



From Tabular Data to Knowledge Graphs

Ernesto Jiménez-Ruiz

Lecturer in Artificial Intelligence

Before we start...

Some important information

- Drop-in session
 - Today from **2:30pm to 3:30pm.**

Some important information

- Drop-in session
 - Today from **2:30pm to 3:30pm.**
- Reading week.
 - Thursday **March 7, 10am-12pm:** Coursework Part 2 discussion + Questions + Lab session.
 - Online (usual zoom room) and recorded.

Data Bites Seminar

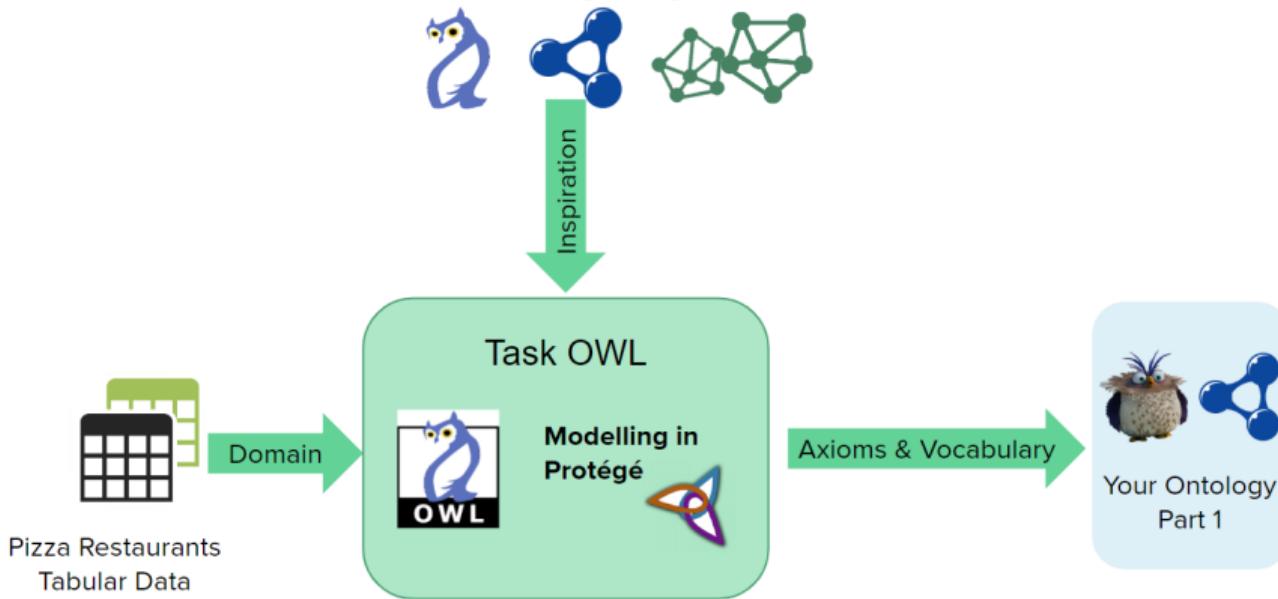
- Seminar on "**Get set for Graphs!**".
- **Semantic Partners Ltd** (<https://www.semanticpartners.com/>)
- They may bring ideas for projects, internships, etc.
- When: Today, 1:30pm
- Where: C322, Tait building
- They offer coffee and cookies.

Interest Group on KGs: The Alan Turing Institute

- UK community: academia, industry, government, etc.
- Next meet-up on March 25, Liverpool
 - **Knowledge Graphs in the Wild**
- <https://www.turing.ac.uk/research/interest-groups/knowledge-graphs>

Coursework Part 1: Ontology modelling (i)

External Resources, KGs, Ontologies (e.g., Pizza.owl, DBpedia)



(*) Ontology CW2 = Model solution \neq Your Ontology Part 1

Coursework Part 1: Ontology modelling (ii)

- Part 1 (20%): **creation of an ontology** that covers the knowledge of a given domain. **Deadline:** Sunday, 3 March 2024, 5:00 PM
- Please **follow the instructions** in the submission guidelines.
- Please **do not submit last minute!** Nor last second!

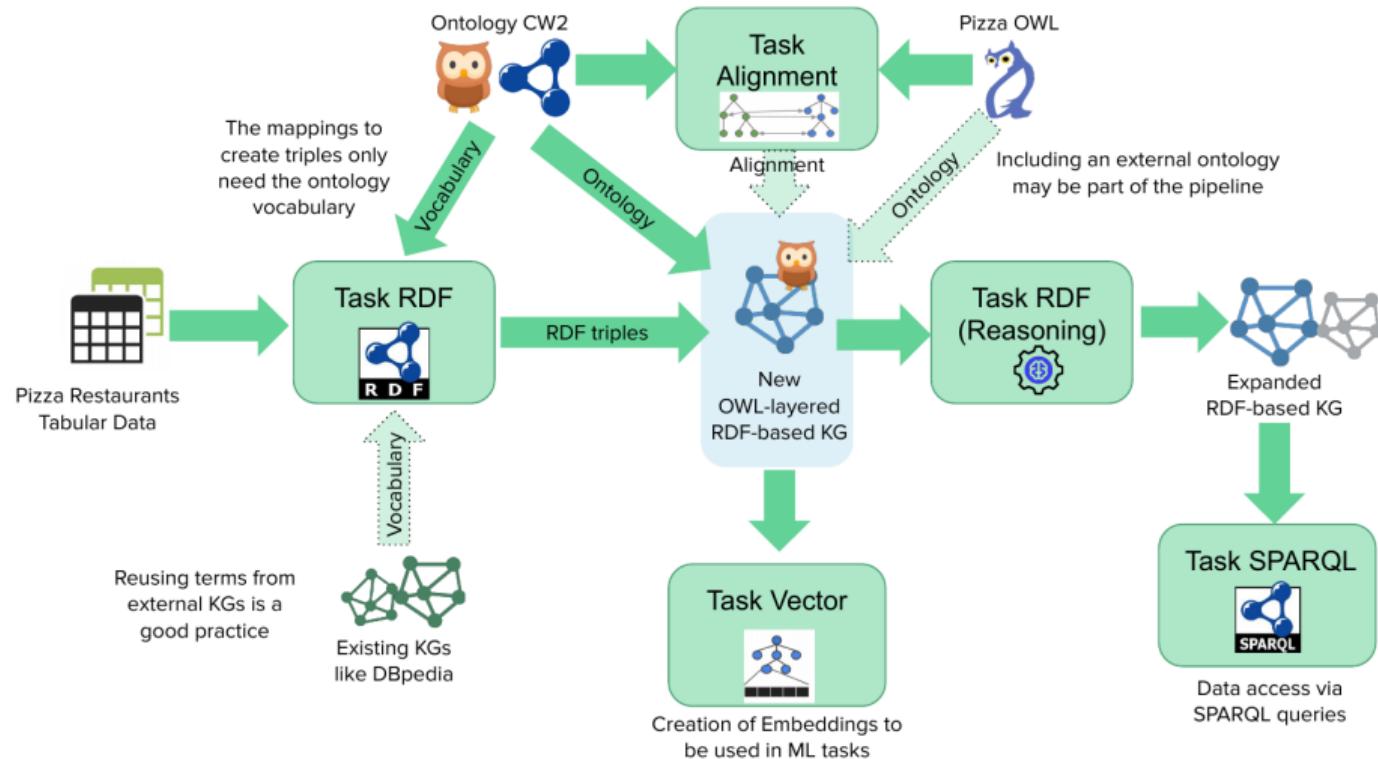
Where are we? Module organization.

- ✓ Introduction: Becoming a knowledge scientist.
- ✓ RDF-based knowledge graphs.
- ✓ OWL ontology language. Focus on modelling.
- ✓ SPARQL 1.0 Query Language.

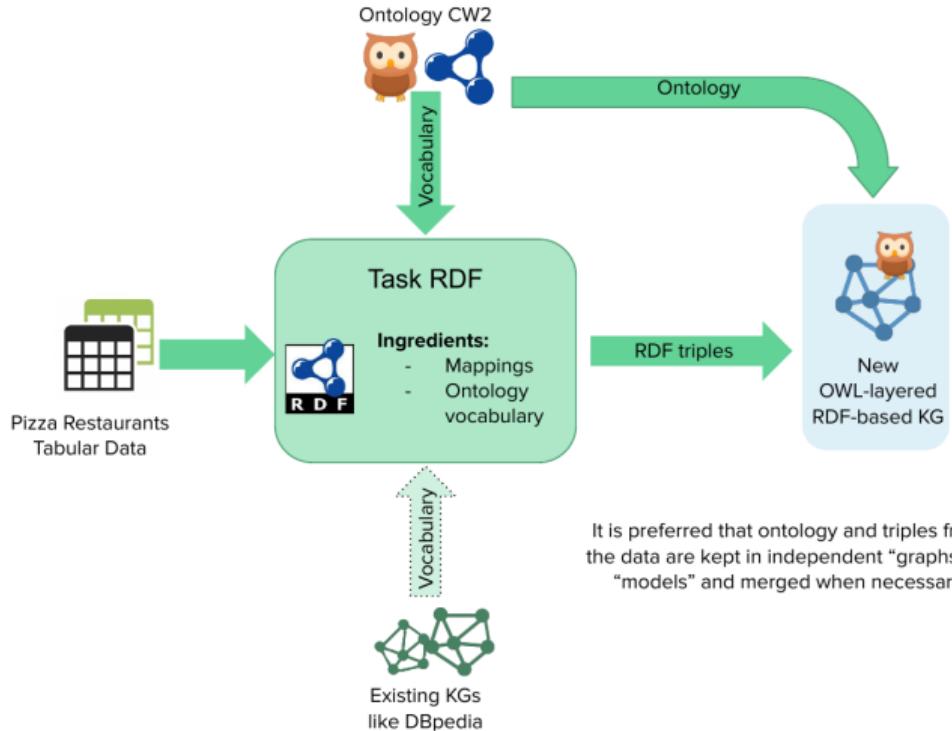
5. **From tabular data to KG.**(Today)
6. RDFS Semantics and OWL 2 profiles.
7. Ontology Alignment.
8. Ontology (KG) Embeddings and Machine Learning.
9. SPARQL 1.1 and Graph Database solutions.
10. (Large) Language Models and KGs. (Seminar)

The global picture

The global picture (Coursework Part 2)



Global picture: Task RDF (Today)



Data Science Bottleneck

Data Science Bottleneck



Big Data Borat

@BigDataBorat



Follow

In Data Science, 80% of time spent prepare
data, 20% of time spent complain about need
for prepare data.

Data Science Bottleneck



Big Data Borat

@BigDataBorat



Follow

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

One of the main problems is the lack of understanding of the data and the domain.

