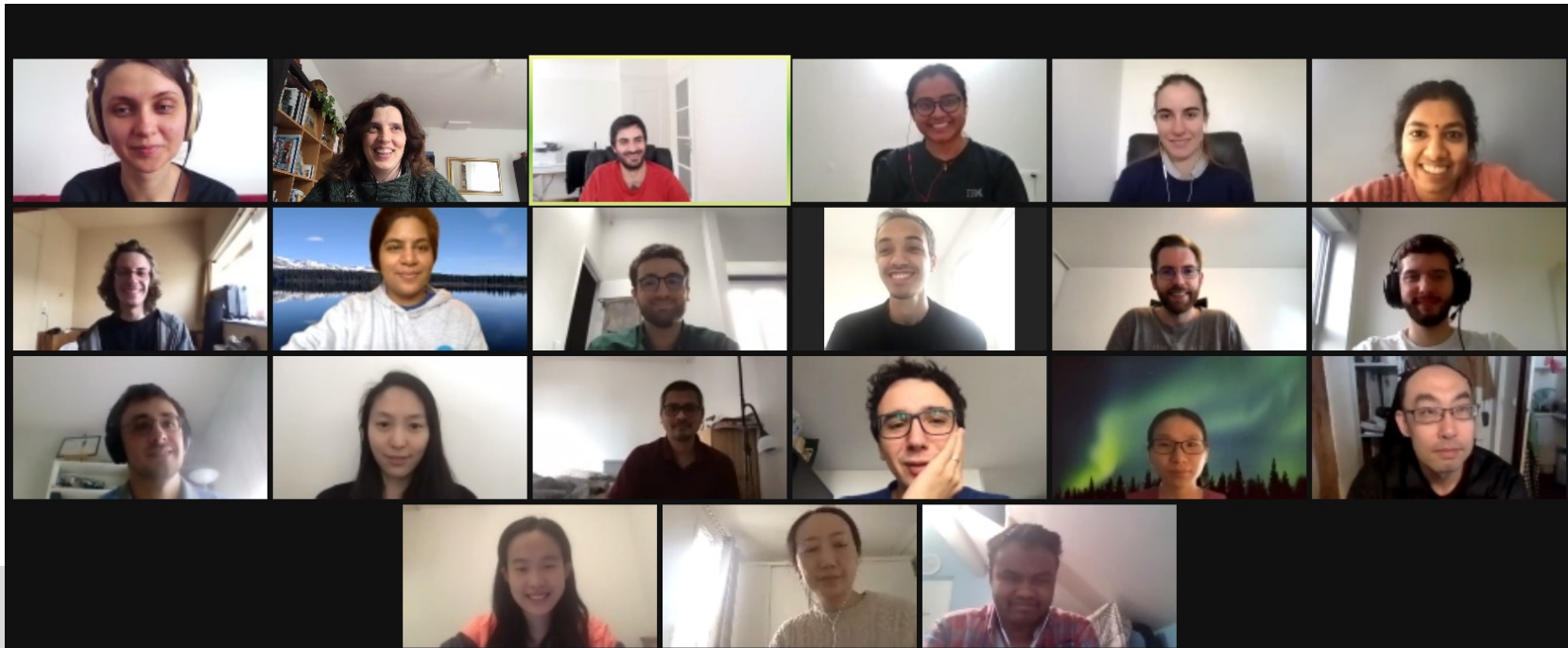**Ecole Polytechnique**: created in 1794 to train military engineers

Among the professors: Ampère, Fourier, Monge, Laplace, Cauchy, Becquerel; 2018 Nobel prize in physics…

# CEDAR team: Inria and Polytechnique

Created in **2016** (I. Manolescu moved from U. Paris Sud, Yanlei Diao from U. Massachussets at Amherst)

Junior faculty: Angelos Anadiotis (EPFL), Oana Balalau (MPI Saarbrucken)

# Motivation:
# Why journalism?

# Democratic societies crucially need the press

❑ To debate and express dissent


Socialist Romania, 1984

❑ To analyze, confim or refute public statements — Fact-checking

(Data) journalism

❑ To expose and explain society functioning



Ioana Manolescu, Inria and Institut Polytechnique de Paris        Trustworthy Data Science and AI, SFU   April 2021

# Fact-checking

Not everyone agrees, however, that Democrats are not flip-flopping on the issue.

Mark Krikorian, executive director of the Center for Immigration Studies, a think tank that advocates for lower immigration, said that because the public doesn't know exactly what border barriers the Trump administration wants to build, Mulvaney's statement is not an "exact" comparison. But, he said, to dismiss it simply on that basis would be "tendentiously literal."

"The fact is that, other than the 'Mexico will pay for it' stuff, Trump is simply channeling the 2006 Secure Fence Act, and Schumer et al. who voted for it out of political calculation are indeed hypocrites for opposing the attempt to finally bring that law to fruition," Krikorian told us via email.

At the surface level, it is true in a broad sense that Democrats including Schumer, Obama and Clinton have in the past supported border fencing. All three voted for the Secure Fence Act of 2006, and all three supported the 2013 Senate immigration overhaul that passed the Senate, and which called for tougher border security including some additional fencing. But to claim that those measures are the same as what Trump is proposing is a stretch.

