

Inria

All your digital data sources in one place:
the **SourcesSay** ANR/DGA project



<https://sourcessay.inria.fr>

Ioana Manolescu

CEDAR team

Inria Saclay-Île-de-France, Institut Polytechnique de Paris

[@ioanamanol](#) [@cedarinrialix](#)





Why should computer scientists focus on tools for journalists?

Bad memories: Romania, 1989



Bad memories: Romania, 1989



Ceaușescu re-elected
at the 14th congress!

Bad memories: Romania, 1989



Ceaușescu re-elected
at the 14th congress!

He had held power
since 1965.

Bad memories: Romania, 1985



1990: things got better



... kind of



1000 dead (approx.)
No one convicted.

Democratic societies crucially need the press

- ❑ To debate and express dissent
- ❑ To analyze, confirm or refute public statements
- ❑ To expose and explain society functioning



Socialist Romania, 1984

Fact-checking

(Data) journalism



Democratic societies crucially need the press

❑ To debate and express dissent



Socialist Romania, 1984

❑ To analyze, confirm or refute public statements

Fact-checking

❑ To expose and explain society functioning

Open-source
intelligence

bellingscat



Information, **mis**information, **dis**information and **data**

How do we analyze statements?

A statement is made



How do we analyze statements?

A statement is made



What should we think of this?



How do we analyze statements?

A statement is made



What should we think of this?



Information?
Misinformation?
Dysinformation?

How do we analyze statements?

A statement is made



What should we think of this?



Information?
Misinformation?
Dysinformation?



Evaluate
source authority
(past statements...)

How do we analyze statements?

A statement is made

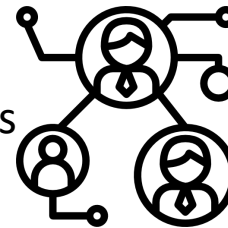


What should we think of this?



Information?
Misinformation?
Dysinformation?

Evaluate
source
connections



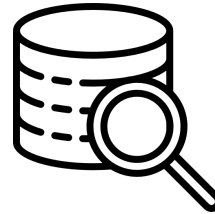
Evaluate
source authority
(past statements...)

How do we analyze statements?

A statement is made



Consult
reference
data

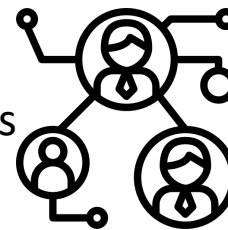


What should we think of this?



Information?
Misinformation?
Dysinformation?

Evaluate
source
connections



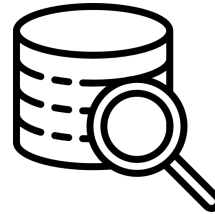
Evaluate
source authority
(past statements...)

How do we analyze statements?

A statement is made



Consult
reference
data

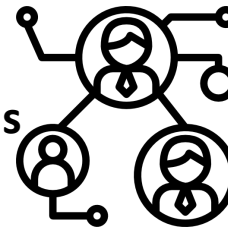


What should we think of this?



Information?
Misinformation?
Dysinformation?

Evaluate
source
connections



Evaluate
source authority
(past statements...)

How do we analyze statements?

A statement is made



Information?

Misinformation?

Dysinformation?



source
connections



source past
statements



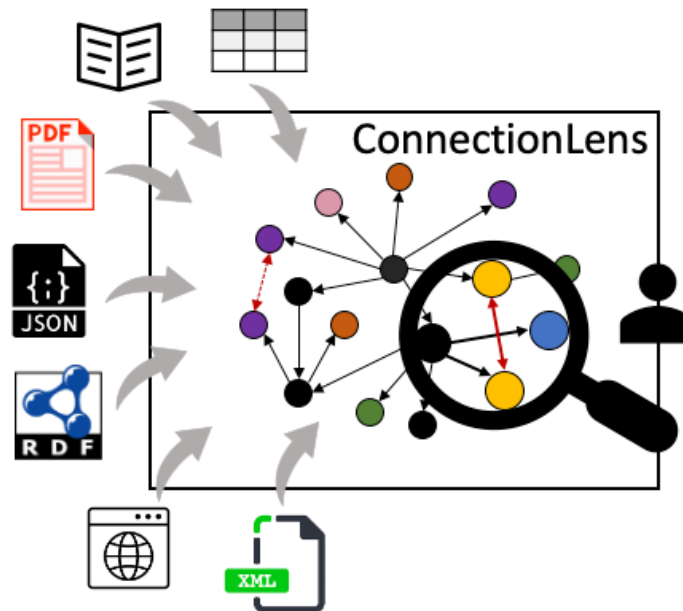
fact-checks



reference
data

ConnectionLens: graph-based integration of heterogeneous data sources

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

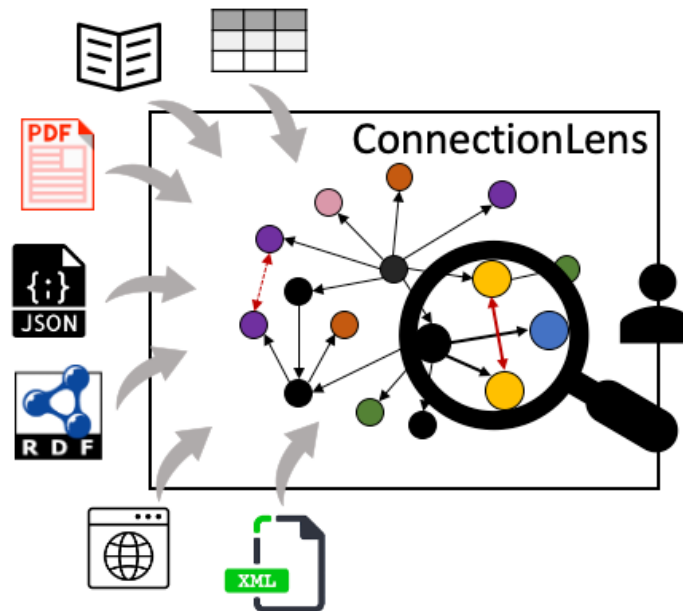


ConnectionLens: graph-based integration of heterogeneous data sources

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

Data ingestion:

- Web pages
- Documents
- Data:
 - Social media content
 - (Semi)structured data: CSV, RDF, XML, JSON...