Contribution of Semantics in Data Wrangling Challenges

- *Data parsing*, e.g. converting csv's or tables.
- (+++) *Data dictionary*: basic types and semantic types.
- (++) *Data integration* from multiple sources (foreign key discovery).
- (++) *Entity resolution*: duplication and record linkage.
- (+) *Format variability*: e.g. for dates and names.
- (+) *Structural variability* in the data.
- (++) Identifying and repairing *missing data*.
- (+) *Anomaly detection* and repair.
- (+++) **Metadata/contextual information**. (Semantic) data governance.

AIDA Project: <https://www.turing.ac.uk/research/research-projects/artificial-intelligence-data-analytics-aida>

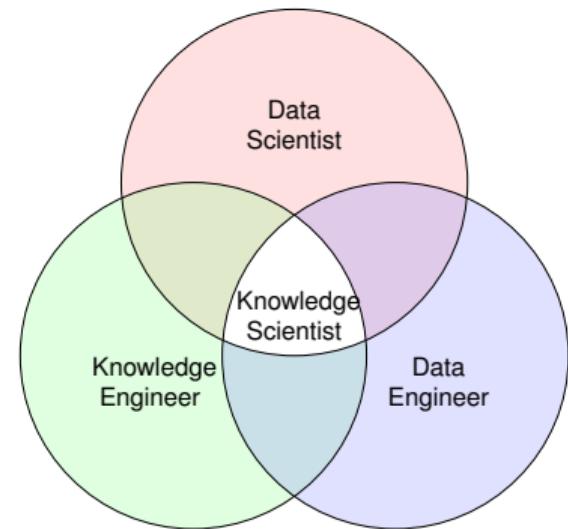
Semantic Understanding of Tabular Data

- **Semi-automatic** process.
- Key for an **enhanced transformation** to RDF triples.
- But also for other tasks with independence of a final KG creation.
 - Tabular data in the form of CSV files is the common input format in a **data analytics pipeline**.
 - The **lack of semantics and context in datasets** hinders their usability.
 - Gaining **semantic understanding** will be very valuable for data integration, data cleaning, data mining, machine learning and knowledge discovery tasks.

The Knowledge Scientist

The Knowledge Scientist role (i)

- **Data Engineer (role):** harnesses and collects data.
- **Data Scientist (role):** draws value from data.
- **Knowledge Engineer (role):** encodes domain expertise.
- **Knowledge Scientist (role):** adds context to the data to make it more useful, clean, reliable and ready to be used.



The Knowledge Scientist role (ii)

- Bridges the data and the **business requirements**/questions.
- Outputs a data model (*i.e.*, **knowledge graph**): how business users see the world.
- Drives a **semantic-lifting** of the data (from Data Engineers to Data Scientists)
- Relies on **Semantic Web technology** and skills (e.g., ontology modelling, data/knowledge integration, graph databases)

George Fletcher and others. **Knowledge Scientists: Unlocking the data-driven organization**. 2020

The Knowledge Scientist: DS perspective

- Domain understanding.
- Designing & developing ontologies and KGs together with domain experts.
- Developing data processing pipelines to drive the semantic shifting of the data.
- Application of data analytics and ML with KGs.

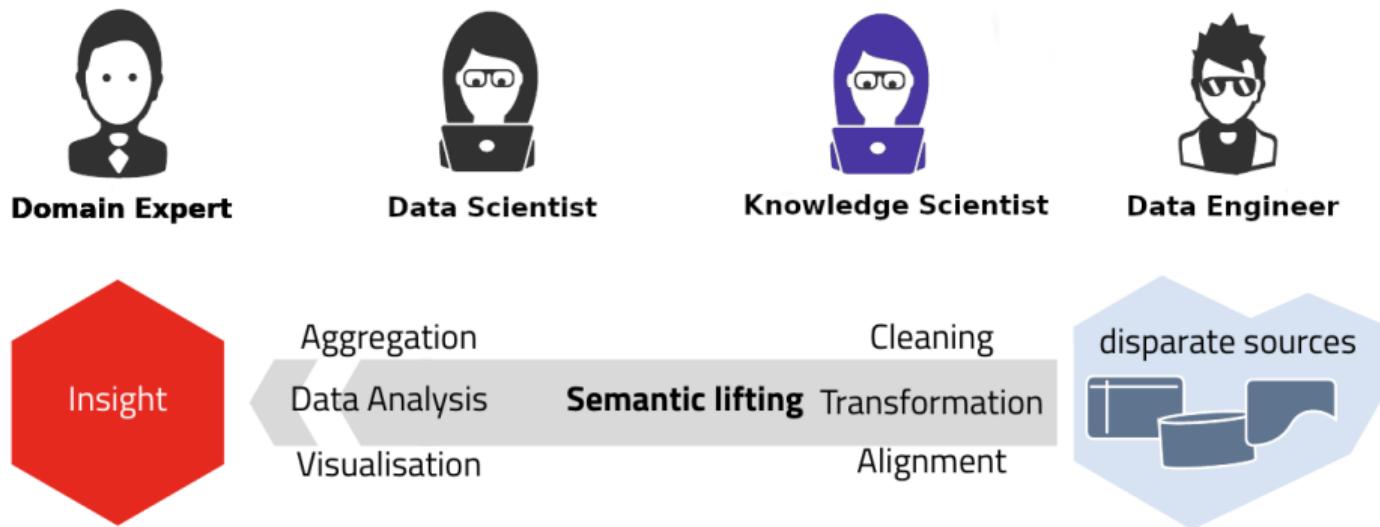
George Fletcher and others. **Knowledge Scientists: Unlocking the data-driven organization.** 2020
Stratos Kontopoulos. **Semantic AI & the Role of the Knowledge Scientist.** 2021

The Knowledge Scientist: SE/CS perspective

- Domain understanding.
- Deploying the semantic infrastructure (e.g., triplestores, reasoning, query formulation, required APIs).
- Integration with other components.
- Semantic infrastructure in the cloud.

George Fletcher and others. **Knowledge Scientists: Unlocking the data-driven organization**. 2020
Stratos Kontopoulos. **Semantic AI & the Role of the Knowledge Scientist**. 2021

The Knowledge Scientist: overview



Adapted from: SIRIUS Centre for Scalable Data Access, <https://sirius-labs.no/>

Why Ontologies and Graphs of Knowledge?

Graph(s) of Knowledge / Knowledge Graphs

- Semantic Web in more **controlled scenarios**, e.g.,
 - **Integrate and orchestrate** data within an organisation
 - Enterprise data as a knowledge graph to **drive products** and make them more “**intelligent**”

Graph(s) of Knowledge / Knowledge Graphs

- Semantic Web in more **controlled scenarios**, e.g.,
 - **Integrate and orchestrate** data within an organisation
 - Enterprise data as a knowledge graph to **drive products** and make them more “**intelligent**”
- **Not new:**
 - Graph data models extensively studied in AI...
 - ...but Google has relaunched the interest on **KGs in industry**

Graph(s) of Knowledge / Knowledge Graphs

- Semantic Web in more **controlled scenarios**, e.g.,
 - **Integrate and orchestrate** data within an organisation
 - Enterprise data as a knowledge graph to **drive products** and make them more “**intelligent**”
- **Not new:**
 - Graph data models extensively studied in AI...
 - ...but Google has relaunched the interest on **KGs in industry**
- Availability of **mature** Semantic Web **technology**
 - Query engines
 - Modelling languages
 - Reasoning

Ontologies and Knowledge Graphs

- Core idea of knowledge graphs is the enhancement of the graph data model with...
- “...an **abstract symbolic representations** of a domain expressed in a formal language”
- In this module: **OWL-layered RDF-based knowledge graphs**

Aidan Hogan and others. **Knowledge Graphs**. 2021 <https://kgbook.org/>.
Pim Borst, Hans Akkermans, and Jan Top. **Engineering ontologies**, 1999.

Why Ontologies and KGSSs?

- Independence of logical/physical schema: **domain model**
- Vocabulary closer to domain experts: **more user-friendly**
- Incomplete and semi-structured data: **flexibility**
- Integration of heterogeneous sources: **unified view**

♠ They can complement tabular data not necessarily substitute.

Why Ontologies and KGs? (Experiences from Industry)

- Knowledge graphs are **not complex**: they are the simplest way to layer human expertise onto stored data.

Juan Sequeda and Ora Lassila. **Designing and Building Enterprise Knowledge Graphs**. 2021

Ora Lassila. **On the broad applicability of Semantic Web technologies**. (COST DKG Talk Series. 2021

<https://www.youtube.com/watch?v=f9wautaqWUs>

Juan Sequeda and Ora Lassila. <https://watch.knowledgegraph.tech/knowledge-espresso/videos/knowledge-espresso-with-ora-lassila-and-juan-sequeda>

Why Ontologies and KGs? (Experiences from Industry)

- Knowledge graphs are **not complex**: they are the simplest way to layer human expertise onto stored data.
- Investing in knowledge graphs is **resilient**: an initial investment that pays off by allowing companies to pre-prepare for unanticipated and unfolding use cases

Juan Sequeda and Ora Lassila. **Designing and Building Enterprise Knowledge Graphs**. 2021

Ora Lassila. **On the broad applicability of Semantic Web technologies**. (COST DKG Talk Series. 2021

<https://www.youtube.com/watch?v=f9wautaqWUs>

Juan Sequeda and Ora Lassila. <https://watch.knowledgegraph.tech/knowledge-espresso/videos/knowledge-espresso-with-ora-lassila-and-juan-sequeda>

Why Ontologies and KGs? (Experiences from Industry)

- Knowledge graphs are **not complex**: they are the simplest way to layer human expertise onto stored data.
- Investing in knowledge graphs is **resilient**: an initial investment that pays off by allowing companies to pre-prepare for unanticipated and unfolding use cases
- **Unifying logical representation** for data and semantics to get out the data science enterprise mess.

Juan Sequeda and Ora Lassila. **Designing and Building Enterprise Knowledge Graphs**. 2021

Ora Lassila. **On the broad applicability of Semantic Web technologies**. (COST DKG Talk Series. 2021

<https://www.youtube.com/watch?v=f9wautaqWUs>

Juan Sequeda and Ora Lassila. <https://watch.knowledgegraph.tech/knowledge-espresso/videos/knowledge-espresso-with-ora-lassila-and-juan-sequeda>

FAIR principles and 5-star data

Why Ontologies and KGs?

5-star Data

- ★ **OL:** make your data available on the Web (in any format) under an open license.

Tim Berners-Lee. **5 * (open) data:** <https://5stardata.info/en/>

5-star Data

- ★ **OL:** make your data available on the Web (in any format) under an open license.
- ★★ **RE:** make the data machine readable (excel instead of an scanned image).