Share The Facts

**Mick Mulvaney**
Director, Office of Management and Budget

MISLEADING

"We don't understand why the Democrats are so wholeheartedly against [President Trump's border wall]. They voted for it in 2006."

FactCheck.org

Fox News Sunday – Sunday, April 23, 2017

SHARE    READ MORE

---

www.factcheck.org/2017/04/democrats-support-border-wall/

**FACTCHECK.ORG** A Project of The Annenberg Public Policy Center

HOME    ARTICLES    ASK A QUESTION    VIRAL SPIRAL    ARCHIVES    ABOUT US    SEARCH    MORE

THE WIRE

## Did Democrats Once Support Border Wall?

By Robert Farley   Posted on April 26, 2017

Like 835    Tweet    Pin it    Share    11

White House Office of Management and Budget Director Mick Mulvaney made an apples-to-oranges comparison when he said he couldn't understand why Democrats opposed supplemental funding for a border wall since many of them were for it back in 2006.

Mulvaney is referring to the Secure Fence Act of 2006, which called for construction of 700 miles of fencing and enhanced surveillance technology, such as unmanned drones, ground-based sensors, satellites, radar coverage and cameras. Sen. Chuck Schumer and then-Sens. Barack Obama and Hillary Clinton were among a bipartisan majority that voted in favor of the legislation, and it was signed into law by President George W. Bush.

In a very general sense, the Democrats named by Mulvaney supported a bill to build more
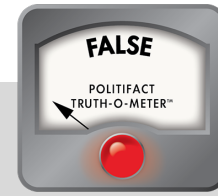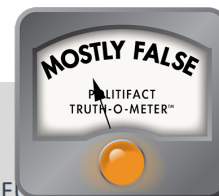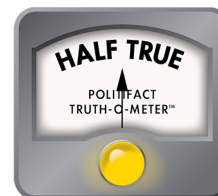
✓ ASK FACTCHECK

Like 953    Tweet    Pin it

Share    98

**Q:** Did the Supreme Court rule that public schools cannot teach students about Islam?

**A:** No. That false claim was spread by a network of fake news websites.

# Data journalism

Panama Papers (International Consortium of Investigative Journalism, ICIJ)

# Projects and collaborations

**Google Award** (2015) with X. Tannier (U. Paris Sud)

**ANR ContentCheck** (2016-2020) with
Sorbonne Université, U. Lyon, U. Rennes 1,
Les Décodeurs (Le Monde) https://contentcheck.inria.fr

**Inria Associated Team WebClaimExplain** (2017-2019),
with AIST Japan (Julien Leblay)

Collaboration with H. Galhardas (University of Lisbon),
A. Anadiotis, O. Balalau (CEDAR), E. Pietriga (ILDA)

**ANR SourcesSay AI Chair** (2020-2023),
with Le Monde and WeDoData https://sourcessay.inria.fr

# Towards automated fact-checking of statistic claims

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

*Inria*

# Most common fact-checking scenarios

❑ "What is the value of metric X in space Y at time T"?

  ❑ X=youth unemployment, Y=Germany, T=2018
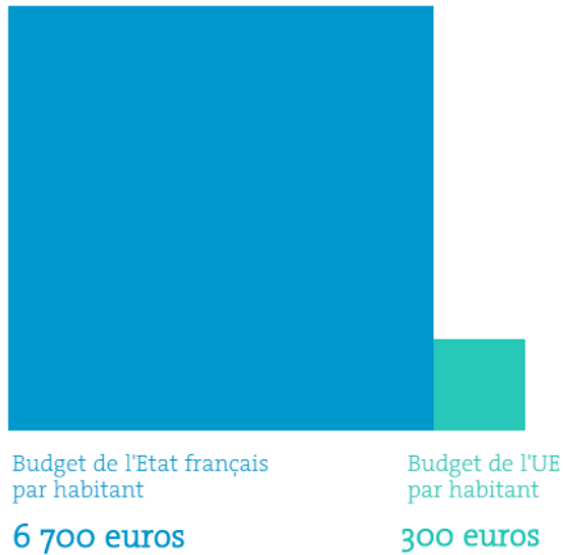  ❑ X=illegal immigrants, Y=Italy, T=[2015-2018]
  ❑ X=budget for research, Y=Canada, T=2020

❑ Comparisons

  ❑ X1 against X2; Y1 against Y2; T1 against T2; temporal trend etc.

# Most common fact-checking scenarios

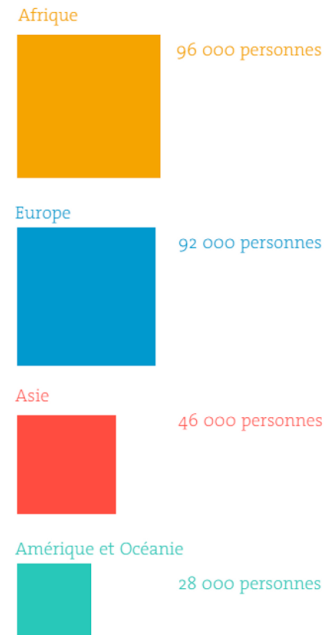**Le budget européen par habitant pèse nettement moins que celui de la France**

Budgets pour l'année 2018 de l'UE et de la France rapportés à leurs populations respectives.

