**KGE - Knowledge Graph Engineering**

# The Reuse Problem

**Fausto Giunchiglia**

# Contents

# KG Purpose

- Both data producer and consumer consider their own objective when building KGs.

- Such an objective implicitly includes the user "point of view", the representation that the user uses to model (a portion of) the world, where the desired information lives.

- We define the user objective, **The Purpose** which will lead the entire KGE process.

# KG Purpose — producers / consumers

- KG producers: publish their data in some catalogues/repositories maximizing reuse. They strive for maximum sharability (use of ontologies, most of the time). Purpose with stronger schema requirements, data defined by the applications generating them, cleaned with as minimal as possible cleaning (EML data level requirements)

- KG consumers: use existing data (and schemas) to produce their own (e.g., market application) data. Purpose with stronger data requirements, schema mainly defined by the application, with as minimal as possible tweaks (EML schema level requirements)

NOTE: Always, the schema is the *means*, data is the *goal*.

# Purpose-specific KGs

- Each user who desires to build a KG, will have her own Purpose.

- That is why, most of the KGs produced are **purpose-specific**.

- Even trying to build general purpose KGs, it is impossible to extract the KG produced from a specific (more/less general/specific) context.

- **Context**: a vision of the world we live in, which can change geographically, socially as well as considering its evolution in time.

# Contents

# Reusable Resources

- The user purpose is reflected on the resources used to build the KG.

- Three categories of reusability are defined for such resources, depending on their relevance for the purpose itself: **Core**, **Common** and **Contextual** resources are used to build KGs.

- In order to define the 3 categories, we will follow the example having the following purpose:

*"The user want build a KG able to support the access to the health facilities in Trento (Italy) and all the medical care that they can offer to the citizens."*

# Common resources

- **Common resources**: this category involves resources carrying information common to several contexts (or domain of interest), thus they can be resources not strictly related to the user's Purpose, but essential to support it in the KG.

- Some example of common resources for the Purpose declared above, can be datasets for:
  - Trento location.
  - Trento public transportation.
  - Trento parking areas.
  - ...

# Core resources

- **Core resources**: this category involves resources carry information about the most important aspects considered by the purpose, information without which it would be impossible to build the KG.

- Some example of core resources for the Purpose declared above, can be datasets for:
    - Trento Hospitals.
    - Trento Pharmacies.
    - Drugs.
    - Patients.
    - Doctors.
    - ...

# Contextual resources

- **Contextual resources**: this last category involves resources which carry specific, possibly unique, information related to the user's Purpose. These are the resources whose main goal is to create added value. If core resources are necessary for a meaningful application, contextual resources are the ones which can make the difference with respect to the competitors.

- Some example of contextual resources for the Purpose declared above, can be datasets for:

    - Drug's components.

    - Medical Prescriptions.

    - Hospital Departments.

    - Doctors Specializations.

    - Patient Allergies.

    - ...

# Reusable Resources

- Depending by its category, a resources can be more or less reusable.

- Common resources are **the most reusable**, due to their usage shared among several contexts (so even among different purposes)

- Core resources are **less reusable** (even less findable) respect the common ones, due to their specificity on the contexts/domain of interest.

- Contextual resources are **the least reusable** due to their specificity on the Purpose. Moreover, **they are often created** from scratch because not available from other sources.

# The Purpose of KG producers / consumers

- KG producers: Focus on common and core data (and schemas as a consequence). contextual data possibly to be dropped and not published.

- KG consumers: Focus on core and contextual data (and schemas as a consequence): maximum reuse of common and core data. Contextual data possibly generated on purpose.

# Contents

# Open data Catalogs

- Where are the reusable resources we need to build KGs ?

- Several projects and open data portal already exist which allow to retrieve useful resources.

- Often such resources are accessible through **Catalogs**. They are open portals collecting information about several resources (i.e. datasets, schemas, ontologies, ... ).

- The catalogs doesn't collect the real resources, but instead the **metadata** describing such resources. (Catalogs are supported by backhand repositories)

- More metadata are associated to a resource, more detailed it is on the catalog, thus by consequence, it will be more findable and even **reusable**.