Tim Berners-Lee. **5 * (open) data:** <https://5stardata.info/en/>

5-star Data

- ★ **OL:** make your data available on the Web (in any format) under an open license.
- ★★ **RE:** make the data machine readable (excel instead of an scanned image).
- ★★★ **OF:** use a non proprietary open format (e.g., CSV).

Tim Berners-Lee. **5 * (open) data:** <https://5stardata.info/en/>

5-star Data

- ★ **OL:** make your data available on the Web (in any format) under an open license.
- ★★ **RE:** make the data machine readable (excel instead of an scanned image).
- ★★★ **OF:** use a non proprietary open format (e.g., CSV).
- ★★★ **URI:** use URIs instead of strings (RDF).

Tim Berners-Lee. **5 * (open) data:** <https://5stardata.info/en/>

5-star Data

- ★ **OL:** make your data available on the Web (in any format) under an open license.
- ★★ **RE:** make the data machine readable (excel instead of an scanned image).
- ★★★ **OF:** use a non proprietary open format (e.g., CSV).
- ★★★★ **URI:** use URLs instead of strings (RDF).
- ★★★★★ **LOD:** link your data to other data to provide extended context.

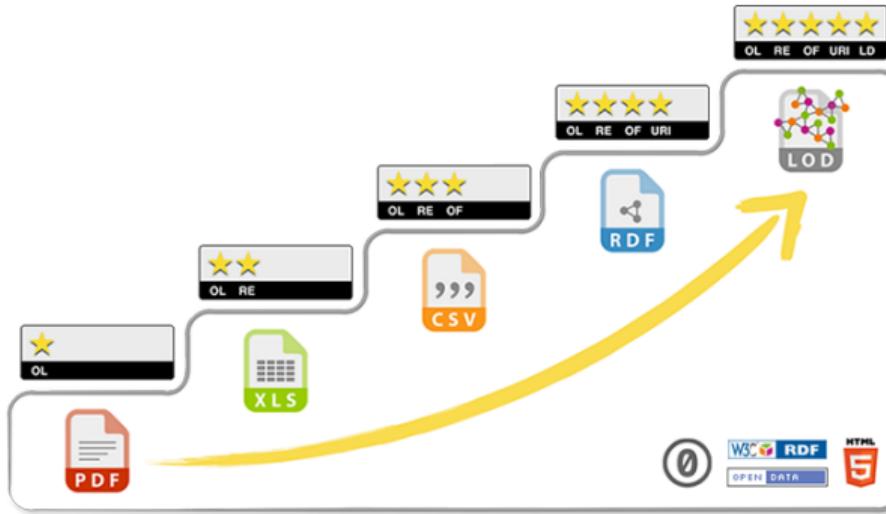
Tim Berners-Lee. **5 * (open) data:** <https://5stardata.info/en/>

5-star Data

- ★ **OL:** make your data available on the Web (in any format) under an open license.
 - ★★ **RE:** make the data machine readable (excel instead of an scanned image).
 - ★★★ **OF:** use a non proprietary open format (e.g., CSV).
 - ★★★★ **URI:** use URIs instead of strings (RDF).
 - ★★★★★ **LOD:** link your data to other data to provide extended context.
- ♠ This could be applied **within an organisation** (intranet), not only for the Web. Ideally with an OL, but at least data accessible by everyone in the organisation.

Tim Berners-Lee. **5 * (open) data:** <https://5stardata.info/en/>

5-star Data



♠ **LOD:** Linked Open Data.

Tim Berners-Lee. 5 * (open) data: <https://5stardata.info/en/>

5-star Data: Technical challenges:

- How to **expose** data (e.g., databases, csv files) as knowledge graphs?
- How to **create** (or reuse) and use (abstract) **knowledge** (i.e., *Ontologies*)?
- How to **align** different knowledge graphs? ♠
- How to check **consistency and trust** of the data and knowledge? ♠

♠ Better with things than with strings

5-star Data: Technical challenges:

- How to **expose** data (e.g., databases, csv files) as knowledge graphs?
 - *RDF (Week 2) and Today's session - 4-5★ data*
- How to **create** (or reuse) and use (abstract) **knowledge** (i.e., *Ontologies*)?
 - *OWL (Week 3)*
- How to **align** different knowledge graphs? ♠
 - *Ontology Alignment (Week 8) - 5★ data*
- How to check **consistency and trust** of the data and knowledge? ♠
 - *Reasoning with RDFS and OWL (Week 7) - 6★ data?*

♠ Better with things than with strings

FAIR Data Principles (i)

F

indable

A

ccessible

I

nteroperable

R

eusable



- ♠ For machines and people.

The FAIR Guiding Principles for scientific data management and stewardship. Nature Scientific Data 2016.

<https://www.go-fair.org/fair-principles/>

FAIR Data Principles (ii)

- **Findable**: (meta)data is assigned a unique identifier and data is described with rich metadata.
- **Accessible**: (meta)data can be accessed via a known protocol. (*)
- **Interoperable**: meta(data) uses a formal language.
- **Reusable**: well described (meta)data with rich provenance.

(*) Not necessarily open. (FAIR \neq OPEN)

FAIR Data Principles (ii)

- **Findable**: (meta)data is assigned a unique identifier and data is described with rich metadata. ([URIs](#))
- **Accessible**: (meta)data can be accessed via a known protocol. ([SPARQL](#)) (*)
- **Interoperable**: meta(data) uses a formal language. ([Knowledge representation with OWL](#))
- **Reusable**: well described (meta)data with rich provenance. ([W3C standards](#))

(*) Not necessarily open. (FAIR ≠ OPEN)

FAIR Data Principles (iii)

- FAIR data is more valuable, easier to find and combine thanks to unique identifiers and a formal shared knowledge representation.
- “KGs must be in want of FAIR data. And FAIR data is in want of KGs”.

Carole Golbe. FAIRy stories: the FAIR Data principles in theory and in practice.

<https://www.slideshare.net/carolegoble/fairy-stories-the-fair-data-principles-in-theory-and-in-practice>

From (Tabular) Data to Knowledge Graphs: Towards 5 ★ and FAIR data

Semantic Table Understanding/Interpretation



Vincenzo Cutrona. **Why Table Understanding Matters.** See the module's invited speakers in moodle.

Semantic Table Understanding/Interpretation



Vincenzo Cutrona. **Why Table Understanding Matters.** See the module's invited speakers in moodle.

Exposing data as an RDF-based Knowledge Graph

- ✓ **End-users' friendly access** to “unfriendly” tabular data.
- ✓ **Pay as you go** (modular) data integration via mappings.

Exposing data as an RDF-based Knowledge Graph

- ✓ End-users' friendly access to "unfriendly" tabular data.
- ✓ Pay as you go (modular) data integration via mappings.
- Option 1: Virtual exposure of data (OBDA)
 - ✓ Data remains in its original format.
 - ✗ Typically only over relational databases.

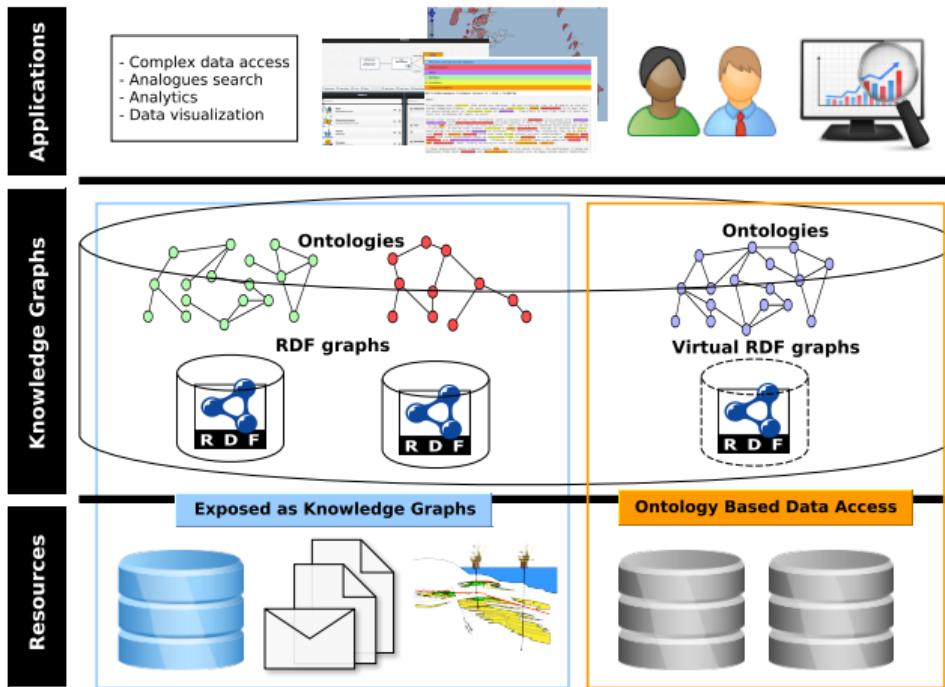
Exposing data as an RDF-based Knowledge Graph

- ✓ End-users' friendly access to "unfriendly" tabular data.
- ✓ Pay as you go (modular) data integration via mappings.
- Option 1: Virtual exposure of data (OBDA)
 - ✓ Data remains in its original format.
 - ✗ Typically only over relational databases.
- Option 2: Data Export/Materialization
 - ✓ Easy to exchange data (RDF).
 - ✓ Integration of data in disparate formats.
 - ✗ Data replication.
 - Due to size or privacy it may not be possible to export the data.

Exposing data as an RDF-based Knowledge Graph

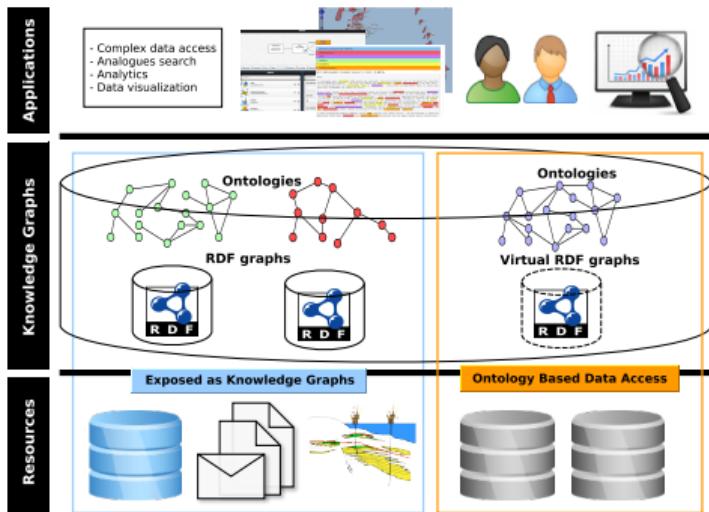
- ✓ End-users' friendly access to "unfriendly" tabular data.
- ✓ Pay as you go (modular) data integration via mappings.
- Option 1: Virtual exposure of data (OBDA)
 - ✓ Data remains in its original format.
 - ✗ Typically only over relational databases.
- Option 2: Data Export/Materialization (In this module)
 - ✓ Easy to exchange data (RDF).
 - ✓ Integration of data in disparate formats.
 - ✗ Data replication.
 - Due to size or privacy it may not be possible to export the data.