Budget de l'Etat français par habitant
**6 700 euros**

Budget de l'UE par habitant
**300 euros**

Ⓓ LeMonde.fr/lesdecodeurs

**Parmi les immigrés en France, presque autant d'Européens que d'Africains**

Lieu de naissance des personnes entrées sur le territoire national en 2017

Afrique
96 000 personnes

Europe
92 000 personnes

Asie
46 000 personnes

Amérique et Océanie
28 000 personnes

Source : Insee

Ⓓ LeMonde.fr/lesdecodeurs

*Inria*

# Fact-checking as a content management problem



Media content → Claim to be checked (text or data)

User (journalist, expert, crowd worker)

Media context → Verification tool (query, match, source search…) → Analysis result + proof « True / rather true / rather false / false

See sources: http://dataref.com… »

Reference information source 1

Reference information source 2

…

Reference information source n

**[WWW2018]** "A Content Mgmt. Perspective on Fact-Checking" **[WWW2018, VLDB 2018 tutorial]** "Computational FC: problems, state of the art, and perspectives"

# Facilitating **statistical** fact-checking

**INSEE**: French national institute of statistics

❑ Publish valuable statistic datasets about economy, health, education etc. yearly or per quarter

❑ Web pages + statistic information as tabular files (mostly Excel)

**UN, OCDE, IMF**: SDMX databases

❑ Special Data Warehouse-style format for describing the data
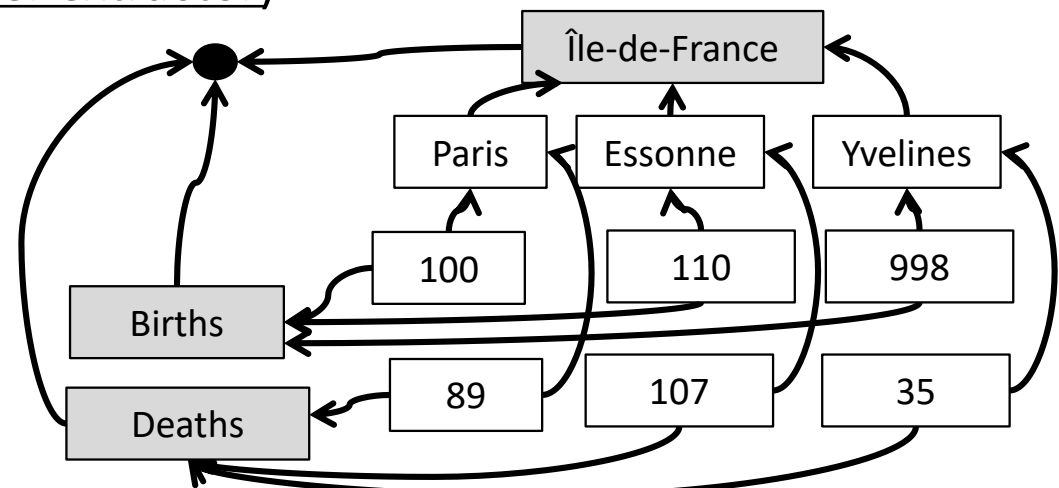
How to facilitate using them as reference data sources?

# Fact-checking using INSEE data [SBD2017]

1. Crawled complete INSEE publication site nightly, gather HTML+XLS files (https://gitlab.inria.fr/cedar/insee-crawler)

2. Extract data from all statistic cells into RDF, preserving the connections between the cells (https://gitlab.inria.fr/cedar/excel-extractor)

| | Ile-de-France | | |
|---|---|---|---|
| | Paris | Essonne | Yvelines |
| Births | 100 | 110 | 98 |
| Deaths | 89 | 107 | 35 |



Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

# Fact-checking using INSEE data [WebDB2018]

**3. Keyword search** algorithm on resulting RDF graph which, given *"Créations d'entreprises en France en 2015"* , returns:

- ❑ Searches for line and column
- ❑ Returns **cell at their intersection**, if possible
- ❑ Otherwise, column or line
- ❑ Otherwise, spreadsheet
- ❑ Always with provenance (link to INSEE Web site)

Started similar project on SDMX [NLIWOD2020]

**Créations d'entreprises dans quelques pays de l'Union européenne en 2015**

*en %*

| Pays | Taux de création |
|------|------------------|
| Allemagne | 7,1 |
| Belgique | 6,2 |
| Espagne | 9,5 |
| France (1) | 9,5 |
| Italie | 7,5 |
| Pays-Bas | 10,1 |
| Portugal | 15,7 |
| République tchèque | 8,2 |
| Royaume-Uni | 14,3 |

# Fact-checking using INSEE data [NLDB2019]

4. **Statistic claim extraction** algorithm which reads incoming tweets and identifies: Measure, Quantity, [Where], [When], [Other dimensions]

- ❑ Formulates keyword query with Measure, [Where], [When], [Other dimensions] for the INSEE search algoritm

- ❑ Leave it to the user to decide how to interpret the difference between claimed and found values

*Inría*

# Integrating heterogeneous journalistic datasets

https://team.inria.fr/cedar/connectionlens/

*Inria*

# A project discussed with Les Décodeurs: fake new detection and propagation on Twitter

**Online fact-checks**: (semi)structured data sources (JSON, XML) listing

- Link to claim (media, social network etc.), **claim author**

- **Fact-check**, containing: analysis (details), final assessment, fc author, date, institution