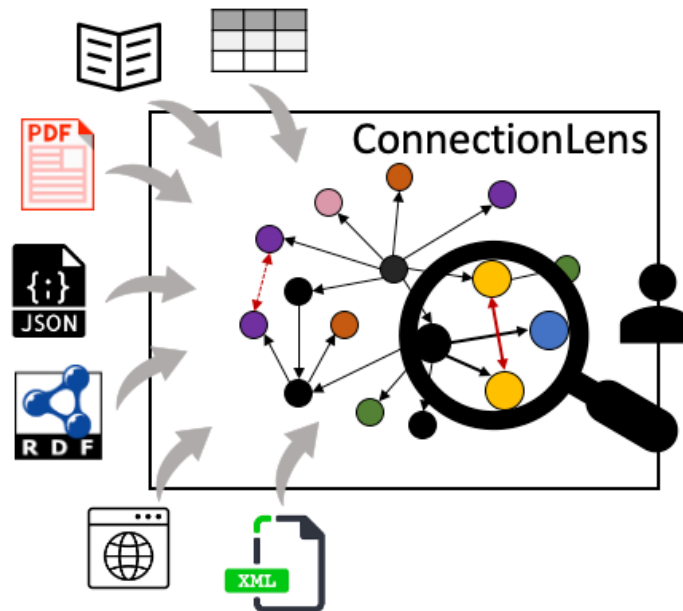


ConnectionLens: graph-based integration of heterogeneous data sources

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

Data ingestion:

- Web pages
- Documents
- Data:
 - Social media content
 - (Semi)structured data: CSV, RDF, XML, JSON...



Data enrichment:

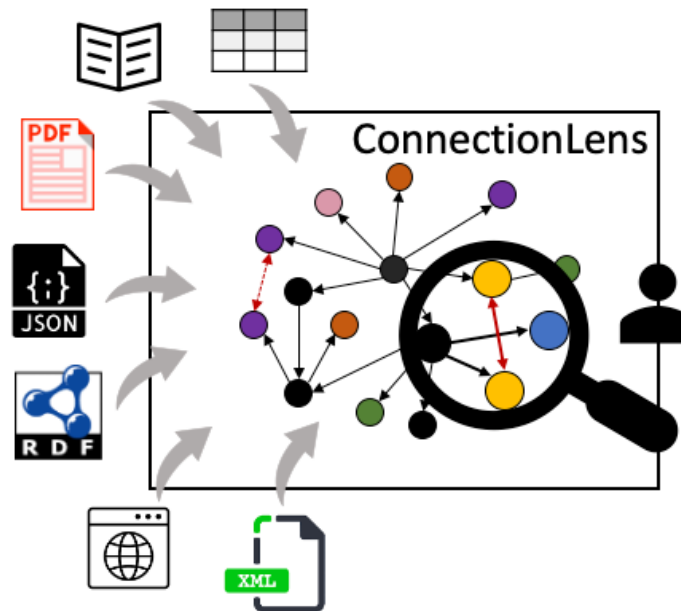
- Extraction of **entities** from text (people, places, dates, organizations, emails, URIs, ...)
- Extraction of **relationships** from text

ConnectionLens: graph-based integration of heterogeneous data sources

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

Data ingestion:

- Web pages
- Documents
- Data:
 - Social media content
 - (Semi)structured data: CSV, RDF, XML, JSON...



Data enrichment:

- Extraction of **entities** from text (people, places, dates, organizations, emails, URIs, ...)
- Extraction of **relationships** from text
- Disambiguation of entity names based on context, e.g., « Hollande »

ConnectionLens: graph-based integration of heterogeneous data sources

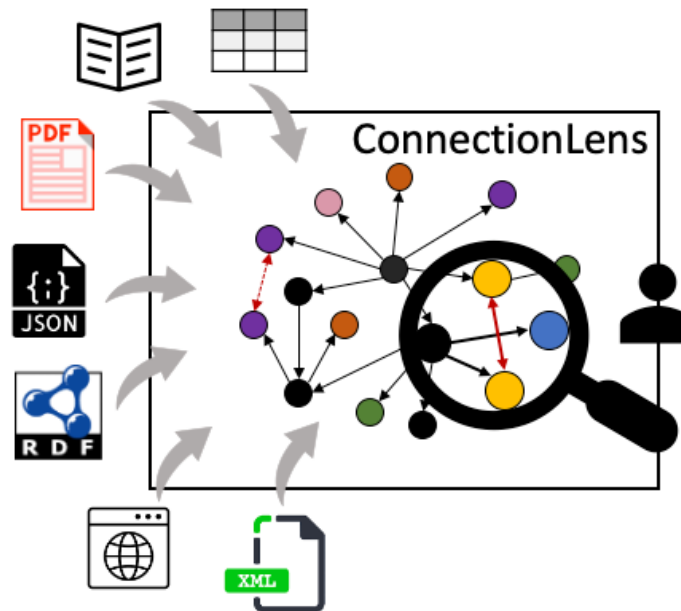
<https://team.inria.fr/cedar/connectionlens/>

Data ingestion:

- Web pages
- Documents
- Data

Data enrichment:

- Entity extraction
- Relationship extraction
- Disambiguation



Data exploitation (1):

- I don't know much about « data formats ». What does this dataset contain?

ConnectionLens: graph-based integration of heterogeneous data sources

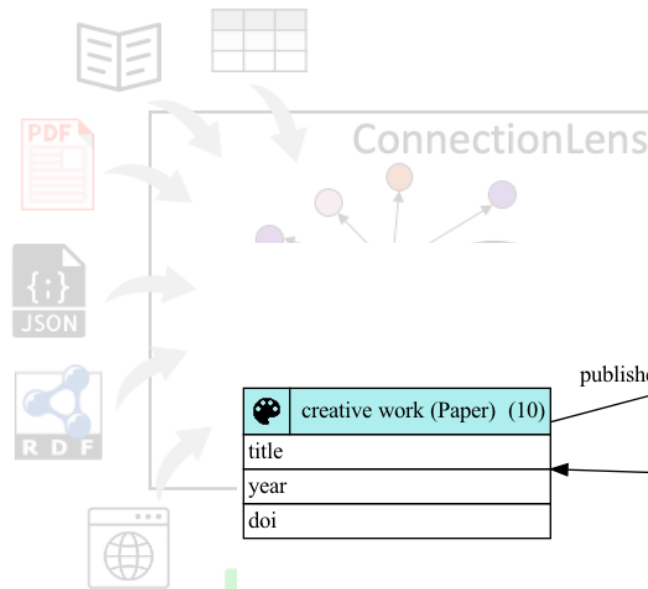
<https://team.inria.fr/cedar/connectionlens/>

Data ingestion:

- Web pages
- Documents
- Data

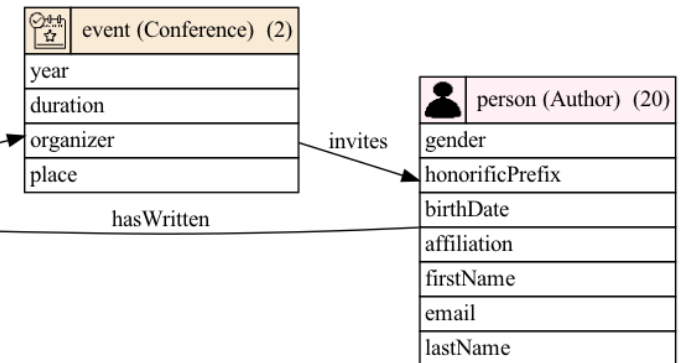
Data enrichment:

- Entity extraction
- Relationship extraction
- Disambiguation



Data exploitation (1):

- I don't know much about « data formats ». What does this dataset contain?