Exposing data as RDF: Architecture



Exposing data as RDF: Ingredients

- **Ontology vocabulary.** Custom and/or given by a public KG.
- **Mappings.** Define a transformation from the tabular data to RDF data.
- **Ontology Axioms (optional)** - ♠



♣ Ernesto Jimenez-Ruiz and others. **BootOX: Practical Mapping of RDBs to OWL 2.** ISWC 2015

Exposing data as RDF: W3C Mapping Standards

- **Relational Database to RDF:**
 - A Direct Mapping of Relational Data to RDF:
<https://www.w3.org/TR/rdb-direct-mapping/>
 - R2RML: RDB to RDF Mapping Language: <https://www.w3.org/TR/r2rml/>
 - Each mapping involves the creation of a **SQL query** and the transformation of the results to RDF triples.

Exposing data as RDF: W3C Mapping Standards

- **Relational Database to RDF:**
 - A Direct Mapping of Relational Data to RDF:
<https://www.w3.org/TR/rdb-direct-mapping/>
 - R2RML: RDB to RDF Mapping Language: <https://www.w3.org/TR/r2rml/>
 - Each mapping involves the creation of a **SQL query** and the transformation of the results to RDF triples.
- **CSV to RDF:**
 - Direct mapping from CSV to RDF (CSV2RDF): <https://www.w3.org/TR/csv2rdf/>
 - General (declarative) RDF Mapping Language (RML):
<https://rml.io/specs/rml/>
 - A mapping can also be seen as a (**small**) **script** that creates specific RDF triples from the CSV file (e.g., data frame).

Exposing data as RDF: W3C Mapping Standards

- **Relational Database to RDF:**
 - A Direct Mapping of Relational Data to RDF:
<https://www.w3.org/TR/rdb-direct-mapping/>
 - R2RML: RDB to RDF Mapping Language: <https://www.w3.org/TR/r2rml/>
 - Each mapping involves the creation of a **SQL query** and the transformation of the results to RDF triples.
- **CSV to RDF:**
 - Direct mapping from CSV to RDF (CSV2RDF): <https://www.w3.org/TR/csv2rdf/>
 - General (declarative) RDF Mapping Language (RML):
<https://rml.io/specs/rml/>
 - A mapping can also be seen as a **(small) script** that creates specific RDF triples from the CSV file (e.g., data frame). **In this module.**

Exposing data as RDF: Direct Mapping Example

Automatic triples:

```
ex:row1 ex:col1 "China"  
ex:row1 ex:col2 "Beijing"  
ex:row2 ex:col1 "Indonesia"  
ex:row2 ex:col2 "Jakarta"  
...
```

China	Beijing
Indonesia	Jakarta
Congo	Kinshasa
Brazil	
Congo	Brazzaville

(*) ex: is the prefix defined for the namespace <http://example.org/>

Exposing data as RDF: Direct Mapping Example

Automatic triples:

```
ex:row1 ex:col1 "China"  
ex:row1 ex:col2 "Beijing"  
ex:row2 ex:col1 "Indonesia"  
ex:row2 ex:col2 "Jakarta"  
...
```

(we can probably do better)

China	Beijing
Indonesia	Jakarta
Congo	Kinshasa
Brazil	
Congo	Brazzaville

(*) ex: is the prefix defined for the namespace <http://example.org/>

Exposing data as RDF: Enhanced Mapping/Transformation (i)

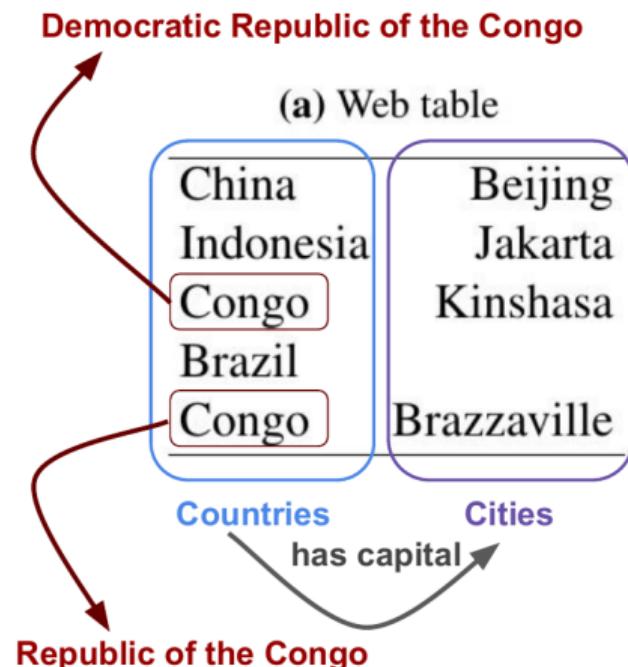
- If we know the **semantics** of the data.
- **Potential automatic triples:**

ex:China rdf:type ex:Country

ex:Beijing rdf:type ex:City

ex:China ex:hasCapital ex:Beijing

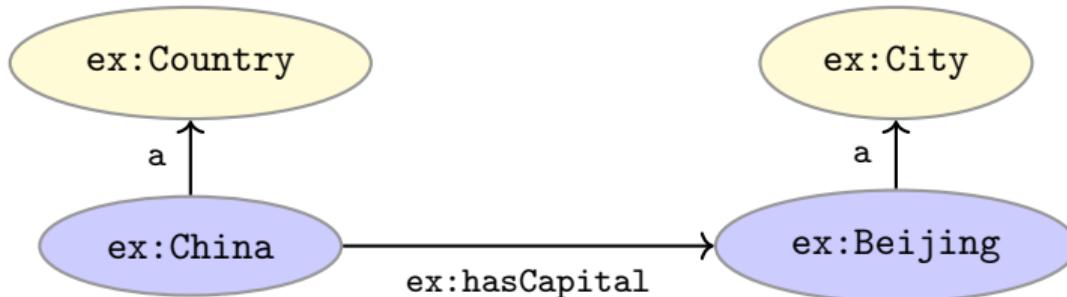
...



Exposing data as RDF: Enhanced Mapping/Transformation (ii)

Return capital of China (for the KG \mathcal{G} below):

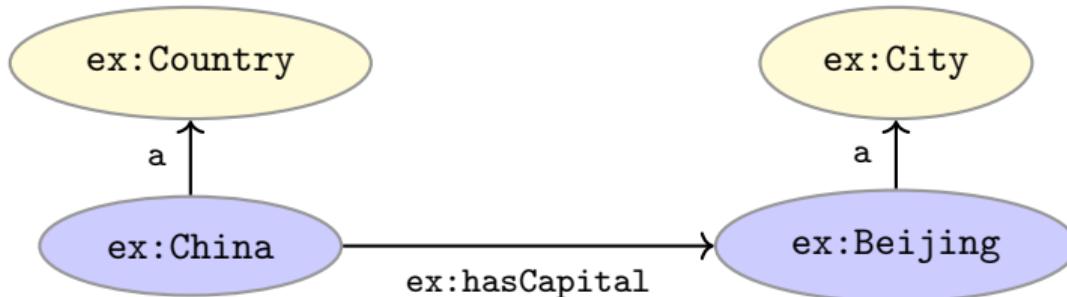
```
PREFIX ex: <http://example.org/>
SELECT DISTINCT ?capital WHERE {
    ex:China ex:hasCapital ?capital .
}
```



Exposing data as RDF: Enhanced Mapping/Transformation (ii)

Return capital of China (for the KG \mathcal{G} below): Query Result= {ex:Beijing}

```
PREFIX ex: <http://example.org/>
SELECT DISTINCT ?capital WHERE {
    ex:China ex:hasCapital ?capital .
}
```



Exposing data as RDF: Mappings

Mapping or Transformation $\varphi \rightsquigarrow \psi$

- φ : query over database or CSV extraction
- ψ : RDF triple template

Exposing data as RDF: Mappings

Mapping or Transformation $\varphi \rightsquigarrow \psi$

- φ : query over database or CSV extraction
- ψ : RDF triple template
- RDB to RDF mapping:

```
SELECT col1 FROM TABLE ~> ex:{col1} rdf:type ex:Country
```

Exposing data as RDF: Mappings

Mapping or Transformation $\varphi \rightsquigarrow \psi$

- φ : query over database or CSV extraction
- ψ : RDF triple template
- RDB to RDF mapping:

```
SELECT col1 FROM TABLE ~> ex:{col1} rdf:type ex:Country
```

- CSV to RDF mapping:

```
for value in data_frame[col1]:                                ( $\varphi$ )
    subject = "ex:" + value      #e.g., ex:China
    create_triple(subject rdf:type ex:Country)   ( $\psi$ )
```

Exposing data as RDF: Mappings

Mapping or Transformation $\varphi \rightsquigarrow \psi$

- φ : query over database or CSV extraction
- ψ : RDF triple template
- RDB to RDF mapping:

```
SELECT col1 FROM TABLE ~> ex:{col1} rdf:type ex:Country
```

- **CSV to RDF mapping:** (In this module)

```
for value in data_frame[col1]:                                ( $\varphi$ )
    subject = "ex:" + value      #e.g., ex:China
    create_triple(subject rdf:type ex:Country)   ( $\psi$ )
```

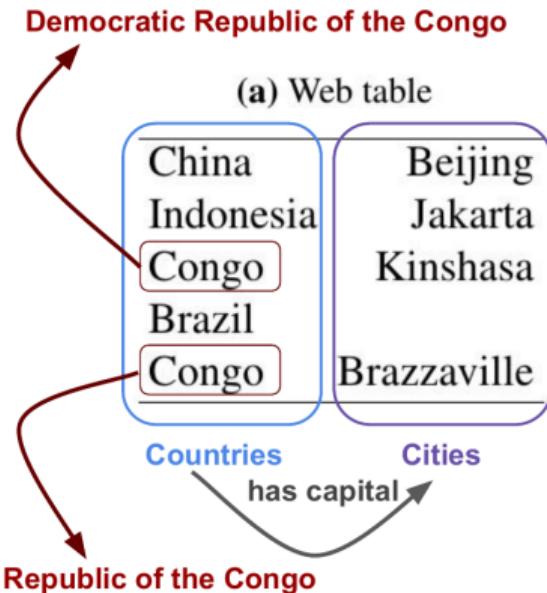
Exposing data as RDF: (Manually defined) Mappings