Among the first published: https://www.lemonde.fr/webservice/decodex/updates

Years later: **ClaimReview** by Google and others (https://www.claimreviewproject.com/)

Inria

# A project discussed with Les Décodeurs: fake new detection and propagation on Twitter
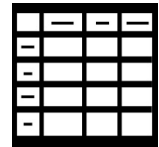
**Online fact-checks**: (semi)structured data sources (JSON, XML) listing

- Link to claim (media, social network etc.), **claim author**

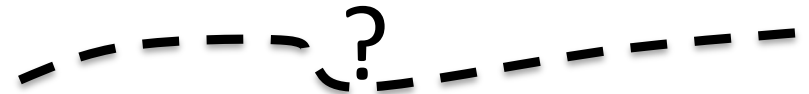- **Fact-check**, containing: analysis (details), final assessment, fc author, date

**Décodeurs'** database of French public figures (Excel)

- First name, last name, Twitter ID, position, political party when known
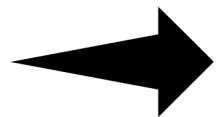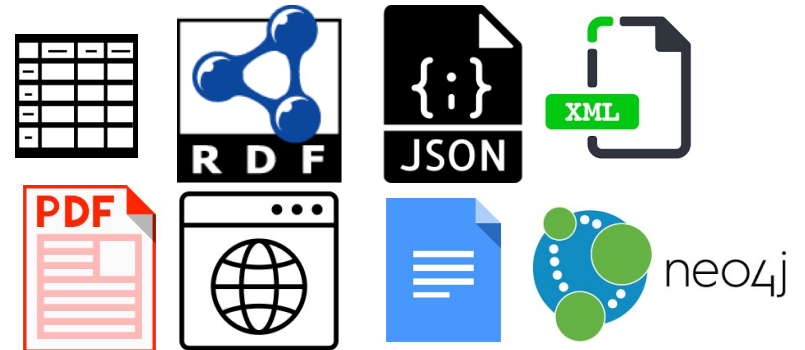
**Question**: When does a fake news post first cross into a supposedly legitimate community (e.g. members of the Parliament)?

- Looking for tweets connected to a fake news author, and to a community member; both connections are arbitrary paths (chains of author/likes/retweets/inParty/…)

# Graph-based integration of heterogeneous data sources

❑ The sources are not RDF. They can be (semi)structured, or unstructured (text).

❑ The sources may be very dynamic (projects started and abandoned as per news cycle and data availability).

❑ There is no schema. Data producers often uncollaborative.

❑ For most journalists, databases do not come naturally, and IT support is limited. They know keyword-based search...

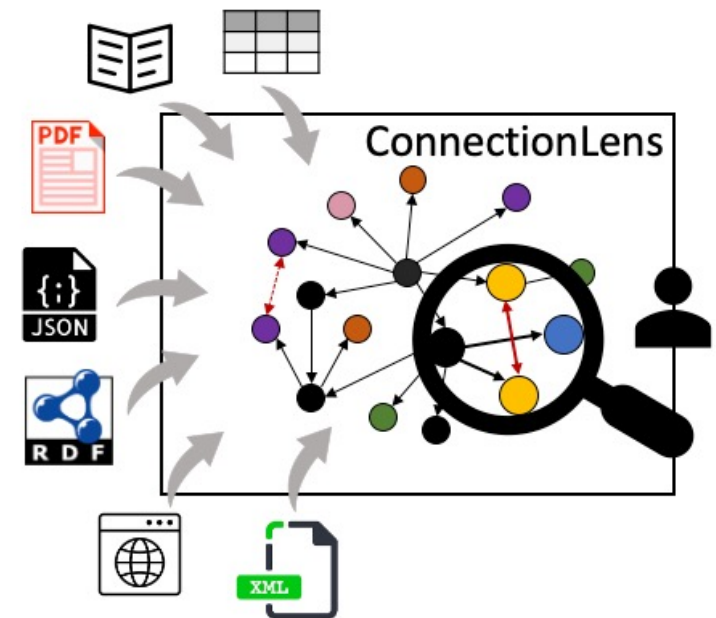➡ **Integrate heterogeneous sources within a graph, query w/ keywords**

# ConnectionLens: graph-based integration of heterogeneous data sources
https://team.inria.fr/cedar/connectionlens/

**Joint work with**:  J. Leblay (AIST Japan),
H. Galhardas and C. Conceiçao (U. Portugal),
A. Anadiotis, O. Balalau, N. Barret, T. Bouganim,
F. Chimienti, M.-Y. Haddad, T. Merabti,
P. Upadhyay (CEDAR) + interns

S.Horel (Le Monde, European Press Prize
          "Investigative Reporting Award 2018")

Ongoing work in ANR/DGA AI Chair SourcesSay
(https://sourcessay.inria.fr), DIM RFSI

# ConnectionLens principles [Chanial et al., 2018]

Integrate **any kind of data** into a **graph**

**Extract entities** from any text node (regardless of the model of the data source where the text comes from)

❑ Same entity in two different text nodes =
   link among the text nodes (*densification* of the graph)

The graph is **heterogeneous** and **irregular** →

Query it through **keywords**: find trees that connect 1 node matching each kwd

❑ Closely related to the Group Steiner Tree Problem (GSTP)