ConnectionLens: graph-based integration of heterogeneous data sources

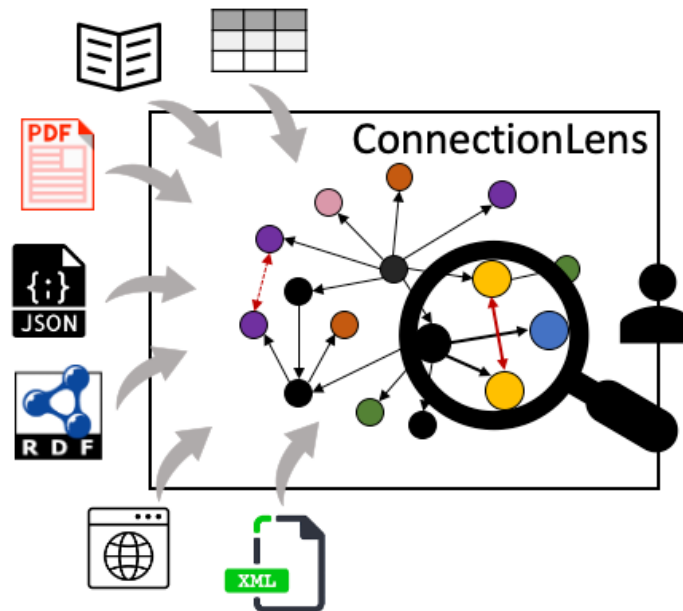
<https://team.inria.fr/cedar/connectionlens/>

Data ingestion:

- Web pages
- Documents
- Data

Data enrichment:

- Entity extraction
- Relationship extraction
- Disambiguation



Data exploitation (1):

- I don't know much about « data formats ». What does this dataset contain?
- Entity search, statistics, co-occurrence
 - « Which organizations (dates, people, ...) are found in connection with A.Kohler? »
- Connections between entities
 - « How is A.Kohler connected to the STX and MSC companies? »

ConnectionLens: graph-based integration of heterogeneous data sources

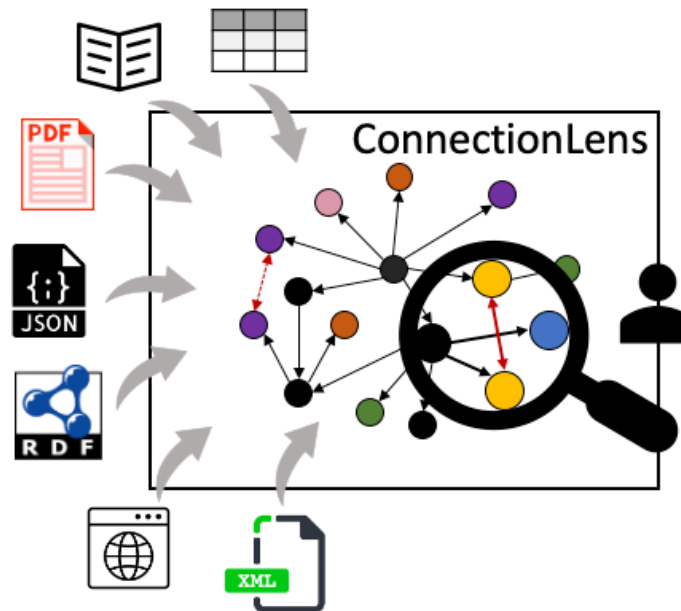
<https://team.inria.fr/cedar/connectionlens/>

Data ingestion:

- Web pages
- Documents
- Data

Data enrichment:

- Entity extraction
- Relationship extraction
- Disambiguation



Data exploitation (2):

- What are the connections between people and organizations in this dataset?
 - Which people (organizations)?
 - Which connections? (« worked for », « funded by »...)

Sort queries by length | Sort queries by number of associated data paths | Hide/show queries without associated data paths

ORGANIZATION — label — LaunchSite — country — #val — agency — Spacecraft — description — PERSON (1362 data paths)

ID	ORGANIZATION	label	LaunchSite	country	#val	agency	Spacecraft	description	PERSON
0	Edwards Air Force Base, United States	http://www.w3.org/2000/01/rdf-schema#label	http://data.kasabi.com/dataset/nasa/launchsite/edwardsairforcebase	http://purl.org/net/schemas/space/country	United States	http://purl.org/net/schemas/space/agency	http://data.kasabi.com/dataset/nasa/spacecraft/1968-089A	http://purl.org/dc/elements/1.1/description	Apollo 7 was the first crewed flight of the Apollo spacecraft, with astronauts Walter Schirra, Jr, Donn Eisele, and Walter Cunningham on board. The primary objectives of the Earth orbiting mission were to demonstrate Command and Service Module (CSM), crew, launch vehicle, and mission support facilities performance and to demonstrate CSM rendezvous capability. Two photographic experiments and three medical experiments were planned.
1	Edwards Air Force Base, United States	http://www.w3.org/2000/01/rdf-schema#label	http://data.kasabi.com/dataset/nasa/launchsite/edwardsairforcebase	http://purl.org/net/schemas/space/country	United States	http://purl.org/net/schemas/space/agency	http://data.kasabi.com/dataset/nasa/spacecraft/1968-089A	http://purl.org/dc/elements/1.1/description	Apollo 7 was the first crewed flight of the Apollo spacecraft, with astronauts Walter Schirra, Jr , Donn Eisele, and Walter Cunningham on board. The primary objectives of the Earth orbiting mission were to demonstrate Command and Service Module (CSM), crew, launch vehicle, and mission support facilities performance and to demonstrate CSM rendezvous capability. Two photographic experiments and three medical experiments were planned.
3	Edwards Air Force Base, United States	http://www.w3.org/2000/01/rdf-schema#label	http://data.kasabi.com/dataset/nasa/launchsite/edwardsairforcebase	http://purl.org/net/schemas/space/country	United States	http://purl.org/net/schemas/space/agency	http://data.kasabi.com/dataset/nasa/spacecraft/1968-089A	http://purl.org/dc/elements/1.1/description	Apollo 7 was the first crewed flight of the Apollo spacecraft, with astronauts Walter Schirra, Jr, Donn Eisele , and Walter Cunningham on board. The primary objectives of the Earth orbiting mission were to demonstrate Command and Service Module (CSM), crew, launch vehicle, and mission support facilities performance and to demonstrate CSM rendezvous capability. Two photographic experiments and three medical experiments were planned.
5	Edwards Air Force Base, United States	http://www.w3.org/2000/01/rdf-schema#label	http://data.kasabi.com/dataset/nasa/launchsite/edwardsairforcebase	http://purl.org/net/schemas/space/country	United States	http://purl.org/net/schemas/space/agency	http://data.kasabi.com/dataset/nasa/spacecraft/1990-019A	http://purl.org/dc/elements/1.1/description	STS-36 was the sixth flight of the shuttle Atlantis. On board were John Creighton , John Casper, David Hilmers, Richard Mullane and Pierre Thout. The primary purpose for this mission was to launch a spacecraft for the US Department of Defense. Mission duration was 106 hours 18 minutes 22 seconds.
6	Edwards Air Force Base, United States	http://www.w3.org/2000/01/rdf-schema#label	http://data.kasabi.com/dataset/nasa/launchsite/edwardsairforcebase	http://purl.org/net/schemas/space/country	United States	http://purl.org/net/schemas/space/agency	http://data.kasabi.com/dataset/nasa/spacecraft/1990-019A	http://purl.org/dc/elements/1.1/description	STS-36 was the sixth flight of the shuttle Atlantis. On board were John Creighton, John Casper, David Hilmers , Richard Mullane and Pierre Thout. The primary purpose for this mission was to launch a spacecraft for the US Department of Defense. Mission duration was 106 hours 18 minutes 22 seconds.
7	Edwards Air Force Base, United States	http://www.w3.org/2000/01/rdf-schema#label	http://data.kasabi.com/dataset/nasa/launchsite/edwardsairforcebase	http://purl.org/net/schemas/space/country	United States	http://purl.org/net/schemas/space/agency	http://data.kasabi.com/dataset/nasa/spacecraft/1990-019A	http://purl.org/dc/elements/1.1/description	STS-36 was the sixth flight of the shuttle Atlantis. On board were John Creighton, John Casper, David Hilmers, Richard Mullane and Pierre Thout. The primary purpose for this mission was to launch a spacecraft for the US Department of Defense. Mission duration was 106 hours 18 minutes 22 seconds.