- Column 1 is composed by Countries.

```
for value in data_frame[col1]:  
    subject = "ex:" + value #e.g., ex:China  
    create_triple(subject rdf:type ex:Country)
```

- Column 2 entities are the capitals of column 1 entities.

```
for row in data_frame:  
    subj = "ex:" + row[col1] #e.g., ex:China  
    obj = "ex:" + row[col2] #e.g., ex:Beijing  
    create_triple(subj ex:hasCapital obj)
```



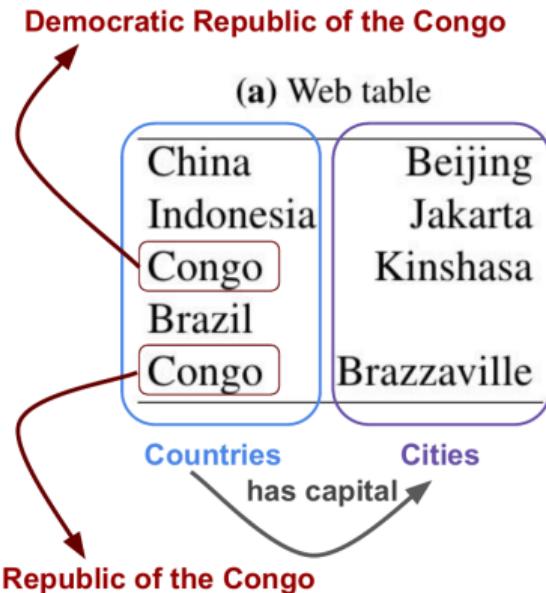
Exposing data as RDF: (Manually defined) Mappings

- Column 1 is composed by Countries.

```
for value in data_frame[col1]:  
    subject = "ex:" + value #e.g., ex:China  
    create_triple(subject rdf:type ex:Country)
```

- Column 2 entities are the capitals of column 1 entities.

```
for row in data_frame:  
    subj = "ex:" + row[col1] #e.g., ex:China  
    obj = "ex:" + row[col2] #e.g., ex:Beijing  
    create_triple(subj ex:hasCapital obj)
```



Target in Coursework Part 2

Semantic Understanding of Tabular Data: Automating the Process (suitable for a MSc project)

Adding Semantics to Tabular Data: Basic Tasks

- Matching a cell to a KG entity (**CEA task** - Cell-Entity Annotation)
- Assigning a semantic type (e.g., a KG class) to an (entity) column (**CTA task** - Column-Type Annotation)
- Assigning a KG property to the relationship between two columns (**CPA task** - Columns-Property Annotation)

Ernesto Jiménez-Ruiz and others. **SemTab 2019: Resources to Benchmark Tabular Data to Knowledge Graph Matching Systems.** ESWC 2020

Adding Semantics to Tabular Data: Basic Tasks

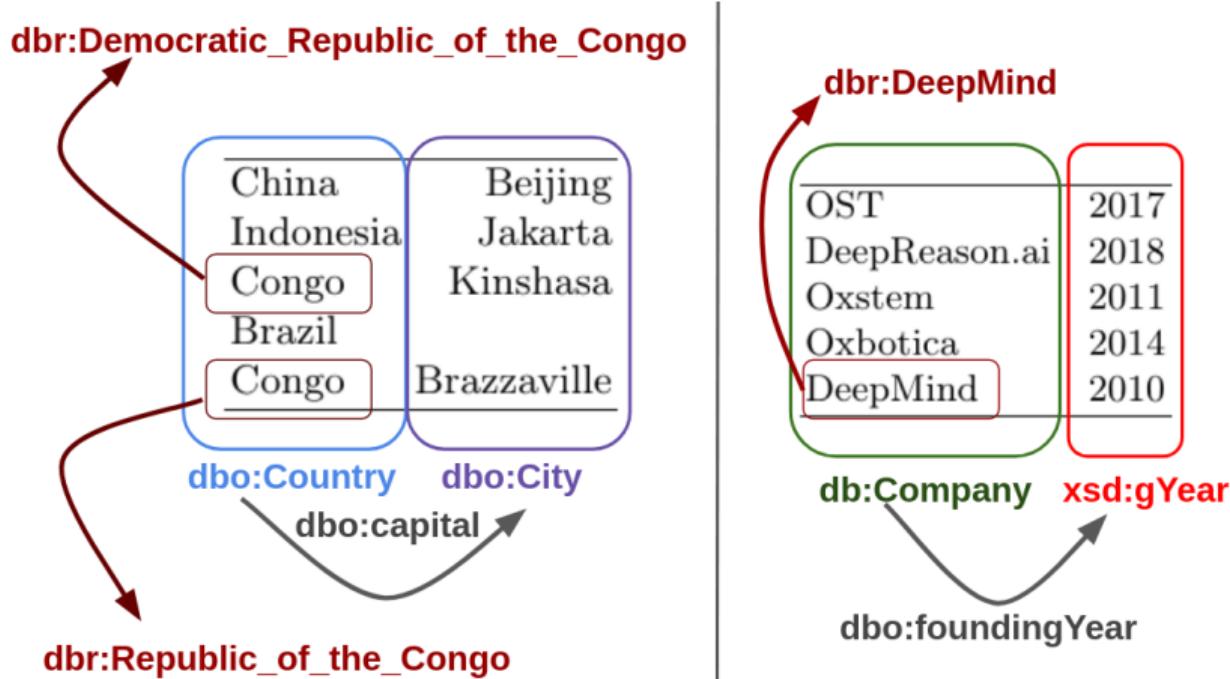
- Matching a cell to a KG entity (**CEA task** - Cell-Entity Annotation)
- Assigning a semantic type (e.g., a KG class) to an (entity) column (**CTA task** - Column-Type Annotation)
- Assigning a KG property to the relationship between two columns (**CPA task** - Columns-Property Annotation)

† For a semi-automatic process, we assume the existence of a (possibly incomplete) **Knowledge Graph (KG)** relevant to the domain.

‡ When transforming to RDF, if no KG matching then create a fresh entity URI.

Ernesto Jiménez-Ruiz and others. **SemTab 2019: Resources to Benchmark Tabular Data to Knowledge Graph Matching Systems**. ESWC 2020

Adding Semantics to Tabular Data: Basic Tasks (with DBpedia)



SemTab Challenge

- **Systematic evaluation** of Tabular Data to KG matching systems.
- Evaluates the **three basic tasks**: CTA, CEA and CPA.
- Relies on:
 - an **automatic** dataset generator, and
 - **manually curated datasets**.
- **Target KGs**: DBpedia, Wikidata and Schema.org.
- Co-organised and sponsored by **IBM Research**.

SemTab: Semantic Web Challenge on Tabular Data to Knowledge Graph Matching. Collocated with the International Semantic Web Conference (ISWC): <http://www.cs.ox.ac.uk/isg/challenges/sem-tab/>

Semantic Understanding of Tabular Data: Techniques

Common Techniques

- **Pre-processing**: spelling error, stopwords, unicode fixing, etc.
- **Regular expressions** to identify data formats (e.g., numbers, phones, dates, names).
- **NLP techniques** to process textually-rich fields.

Common Techniques

- **Pre-processing**: spelling error, stopwords, unicode fixing, etc.
- **Regular expressions** to identify data formats (e.g., numbers, phones, dates, names).
- **NLP techniques** to process textually-rich fields.
- **Fuzzy search** over a KG
 - Via online look-up services
 - Or local indexes
- Access to the **KG's SPARQL Endpoint** (local or online)

Common Techniques

- **Pre-processing**: spelling error, stopwords, unicode fixing, etc.
- **Regular expressions** to identify data formats (e.g., numbers, phones, dates, names).
- **NLP techniques** to process textually-rich fields.
- **Fuzzy search** over a KG
 - Via online look-up services
 - Or local indexes
- Access to the **KG's SPARQL Endpoint** (local or online)
- **Lexical similarity** (e.g., Levenshtein)
- Word and KG **embeddings**. And **LLMs!**

Common Knowledge Graphs

Wikidata: <https://www.wikidata.org/>

- >100 million entities
- Free and public (anyone can edit)

DBpedia: <https://dbpedia.org/> (**Extracted from Wikipedia**)

- >100 million entities
- >900 million triples

Google KG: <https://developers.google.com/knowledge-graph>

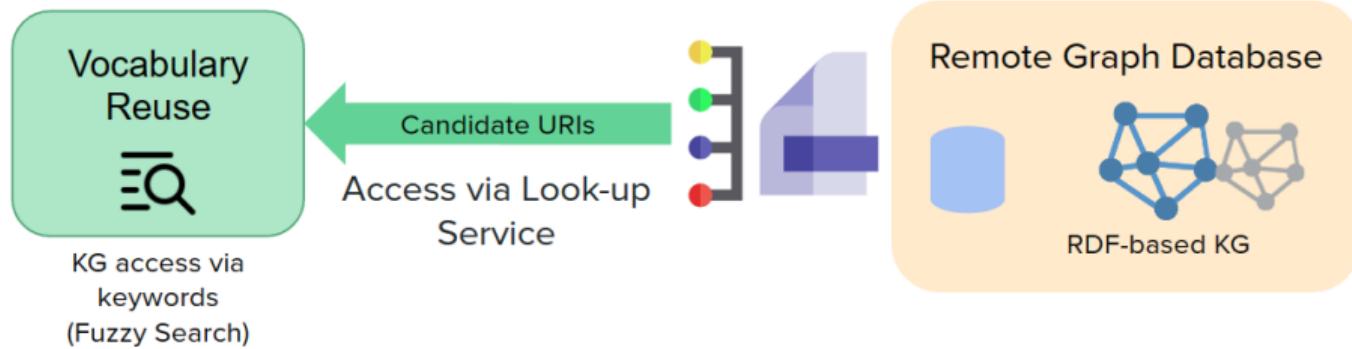
- Private, only accessible via look-up
- >5 billion entities

Fuzzy Search: KG look-up Services

- Given a string (e.g., “Congo”)
- Return a set of candidate KG entities, e.g.,
`http://dbpedia.org/resource/Republic_of_the_Congo`
`http://dbpedia.org/resource/Congo_River`
- Typical starting point for CEA and CTA tasks
- DBPedia, Wikidata and Google KG provide look-up services via a REST API.
- Some systems have built their own local index for fuzzy search.