*Inría*

# ConnectionLens principles [Chanial et al., 2018]

Integrate **any kind of data** into a **graph**

**Extract entities** from any text node (regardless of the model of the data source where the text comes from)

❑ Same entity in two different text nodes =
link among the text nodes (*densification* of the graph)

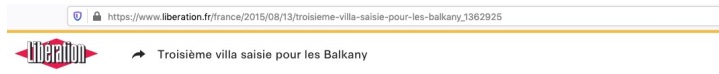The graph is **heterogeneous** and **irregular** →

Quer

❑ Cl

Rest of the talk based on state of the project 2+ years later:
https://arxiv.org/abs/2007.12488 [BDA 2020]
https://arxiv.org/abs/2009.04283 [BDA 2020]
https://arxiv.org/abs/2012.08830 [Invited to Elsevier Information Systems, under minor revision]
https://arxiv.org/abs/2102.04141

*Inria*

# ConnectionLens graph construction

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

*Inria*

# The Balkany and their African connections

# The Balkany and their African connections

# The Balkanys and their African connections

Public officials transparency high authority (CSV)

| Name | Owner | Location | Type |
|------|-------|----------|------|
| Dar Gyucy | P. Balkany | Marrakech | Real Estate |
| Moulin Cossy | I. Balkany | Giverny | Real Estate |

# The Balkanys and their African connections

Public officials transparency high authority (CSV)

| Name | Owner | Location | Type |
|------|-------|----------|------|
| Dar Gyucy | P. Balkany | Marrakech | Real Estate |
| Moulin Cossy | I. Balkany | Giverny | Real Estate |

National Directory of Elected Officials (JSON)

```
[{
    name: "Levallois-Perret",
    mayor: "P. Balkany",
    city-council: [
      {name: "I. Balkany"},
       …
    ]
}, …]
```

# The Balkanys and their African connections

Public officials transparency high authority (CSV)

| Name | Owner | Location | Type |
|------|-------|----------|------|
| Dar Gyucy | P. Balkany | Marrakech | Real Estate |
| Moulin Cossy | I. Balkany | Giverny | Real Estate |

dbpedia.org (RDF)

```
{
dbr:Marrakech
  dbr:name      "Marrakech"
  rdf:type      dbo:City ;
  dbo:country   dbr:Morrocco .
dbr:Morocco
  dbr:name      "Morocco"
  rdf:type      dbo:Country
  dbo:locatedIn dbr:Africa .
dbr:CentralAfricanRepublic
  dbr:name      "Central African Republic"
  dbo:locatedIn dbr:Africa .
}
```

National Directory of Elected Officials (JSON)

```
[{
    name: "Levallois-Perret",
    mayor: "P. Balkany",
    city-council: [
      {name: "I. Balkany"},
        …
    ]
}, …]
```

*Inria*

# The Balkanys and their African connections

Public officials transparency high authority (CSV)

| Name | Owner | Location | Type |
|------|-------|----------|------|
| Dar Gyucy | P. Balkany | Marrakech | Real Estate |
| Moulin Cossy | I. Balkany | Giverny | Real Estate |

dbpedia.org (RDF)

```
{
dbr:Marrakech
  dbr:name       "Marrakech"
  rdf:type       dbo:City ;
  dbo:country    dbr:Morrocco .
dbr:Morocco
  dbr:name       "Morocco"
  rdf:type       dbo:Country
  dbo:locatedIn  dbr:Africa .
dbr:CentralAfricanRepublic
  dbr:name       "Central African Republic"
  dbo:locatedIn  dbr:Africa .
}
```

National Directory of Elected Officials (JSON)

```
[{
    name: "Levallois-Perret",
    mayor: "P. Balkany",
    city-council: [
      {name: "I. Balkany"},
       …
    ]
}, …]
```

Libération – Nov. 13, 2014 (Text)

**Balkany mineur de fonds**

L'élu de Levallois-Perret est soupçonné d'avoir touché 5 millions de dollars de commission en 2009 grâce à son rôle d'intermédiaire entre Areva et la Centrafrique dans le dossier Uramin. […]

*Inria*

# How is Levallois-Perret connected to Africa and "real estate"?

Public officials transparency high authority (CSV)

| Name | Owner | Location | Type |
|------|-------|----------|------|
| Dar Gyucy | P. Balkany | Marrakech | Real Estate |
| Moulin Cossy | I. Balkany | Giverny | Real Estate |

dbpedia.org (RDF)

```
{
dbr:Marrakech
  dbr:name     "Marrakech"
  rdf:type     dbo:City ;
  dbo:country  dbr:Morrocco .
dbr:Morocco
  dbr:name     "Morocco"
  rdf:type     dbo:Country
  dbo:locatedIn dbr:Africa .
dbr:CentralAfricanRepublic
  dbr:name     "Central African Republic"
  dbo:locatedIn dbr:Africa .
}
```

National Directory of Elected Officials (JSON)

```
[{
    name: "Levallois-Perret",
    mayor: "P. Balkany",
    city-council: [
      {name: "I. Balkany"},
       ...
    ]
}, ...]
```

Libération – Nov. 13, 2014 (Text)