ConnectionLens: graph-based integration of heterogeneous data sources

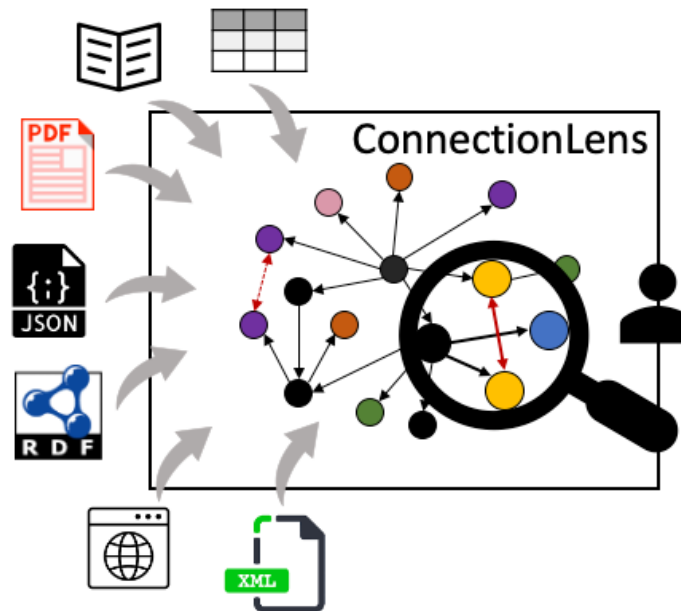
<https://team.inria.fr/cedar/connectionlens/>

Data ingestion:

- Web pages
- Documents
- Data

Data enrichment:

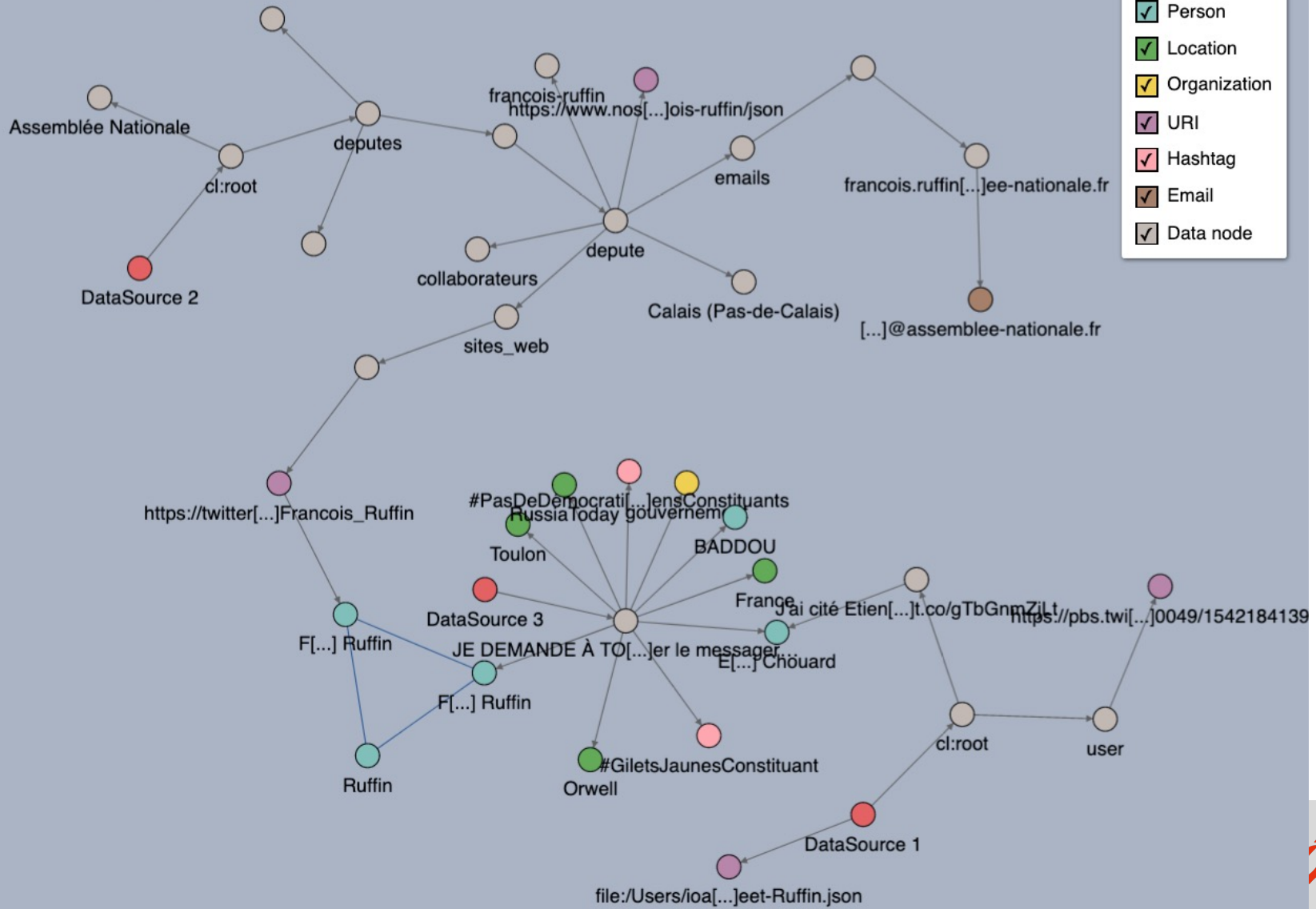
- Entity extraction
- Relationship extraction
- Disambiguation