GitHub repositories: <https://github.com/city-knowledge-graphs>

KG look-up Services

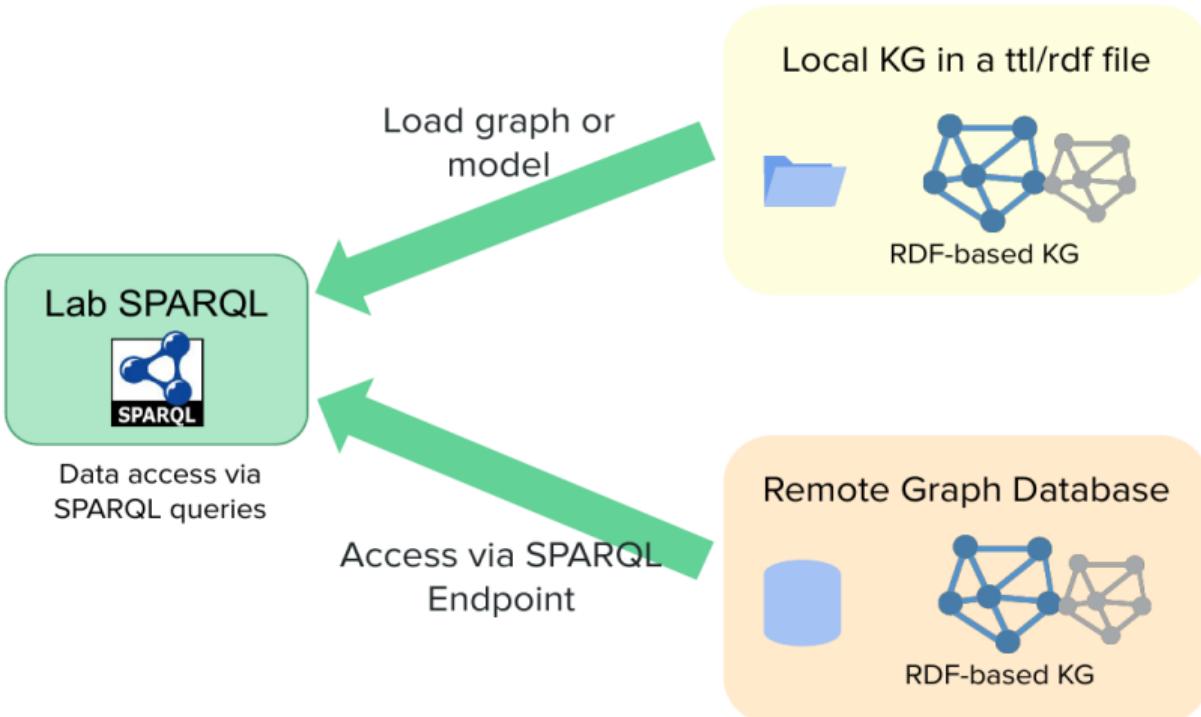


Access to KG SPARQL Endpoint

- Get additional **contextual information**:
 - Additional type information (e.g., dbr:London rdf:type dbo:City)
 - Relationships (e.g., dbr:London dbo:country dbr:United_Kingdom)
 - Labels (e.g., dbr:London rdfs:label "London")
 - Members of a given class.
- Access via **SPARQL queries** (no fuzzy search)
- Typically required for:
 - the **CPA task**
 - **disambiguation** in CTA and CEA tasks

GitHub repositories: <https://github.com/city-knowledge-graphs>

SPARQL: local and remote KG access



Lexical Processing and Similarity

- **Datatype prediction**, e.g., ptype:
<https://github.com/alan-turing-institute/ptype>
- **Spelling corrector**: <https://norvig.com/spell-correct.html>

Lexical Processing and Similarity

- **Datatype prediction**, e.g., ptype:
<https://github.com/alan-turing-institute/ptype>
- **Spelling corrector**: <https://norvig.com/spell-correct.html>
- **Lexical similarity (as in today's lab)**:
 - Levenshtein distance:
levenshtein('Congo', 'Republic of Congo')=12
 - Jaro Winkler:
jaro_winkle('Congo', 'Republic of Congo')=0.0
jaro_winkle('Congo', 'Congo Republic')=0.893
 - I-Sub:
isub('Congo', 'Republic of Congo')=0.727

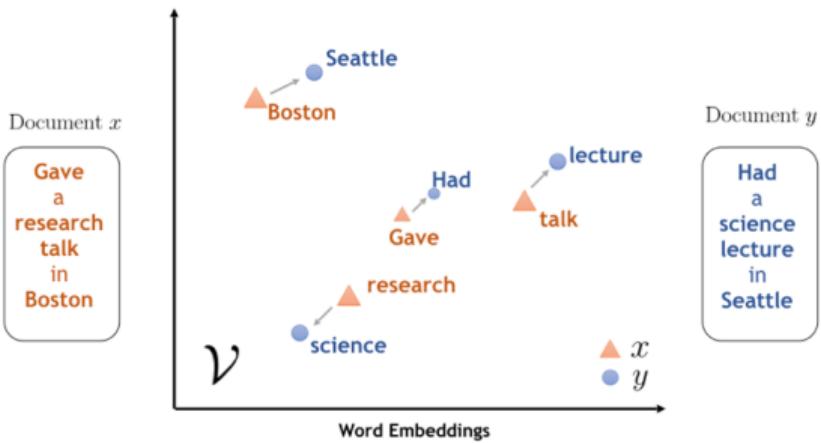
Word and KG Embeddings: Capturing Context

- **Embeddings:** representation in the form of a real-valued vector ([more in week 8](#)).
- Very useful to **capture the meaning/semantics** of a word (or a KG entity).
- Comparison among vectors via **Cosine similarity** or **Euclidean distance** (e.g., between vectors for ‘Congo’ and ‘Republic of Congo’)

Word and KG Embeddings: Capturing Context

- **Embeddings:** representation in the form of a real-valued vector ([more in week 8](#)).
- Very useful to **capture the meaning/semantics** of a word (or a KG entity).
- Comparison among vectors via **Cosine similarity** or **Euclidean distance** (e.g., between vectors for ‘Congo’ and ‘Republic of Congo’)
- **Precomputed word embeddings:**
 - <https://wikipedia2vec.github.io/wikipedia2vec/pretrained/>
 - <https://fasttext.cc/docs/en/pretrained-vectors.html>
 - <https://www.analyticsvidhya.com/blog/2020/03/pretrained-word-embeddings-nlp/>
- **Precomputed KG embeddings:**
 - Wikidata: <http://139.129.163.161/index/toolkits#pretrained-embeddings>

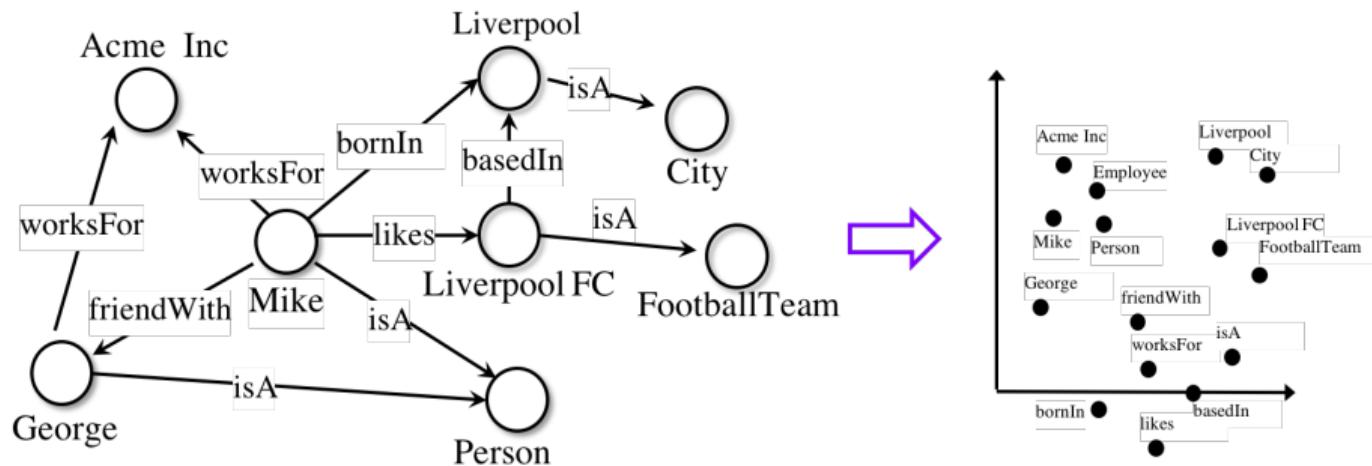
Word Embeddings: Example



Systems like Word2Vec require a corpus of documents as training.

Example from: https://dsgiitr.com/blogs/word_embeddings/

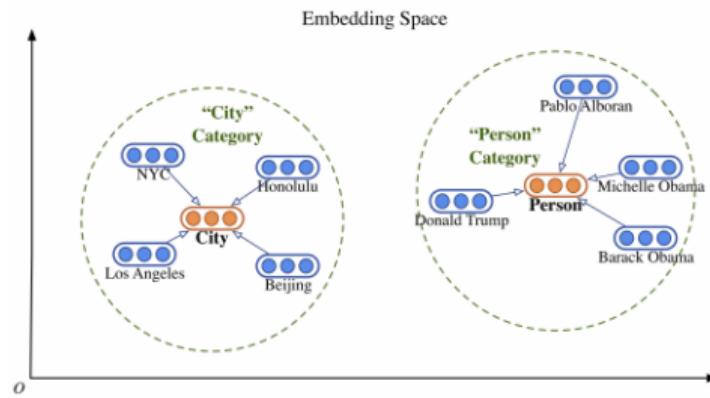
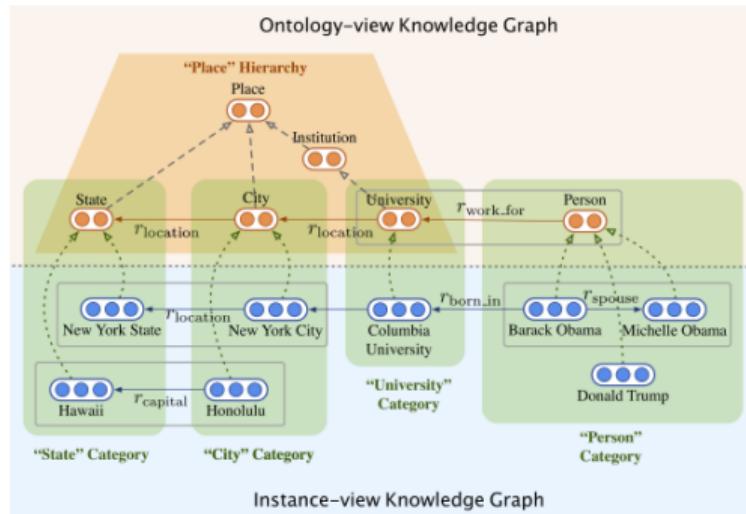
Knowledge Graph Embeddings: Example (i)



KG Embedding Systems exploit the neighbourhood of an entity to calculate its vector.