**Balkany mineur de fonds**

L'élu de Levallois-Perret est soupçonné d'avoir touché 5 millions de dollars de commission en 2009 grâce à son rôle d'intermédiaire entre Areva et la Centrafrique dans le dossier Uramin. […]

Inria

# Idea: integrate all data sources into a heterogeneous graph

Ioana Manolescu, Inria and Institut Polytechnique de Paris      Trustworthy Data Science and AI, SFU   April 2021

*Ínría*

# Graph construction stages

1. **Primary node and edge construction**

   ❑   Direct for XML, JSON, RDF, HTML

   ❑   1 relational tuple=1 node;
       primary keys-foreign keys as links

   ❑   Convert information from PDF into:

   ❑   JSON for text content

   ❑   RDF describing tables

# Graph construction stages

1. **Primary node and edge construction**

   ❑ Direct for XML, JSON, RDF, HTML

   ❑ 1 relational tuple=1 node; PK-FKs as links

   ❑ [Optional] segment text documents

   ❑ Extract information from PDF into: (a) JSON, and (b) RDF describing tables

2. **Entity extraction**

   ❑ From all text nodes of all the sources: **entity node** child of text node

   ❑ [VLDB2018]: based on Stanford NER

   ❑ [BDA2020] Developed and trained new entity extractor from French, based on Flair framework

# Graph construction stages

**2.   Entity extraction**

❑   From all text nodes of all the sources: **entity node** child of text node

❑   [VLDB2018]: based on Stanford NER

❑   [BDA2020] Developed and trained new entity extractor from French, based on Flair framework

**3. Entity disambiguation**

❑   For each recognized entity, e.g., "Hollande" the place or the person?

❑   Built novel disambiguation pipeline for French, based on Ambiverse framework

   ❑   Based on knowledge bases (WikiData, YAGO) and Wikipedia

   ❑   Helpful on well-known entities

Ioana Manolescu, Inria and Institut Polytechnique de Paris        Trustworthy Data Science and AI, SFU   April 2021

# Graph construction stages

**4. Node matching**

- ❑ To create sameAs edges:
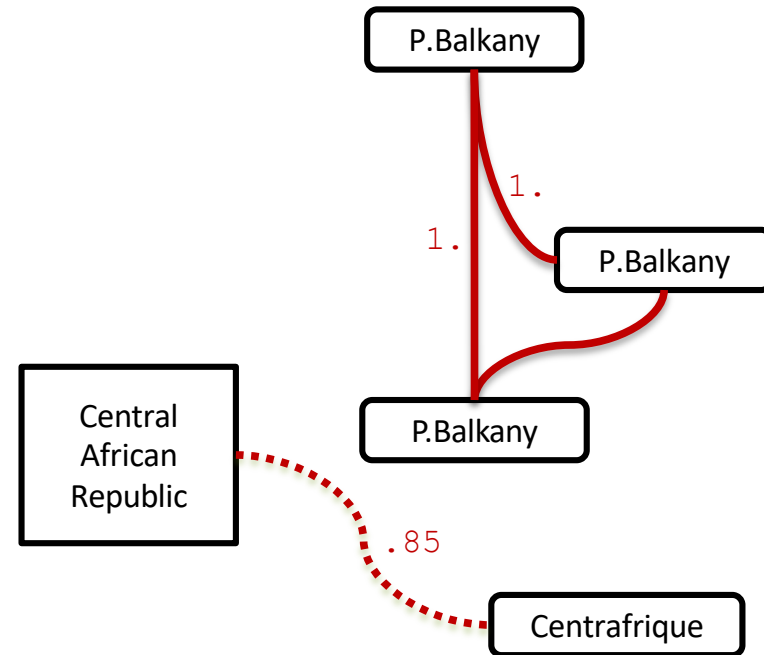  - ❑ Strong sameAs edges: equivalent nodes 1.
  - ❑ Weak sameAs edges: similar nodes .85
- ❑ Appropriate distance functions
- ❑ *New*: more normalization → better matching
- ❑ Remains quadratic at the core ☹ , so...

**Node factorization (heuristic)**: create only one node per label per document (or per graph)



Ioana Manolescu, Inria and Institut Polytechnique de Paris        Trustworthy Data Science and AI, SFU   April 2021

*Inria*

# ConnectionLens graph querying

Ioana Manolescu, Inria and Institut Polytechnique de Paris        Trustworthy Data Science and AI, SFU   April 2021

*Inria*

# Querying problem statement

❑ Given the graph G = (N, E) built out of the datasets D and a query Q= {w1, ...,wm}, return the **k highest-score minimal answer trees**.

❑ An answer tree is a set of edges which (*i*) form a tree (*ii*) contain at least one node whose label matches each keyword wi.