Data exploitation (2):

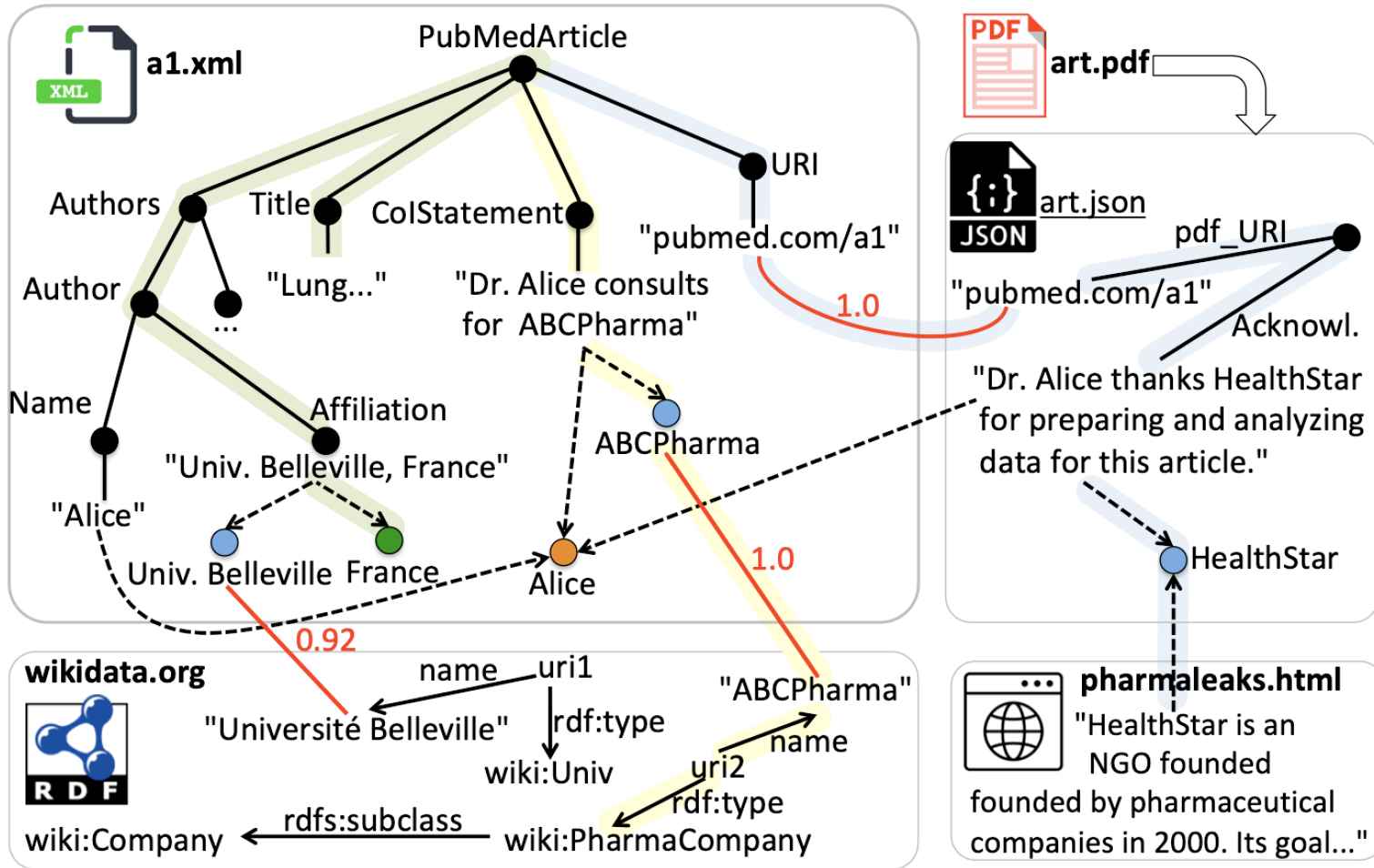
- What are the connections between people and organizations in this dataset?
- Interactive graph exploration
 - « What should I know about A.Kohler? »
- Structured and semistructured querying (requires database skills)