Example from: <https://docs.ampligraph.org/en/1.0.3/>

Knowledge Graph Embeddings: Example (ii)



KG Embedding Systems exploit the neighbourhood of an entity to calculate its vector.

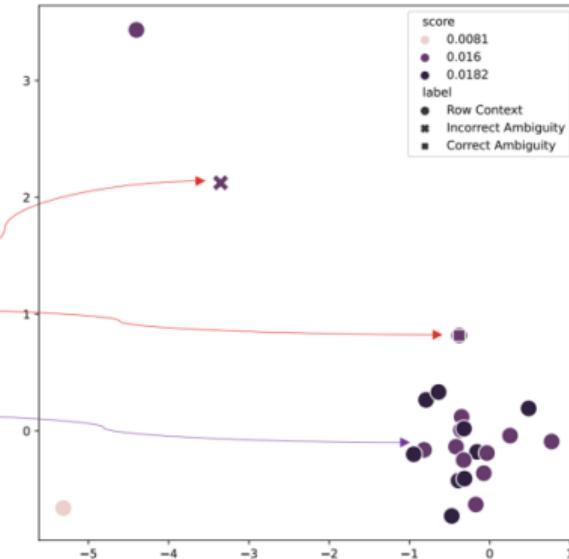
Example from: Universal Representation Learning of Knowledge Bases by Jointly Embedding Instances and Ontological Concepts. KDD 2019.

Knowledge Graph Embeddings: Example (iii)

13:	2002	Enemy Lines : Rebel Dream		
14:	2002	Enemy Lines : Rebel Stand		
15:	2002	Traitor		
16:	2002	Destiny's Way		
17:	2002	Ylesia		e - book

DAGOBAH-SL

```
...  
{'row': 15, 'column': 1,  
'Annotations': [  
    {'id': 'Q21161161', 'score': 0.01600},  
    {'id': 'Q7833036', 'score': 0.01600},  
    {'id': 'Q1536329', 'score': 0.01164}, ..... ]},  
{'row': 16, 'column': 1,  
'Annotations': [  
    {'id': 'Q5265233', 'score': 0.01600},  
    {'id': 'Q60172766', 'score': 0.0102},  
    {'id': 'Q17010392', 'score': 0.0102}, ..... ]},  
...  
-- Ambiguities  
-- Context
```



KG Embeddings can disambiguate Table annotation cases.

Example from: Radar Station: Using KG Embeddings for Semantic Table Interpretation and Entity Disambiguation. ISWC 2022: 498-515

Semantic Understanding of Tabular Data: User Interfaces

OpenRefine

- <https://openrefine.org/>
- Previously known as *Google Refine*.
- **Interface to support** the cleaning and transformation of messy data.
- Includes a **reconciliation service** to link the data with a KG (e.g., Wikidata is the default installation).
- ⚠ In this module we will not use OpenRefine, but perform our own reconciliation programmatically.

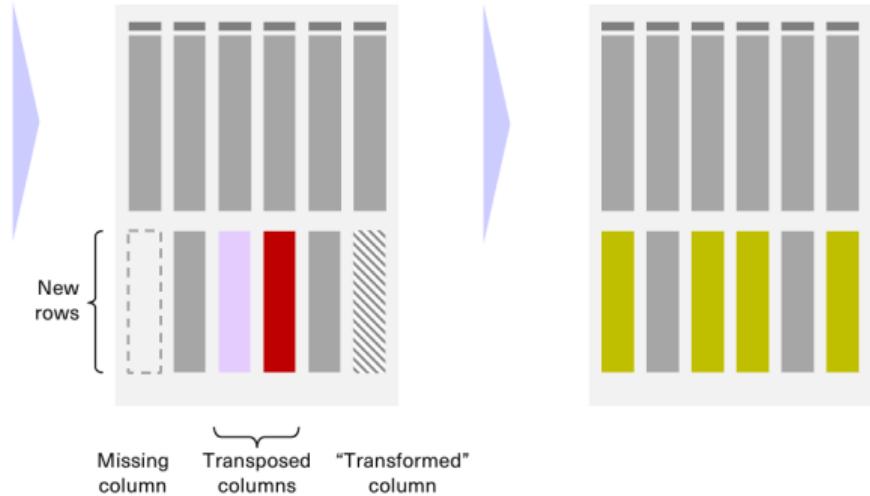
Applications of KGs and Semantics

Examples of Applications of KGs and Semantics

- Data Wrangling (Alan Turing Institute)
- Data Access in Oil & Gas Industry
- Data Access and Prediction in Ecotoxicology
- Semantic Reasoning for Autonomous Vehicles
- FAIR Implementation for Life Science
- Amazon Product Graph

AIDA project: Data Wrangling with DataDiff

- The structure of a dataset may change after an update
- Changes may break the analytical pipeline.
- ✓ Datadiff **identifies and patches** these changes.
- ✗ Limitation: **exhaustive comparison** of columns.
- ✓ Semantic table understanding **may limit the comparison**.



Data Diff: Interpretable, Executable Summaries of Changes in Distributions for Data Wrangling. C. Sutton, T. Hobson, J. Geddes and R. Caruana. In KDD 2018.

Data Access in Oil & Gas Industry

- Data access currently takes **30-70%** of the engineers' time.
- Data cannot be moved from the original sources.
- The EU project Optique advocated for an **Ontology-Based Data Access** (OBDA) process. Requirements:
 - Domain ontology.
 - Mappings to create a virtual KG.

Ontology Based Data Access in Statoil. Journal of Web Semantics, 44, pp. 3-36

<https://openaccess.city.ac.uk/id/eprint/22959/>

Data Access in Oil & Gas Industry: Limitations

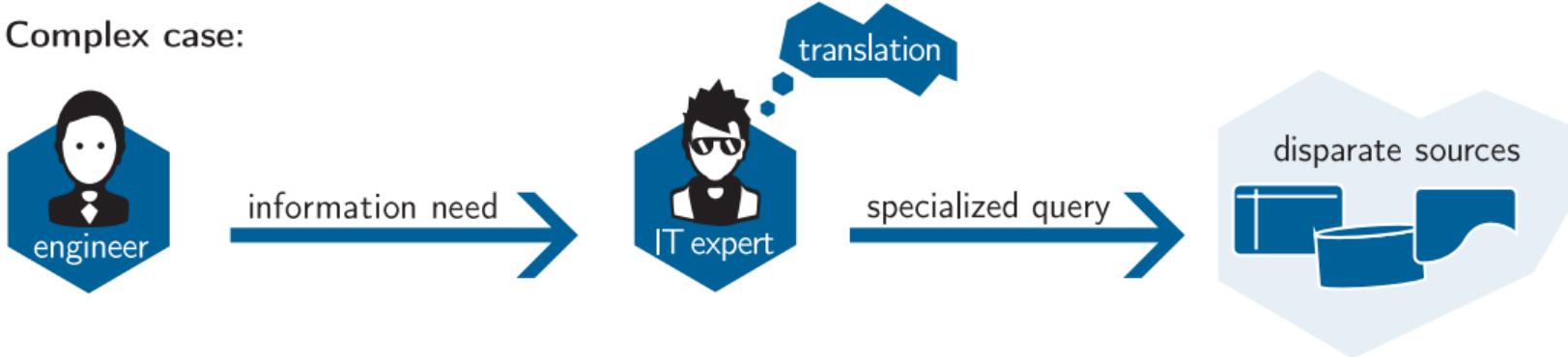
Simple case:



Problem when the information needs fall outside predefined-queries

Data Access in Oil & Gas Industry: Limitations

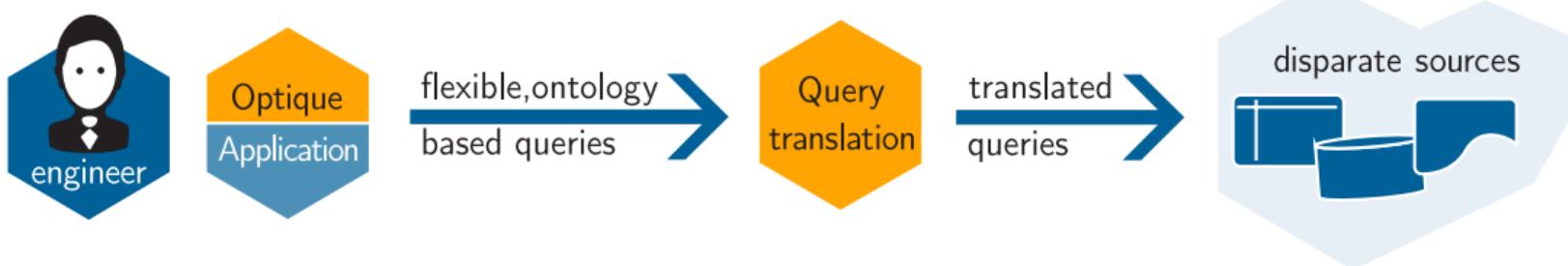
Complex case:



The process may take several days

Data Access in Oil & Gas Industry: Optique Solution

Optique solution

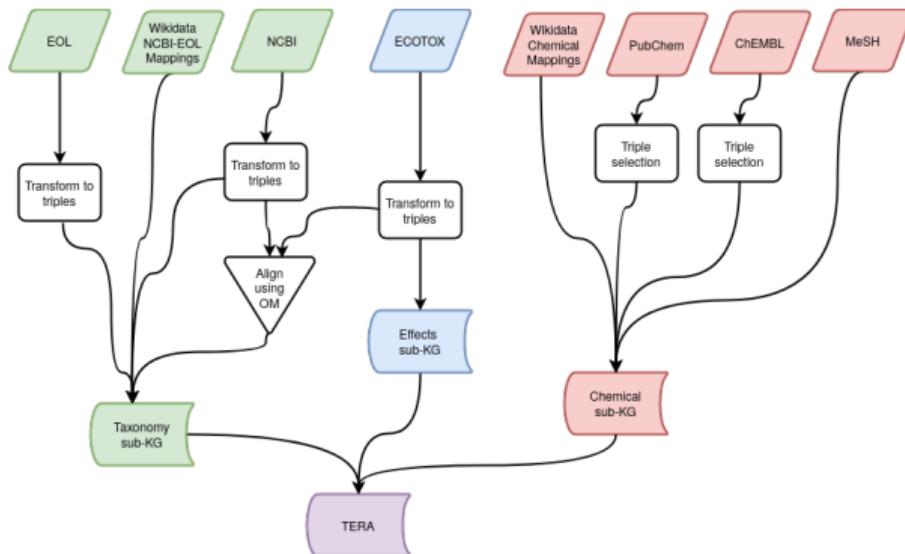


Optique Solution

1. Mediator to create ontology-driven queries (SPARQL).
2. Mediator to translate SPARQL queries into SQL queries.
3. Effort required to create the ontology and maintain the mappings (modular approach).