❑ We are interested in **minimal answer trees**, that is:
   ❑ Removing an edge from the tree should make it lack some keyword(s).
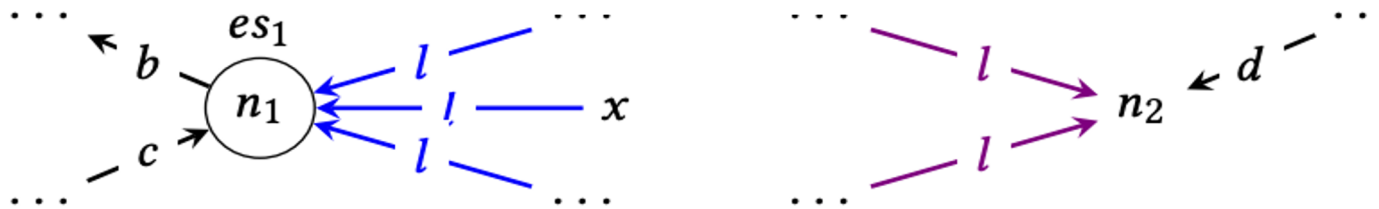   ❑ If a keyword matches more than one nodes in the answer tree, then all these matching nodes must be equivalent.

# Search space and complexity

❑ Problem related to the **(Group) Steiner Tree Problem**
   ❑ Given graph G, and nodes n1, . . . , nm, the Steiner Tree Problem (STP) requires the smallest tree in G that connects all the nodes. Known NP-hard problem in |G|
   ❑ Group STP: start with m groups of nodes
❑ **Differences with our problem:**
   ❑ Each edge can be taken in both directions: exponential increase in search space size
   ❑ We need the k smallest-cost trees, not just one.
   ❑ Score function may be non-monotonous; no optimal sub-structure property
❑ **Large literature on kwd search in text, resp. structured data.**

   ❑ Differ in search space and/or make limitative assumptions on score
❑ **Our approach: enumerate solutions until time-out or max number of solutions reached.**
   ❑ Return best k solutions found

# GAM (Grow and Aggressive Merge) Algorithm

- Builds trees "backward" from the keyword matches
- **GROW** adds an edge to the root of a tree, **MERGE** merges trees with the same root
- Exploration (GROW) order:
  1. Favor trees matching the largest number of query keywords
  2. To break ties, favor smaller trees
  3. To break second tie between (t1, e1), (t2, e2), we prefer the pair with the higher specificity edge.

The specificity of $\quad e = n_1 \xrightarrow{l} n_2 \quad$ is: $\quad s(e) = 2/(N^l_{n_1 \rightarrow} + N^l_{\rightarrow n_2})$



Special measures to handle equivalence clusters efficiently

# Sample query answers

# ConnectionLens architecture and performance

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

# Implementation

❑ Java (220 classes/40K LOC), Python (25 classes/2700 LOC), JS + CSS

❑ Available online: https://gitlab.inria.fr/cedar/connectionlens



Relational DB

ConnectionLens graph construction

Nodes+edges

Parallel, in-mem keyword search

In-memory graph

Keyword search algorithm

text

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

Inria

# Implementation https://arxiv.org/abs/2012.08830

❑ Graph creation time mostly **linear in the size of the data**

❑ Costliest operations involve ML (disambiguation, extraction)

   ❑ **Batch extraction:** 20x speed-up on GPU, 2x speed-up on regular server

   ❑ **Extraction policies** replace or avoid extraction in some parts of the data

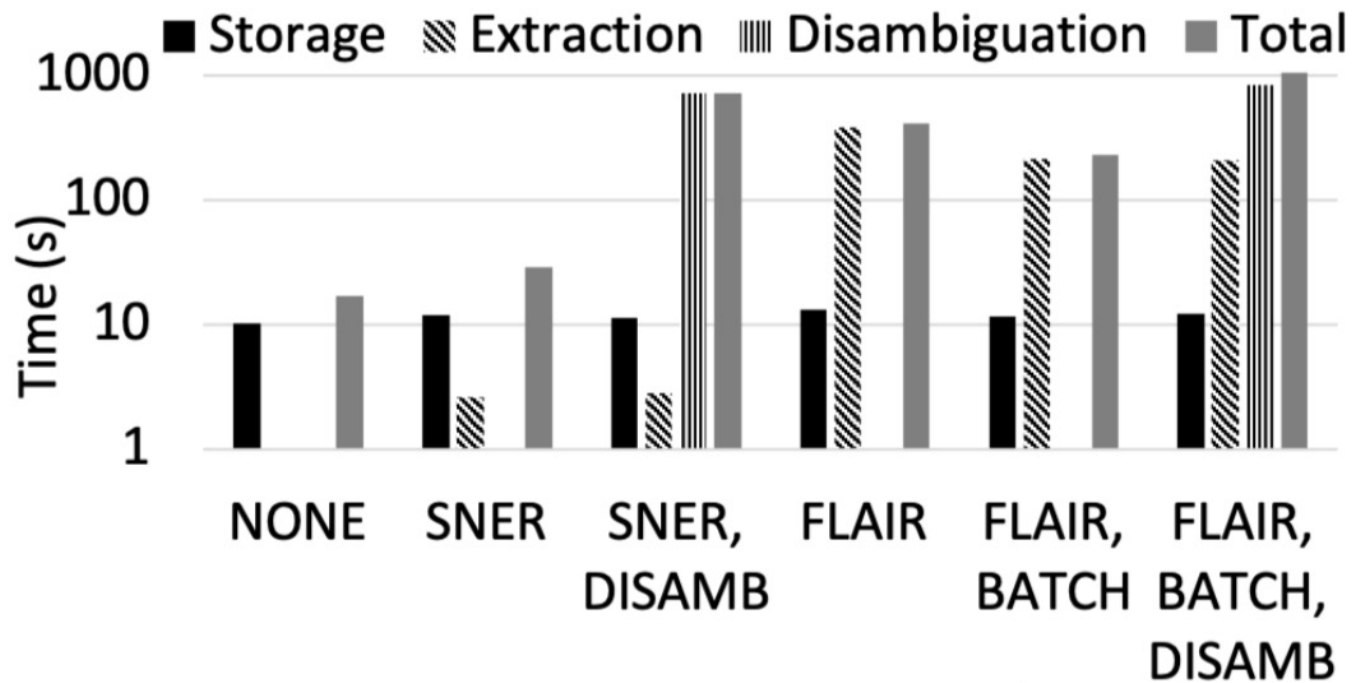# Graph creation performance: storage, extraction, disambiguation https://arxiv.org/abs/2012.08830