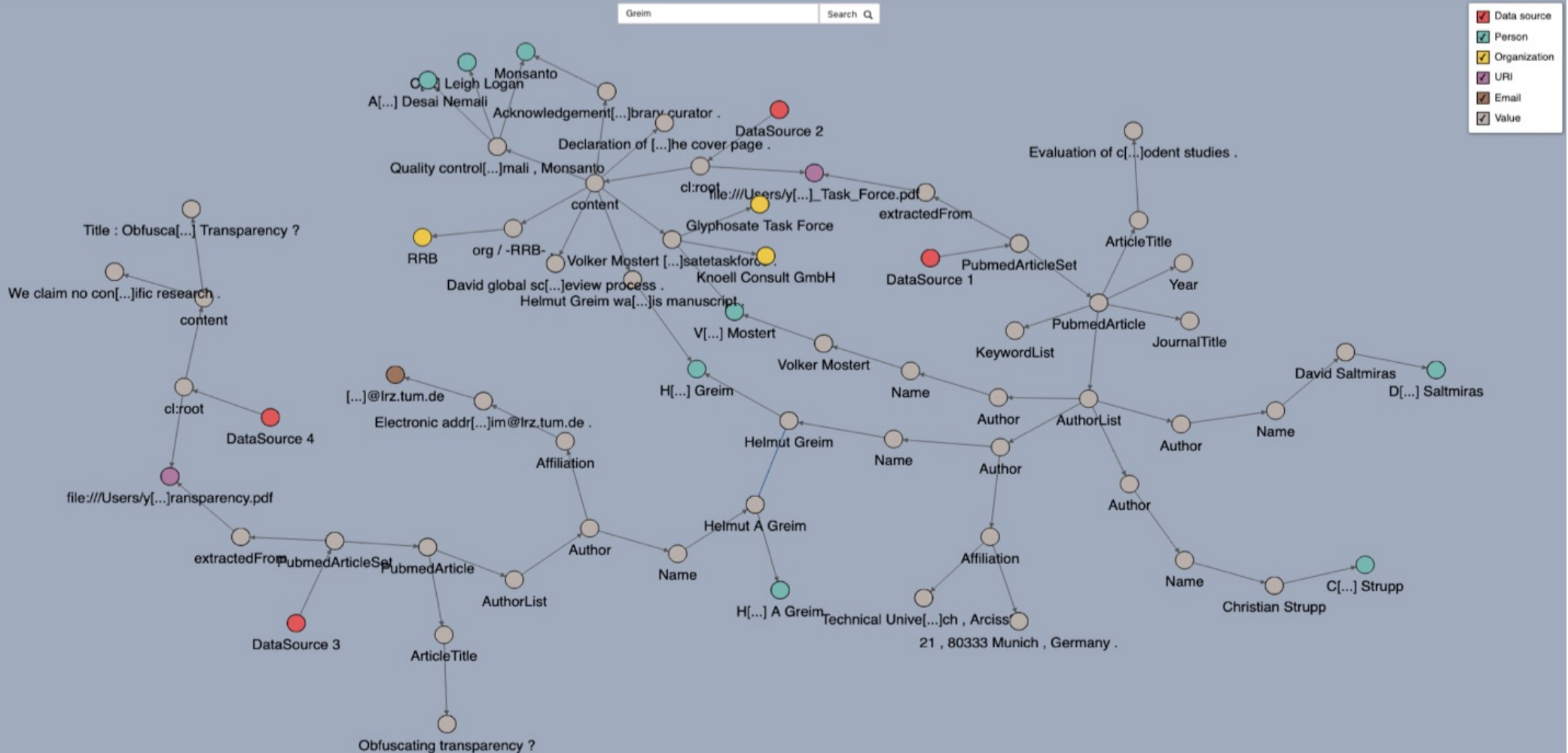
Search Q

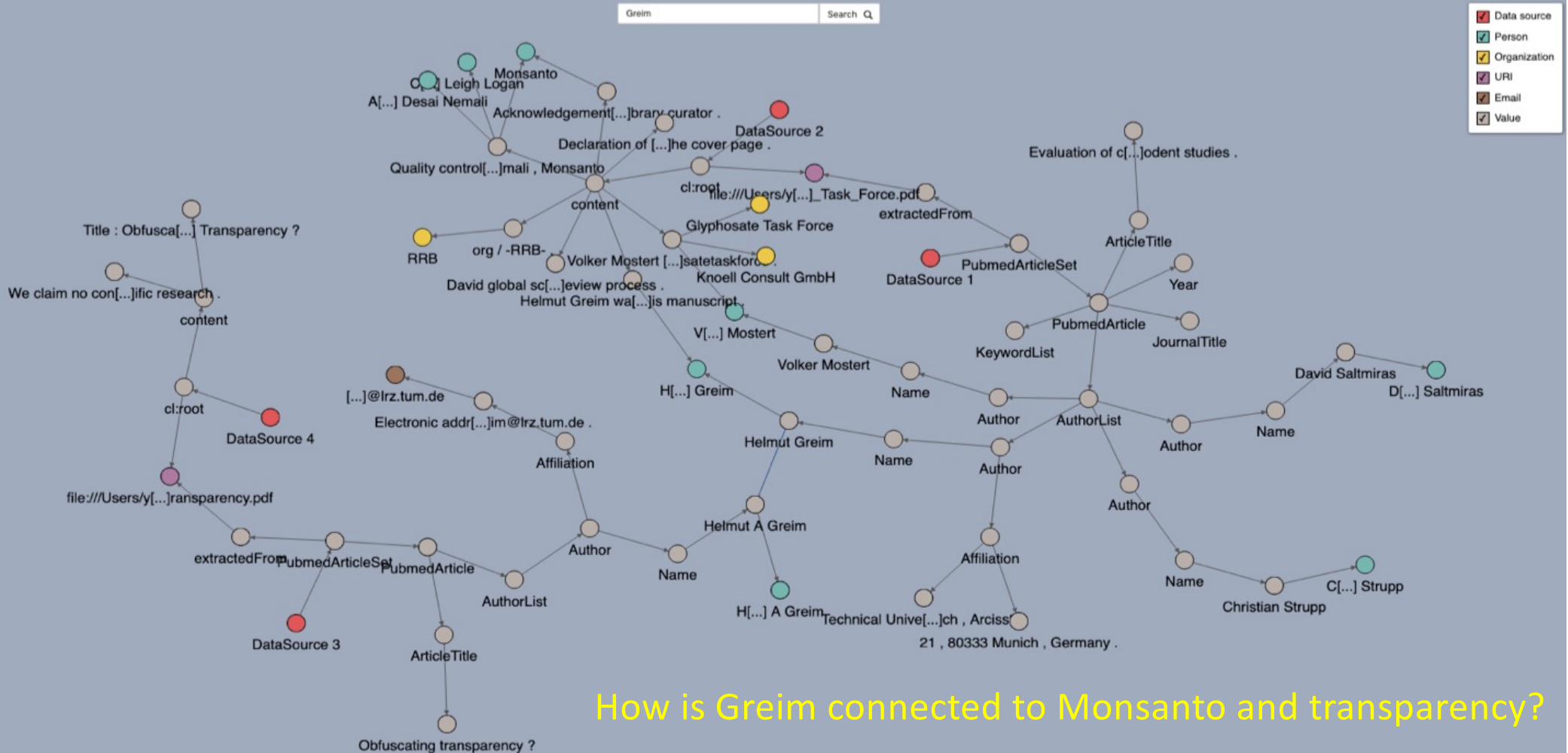




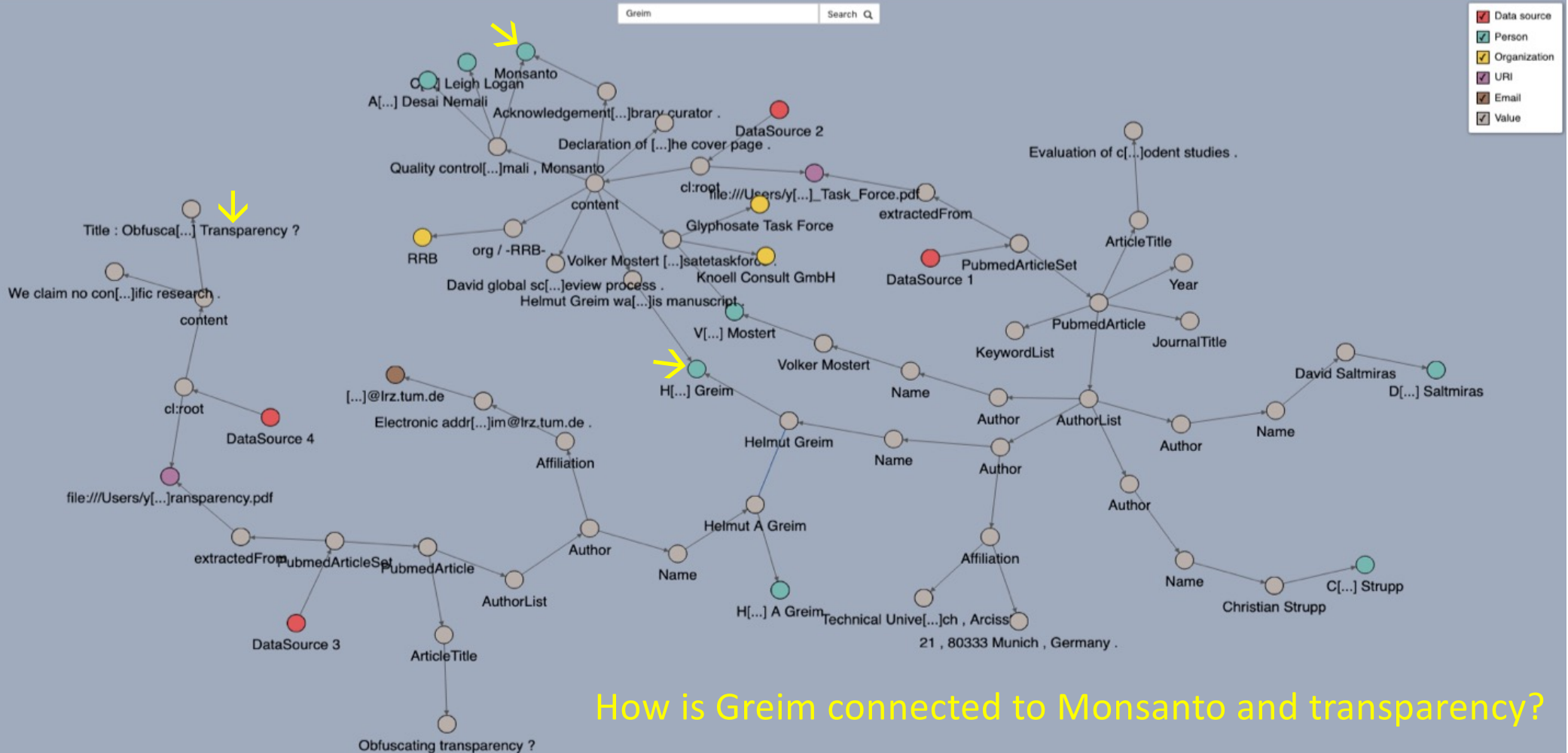
Application: conflicts of interest in the biomedical domain







How is Greim connected to Monsanto and transparency?



How is Greim connected to Monsanto and transparency?

Application: conflicts of interest in the biomedical domain [IEEE DataEngBull 2021, CIKM2021]

Collaboration with Stéphane Horel (Le Monde)

The screenshot shows the CL-LinkingCOIS web application. The browser address bar displays the URL: <https://cl-linkingcois.saclay.inria.fr/?keyword=org%3ANovartis>. The page header includes navigation links: Home, Dashboard, Settings, Help, About, and a user profile dropdown for ioana.manolescu@inria.fr. The main heading is "CL-LinkingCOIS". Below it is a search bar containing "org:Novartis" and a "Search" button. A green badge indicates "315 results".

CoiStatement	PubmedLink
Competing interests : Richard L. Baretto reports grants , personal fees and honoraria for lectures from ThermoFisher , Novartis and ALK Abello outside the submitted work. Mamidipudi Thirumala Krishna received honoraria for lectures from Thermo Fisher and ALK Abello , outside the submitted work.	view the pubmed paper
Conflict of interest : Benjamin Waschk has nothing to disclose. Conflict of interest : Christian Herr has nothing to disclose. Conflict of interest : Christina Magnusser has nothing to disclose. Conflict of interest : Christoph Sinning has nothing to disclose. Conflict of interest : Claus F. Vogelmeier reports grants and personal fees from AstraZeneca , Boehringer Ingelheim , GlaxoSmithKline , Grifols and Novartis , personal fees from CSL Behring , Chiesi , Menarini , Mundipharma , Teva and Cipla , grants from Bayer-Schering , MSD and Pfizer , outside the submitted work. Conflict of interest : Henrik Watz reports personal fees from AstraZeneca , Boehringer Ingelheim , GlaxoSmithKline , BerlinChemie , Chiesi , Novartis and Roche , outside the submitted work. Conflict of interest : Johannes T. Neumann reports personal fees from Abbott Diagnostics and Siemens , outside the submitted work. Conflict of	view the pubmed paper

Application: conflicts of interest in the biomedical domain [IEEE DataEngBull 2021, CIKM2021]

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.

Graph creation performance: storage, extraction, disambiguation [InfSys2022]