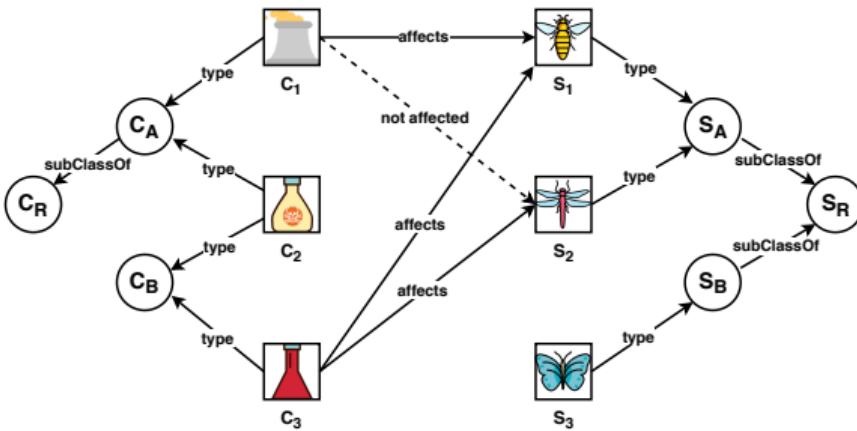
TERA: A KG for Ecotoxicology. Integration and Data Access.

- **Integrates** disparate sources about species, chemicals and effect data.
- Enhances **data access**.



TERA: A KG for Ecotoxicology. Prediction.

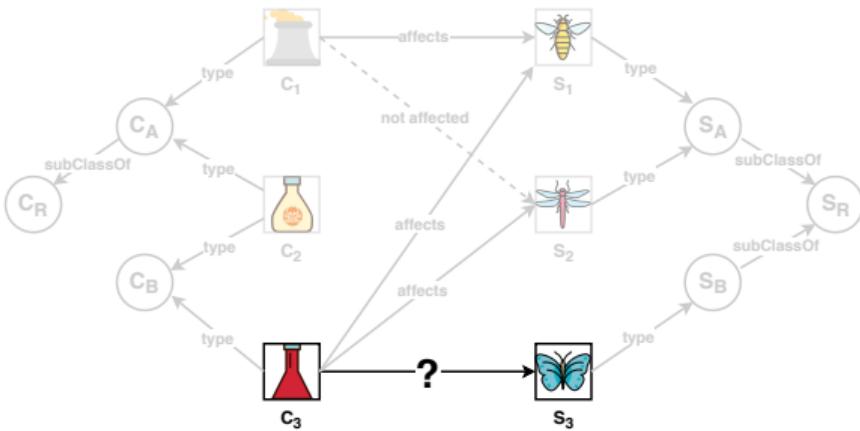
- Drives the **prediction of adverse biological effects** of chemicals via KG embeddings.



Resources and publications: <https://github.com/NIVA-Knowledge-Graph/>

TERA: A KG for Ecotoxicology. Prediction.

- Drives the **prediction of adverse biological effects** of chemicals via KG embeddings.



Resources and publications: <https://github.com/NIVA-Knowledge-Graph/>

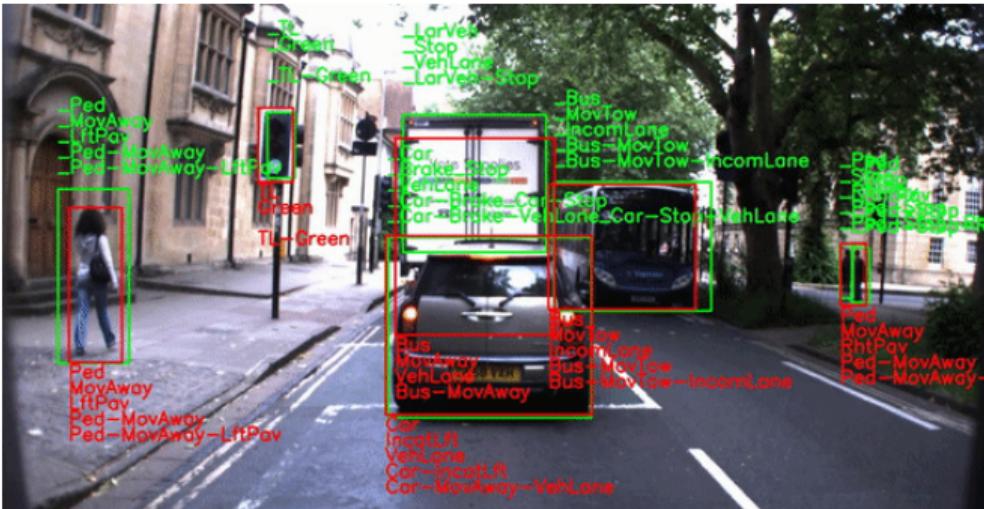
Semantic Reasoning and Ontologies for Autonomous Vehicles



Oxford Semantic Technologies: <https://www.oxfordsemantic.tech/blog/reasonable-vehicles-rule-the-road>

ROAD-R Dataset

Extensions of the ROad event Awareness in Autonomous Driving (ROAD) with Requirements (*i.e.*, background knowledge)



ROAD-R: <https://sites.google.com/view/road-r/>

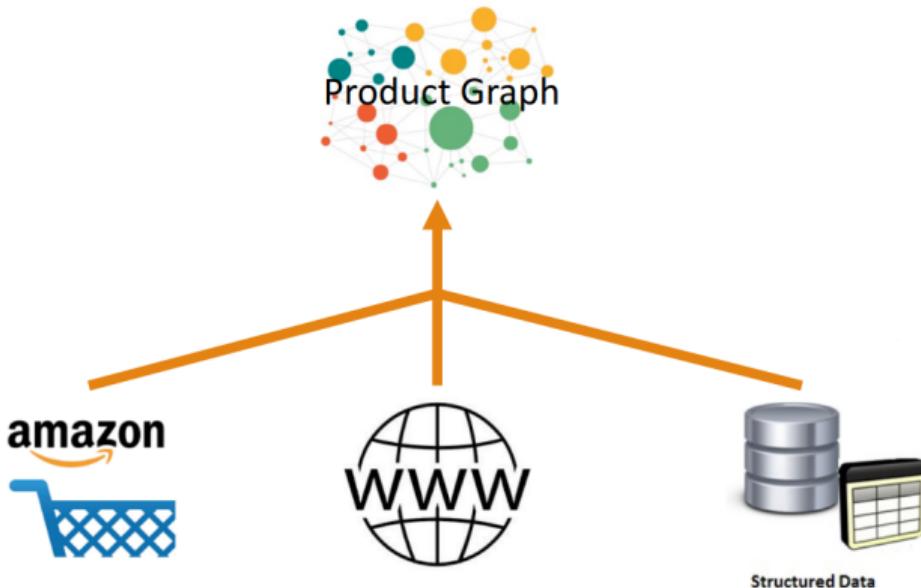
FAIR Implementation for Life Science

Knowledge Graphs play a key role for annotating data.



Pistoia Alliance FAIR Toolkit: <https://fairtoolkit.pistoiaalliance.org/>

Amazon Product (Knowledge) Graph



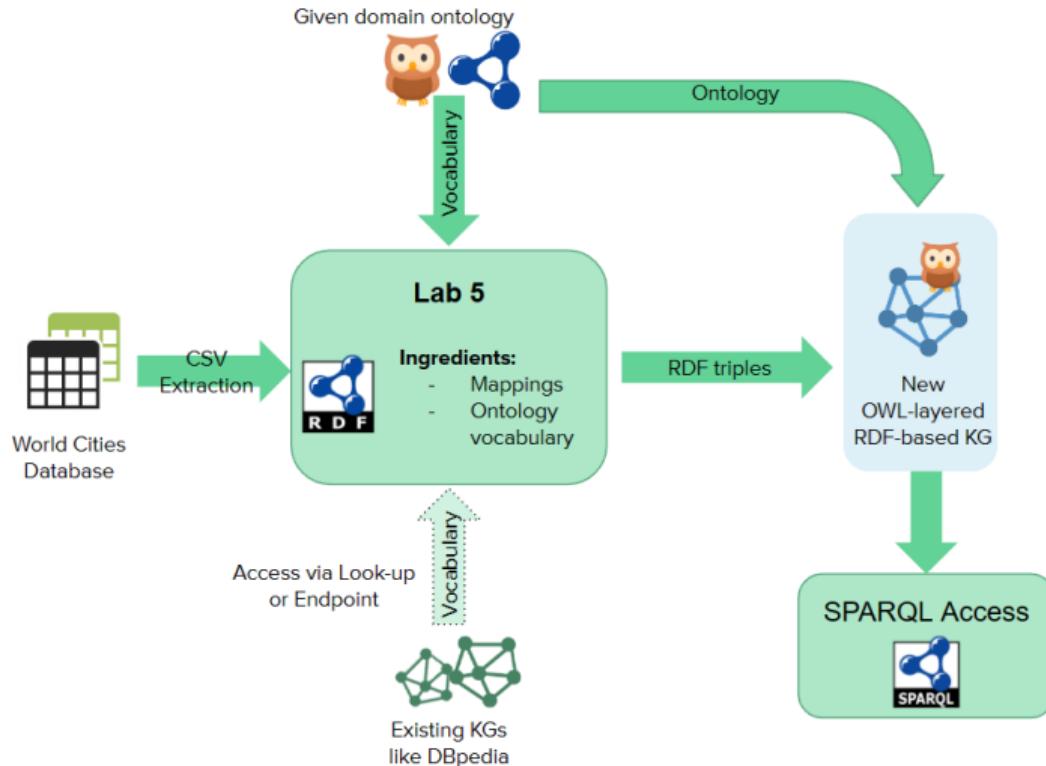
Challenges and Innovations in Building a Product Knowledge Graph: <http://lunadong.com/talks/PG.pdf>
<https://www.amazon.science/blog/building-product-graphs-automatically>

Laboratory: From CSV to a KG

Lab Session 2 Weeks 5, 8 and 10

- **R201 (Franklin Building)**. 48 seats/PCs. 10 min walk.
- **Session 1: Thursday, 11:00-11:50.**
- Session 2: Thursday, 12:00-12:50
 - **C301 on weeks 5, 8 and 10. (*i.e.*, today!)**
- Programming languages: Python and/or Java

Global picture: Lab 5



Lab notes and support Codes

- <https://github.com/city-knowledge-graphs>
- **Solution Lab 2 (Task 2.4)**
- Lookup access.
- SPARQL Endpoint access.
- Lexical similarity
- CSV management