Figure 6: Graph construction time (seconds).

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

# Graph creation performance: batch extraction
https://arxiv.org/abs/2012.08830



Figure 7: YAGO loading time (minutes) using Flair.

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

# Application: conflicts of interest in the biomedical domain https://arxiv.org/abs/2102.04141

Collaboration with Stéphane Horel (Le Monde)

Data: XML, PDF→JSON, HTML

| $|N|$ | $|E|$ | $|N|$ | $|N_P|$ | $|N_O|$ | $|N_L|$ |
|---|---|---|---|---|---|
| XML | 32,028,429 | 19,851,904 | 1,483,631 | 584,734 | 126,629 |
| JSON | 1,025,307 | 432,303 | 75,297 | 7,320 | 4,139 |
| HTML | 246,636 | 185,479 | 3,726 | 7,227 | 320 |
| Total | 33,300,372 | 20,469,686 | 1,562,654 | 665,167 | 131,088 |

**Table 3: Statistics on Conflict of Interest application graph.**

*Inria*

# Application: conflicts of interest in the biomedical domain https://arxiv.org/abs/2102.04141

Collaboration with Stéphane Horel (Le Monde)

Data: XML, PDF→JSON, HTML

| $|N|$ | $|E|$ | $|N|$ | $|N_P|$ | $|N_O|$ | $|N_L|$ |
|---|---|---|---|---|---|
| XML | 32,028,429 | 19,851,904 | 1,483,631 | 584,734 | 126,629 |
| JSON | 1,025,307 | 432,303 | 75,297 | 7,320 | 4,139 |
| HTML | 246,636 | 185,479 | 3,726 | 7,227 | 320 |
| Total | 33,300,372 | 20,469,686 | 1,562,654 | 665,167 | 131,088 |

**Table 3: Statistics on Conflict of Interest application graph.**

# Application: conflicts of interest in the biomedical domain https://arxiv.org/abs/2102.04141

## Collaboration with  Stéphane Horel (Le Monde)



Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

# Graph quality experiments
## https://arxiv.org/abs/2012.08830

❑ PDF extraction accuracy: 63%

❑ F1  score for entity extraction from French:

❑ Flair stacked forward and backward  embeddings with French fastText embeddings: 73%

❑ Spacy: 63%

❑ StanfordNER: 45%

❑ F1 score of disambiguation: 86%

*Inria*

# ConnectionLens in the scientific landscape

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021

*Inria*

# ConnectionLens in the scientific landscape (1)

**Data integration** for structured, semistructured and unstructured data

❑ "Ad-hoc" (combinations of sources to be unioned, joined, or chained)

❑ No schema, ontologies, queries known in advance

❑ Mediator previously developed [Bonaque et al., VLDB 2016] inappropriate due to complexity, lack of structure, and performance → graph warehouse

❑ Lack of structure forces reachability queries instead of join

❑ Price to pay for powerful integration

**Data cleaning** aspects: Similarity links require value or entity matching

❑ Avoid constructing structured objects ("clean tuples"): don't seem necessary

*Inria*

# ConnectionLens in the scientific landscape (2)

**Graph construction**

❑ Users of **entity extraction** modules, trained a model for French

**Keyword search on structured data**

❑ Previously studied for relational, graph, or XML databases

   ❑ Typically assume structure/regularity in the graph

   ❑ Exploit favorable properties of the score function

❑ First keyword search algorithm across heterogeneous data sources,
w/o assumptions on  score, w/o sub-optimal structure prop., w/ bidirectional search

❑ In-memory graph store and parallel query processor (200x speed-up)

*Inria*

# Ongoing work

❑ Extending and improving the in-memory query processor (A. Anadiotis, F. Chimienti, IM)

❑ Relationship extraction  (O. Balalau, M. Mohanty, IM)

❑ Natural language querying of the graph (O. Balalau, P. Upadhyay, IM + PhD in fall 2021)

❑ Improving the quality of graph linking (T. Bouganim, H. Galhardas, IM)

❑ Abstracting CL graphs (N. Barret, H. Galhardas, P. Upadhyay, IM)

❑ Applications:

   ❑ Conflicts of Interests in the biomedical domain (w/ S. Horel and G. Fooks, Aston U., UK)
   ❑ Mediacités (w/ WeDoData)

*Inria*

# 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!

## Thank you

**Questions?**

**ConnectionLens**: https://team.inria.fr/cedar/connectionlens/

Ioana Manolescu, Inria and Institut Polytechnique de Paris          Trustworthy Data Science and AI, SFU   April 2021