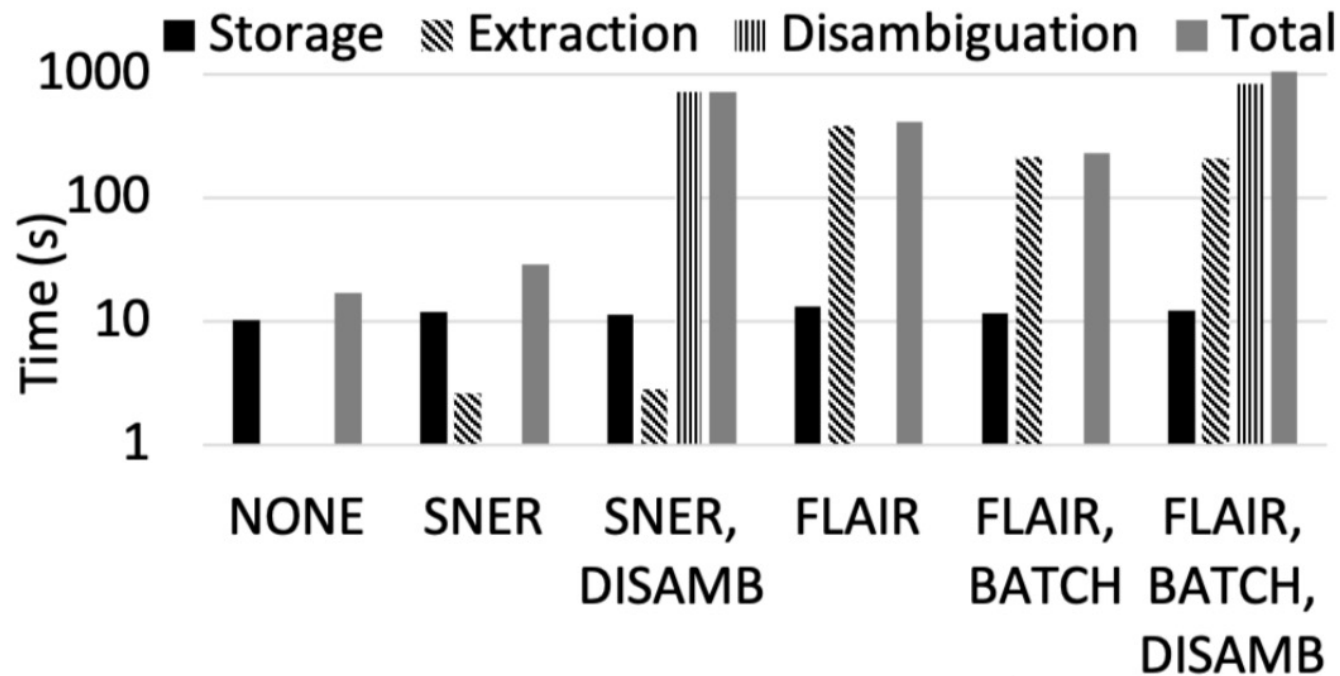


Figure 6: Graph construction time (seconds).

Graph creation performance: batch extraction [InfSys2022]

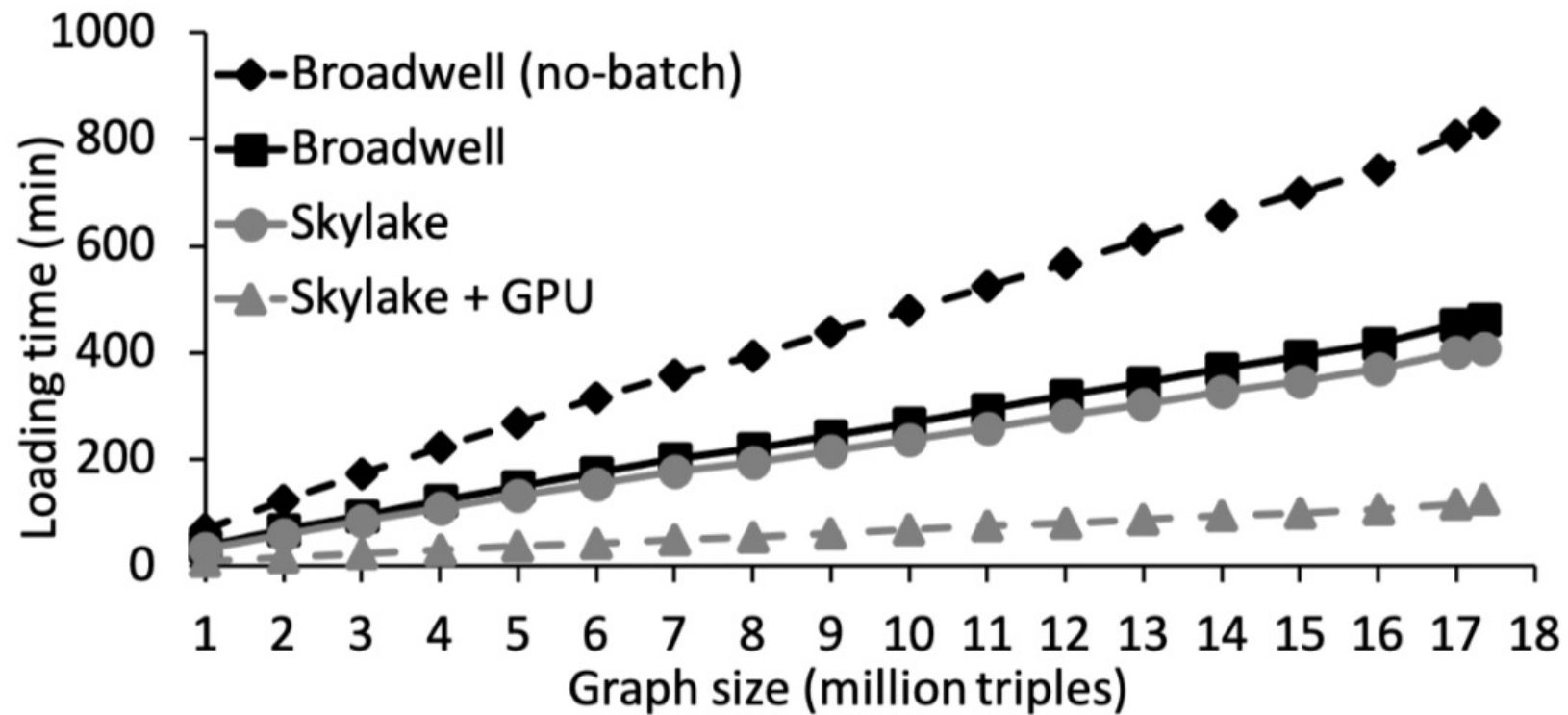


Figure 7: YAGO loading time (minutes) using Flair.

Credits and references <https://sourcessay.inria.fr>

ConnectionLens core architecture: A. Anadiotis, O. Balalau, H. Galhardas, J. Leblay, I. Manolescu, T. Merabti, Y. Haddad, J. You, Y. Youssef [VLDB18, InfSys22, CIKM21demo]

Dataset abstraction: N. Barret, I. Manolescu, P. Upadhyay [CIKM22demo, BDA2022, 2 under submission]

Entity path finding: N. Barret, A. Gauquier, J. Law, I. Manolescu [1 under submission]

Graph exploration: T. Bouganim, I. Manolescu, E. Pietriga [work in progress]

Structure and unstructured querying: A. Anadiotis, Y. Haddad, I. Manolescu, M. Mohanty [IEEE DEBull21, ICDE23, 1 under submission]

Biomedical Conflicts of Interest application: A. Anadiotis, O. Balalau, T. Bouganim, G. Fooks, H. Galhardas, S. Horel, I. Manolescu, T. Mills, C. Pettineo, P. Upadhyay [IEEE DEBull21, EACL23, 1 under preparation]

Looking forward

Perspectives

SourcesSay: 2-3 more years to go

- ❑ Finalizing abstraction + finding interesting entity paths (PhD N. Barret)
- ❑ Exploratory querying (PhD T. Bouganim)

Inria Exploratory Action JoDaIA (“Data and AI for Journalism”) (2023-2026)

- ❑ One axis: identifying and developing more applications based on ConnectionLens
- ❑ *Open to exploring more proposals!*

ELIAS EU project (coord. U. Trento, with O. Goga): Development of learning algorithms to protect and secure democracy

Thank you

Questions?

SourcesSay: <https://sourcessay.inria.fr>

ConnectionLens: <https://team.inria.fr/cedar/connectionlens/>