# Linguistic Resources

A linguistic resource is a dataset which provides data about languages (e.g., meanings, relations between words, ...).
There are two types of mono/multi-lingual resources: (i) online dictionaries and (ii) Wordnet like resources. Wordnets much more useful in data integration as they connect meanings of words in a LKG.
Check the licence (lots of options).

## Example

- Global Wordnet Association
- WordNet
- Open Multilingual WordNet
- Datascientia/UKC (forthcoming)

# Linguistic Resource Repositories



Figure: Global WordNet Association[1]

---

[1] http://globalwordnet.org/

# Linguistic Resource Repositories



Figure: WordNet Home[2]

---

# Linguistic Resource Repositories

Open Multilingual Wordnet

This page provides access to open wordnets in a variety of languages, all linked to the Princeton Wordnet of English (PWN). The goal is to make it easy to use wordnets in multiple languages. The individual wordnets have been made by many different projects and vary greatly in size and accuracy. We have (i) extracted and normalized the data, (ii) linked it to Princeton WordNet 3.0 and (iii) put it in one place. The Open Multilingual Wordnet and its components are open: they can be freely used, modified, and shared by anyone for any purpose. There is a fuller list of wordnets at the Global Wordnet Association's Wordnets in the World page.

If you use these wordnets, please cite the original projects who created them (linked in Table 1), if you got value from this aggregation/normalization, please cite Bond and Paik (2012).

You can access the wordnets through the (python) Natural Language Tool-Kit wordnet interface (NLTK).

We have an extended version with automatically extracted data for over a 150 languages from Wiktionary and the Unicode Common Locale Data Repository (Bond and Foster, 2013).

Documentation, News and Updates

**Search**

We have a simple search interface (search the extended wordnet). It uses the SQL database originally developed by the Japanese Wordnet.

### 34 Open Wordnets Merged

| Wordnet | Lang | Synsets | Words | Senses | Core | Licence | Data | Citation |
|---------|------|---------|-------|--------|------|---------|------|----------|
| Albanet | als | 4,675 | 5,988 | 9,599 | 31% | CC BY 3.0 | als.zip (+xml) | cite:als; (.bib) |
| Arabic WordNet (AWN v2) | arb | 9,916 | 17,785 | 37,335 | 47% | CC BY SA 3.0 | arb.zip (+xml) | cite:arb; (.bib) |
| BulTreeBank Wordnet (BTB-WN) | bul | 4,959 | 6,720 | 8,936 | 99% | CC BY 3.0 | bul.zip (+xml) | cite:bul; (.bib) |
| Chinese Open Wordnet | cmn | 42,312 | 61,533 | 79,809 | 100% | wordnet | cmn.zip (+xml) | cite:cmn; (.bib) |

Figure: Open Multilingual WordNet Home[3]

---

[3]http://compling.hss.ntu.edu.sg/omw/

# Linguistic Resource Repositories



Figure: UKC Home[4]

---

# Knowledge Resources

A Knowledge resource is a dataset which consists of a KB encoding information about schemas (etypes and properties).
KBs of high quality are usually called ontologies. We call them teleologies (meaning by this, ontologies with metadata which empower their practical use in knowledge and data integration).

## Example

- LOV/LOV4IoT
- Schema.org
- DBpedia (schema only)
- Datascientia/liveschema (forthcoming)

# Knowledge Resource Repositories



Figure: Linked Open Vocabulary[5]

---

[5]https://lov.linkeddata.es/dataset/lov/

# Knowledge Resource Repositories



Figure: Schema.org[6]

---

[6]http://www.schema.org/

# Knowledge Resource Repositories



Figure: DBpedia Home[7]

---

# Knowledge Resource Repositories



Figure: DataScientia Home[8]

---

# Data Resources

A data resource is a dataset which consists of data in some format (tabular, unstructured, entities and property values).
*Open Data*: data freely available. Check the licence (lots of options).

## Example

- UK Open Data
- National Bureau of Statistics, China
- data.org
- Opendata Trentino (see, among others, Unitn Open Data)
- Geonames
- Open Street Map
- DBPedia
- Data Hub

# Data Resource Repositories



Figure: Open Data UK[9]

---

[9]https://data.gov.uk/

# Data Resource Repositories



Figure: National Bureau of Statistics, China[10]

___

# Data Resource Repositories



Figure: data.org[11]

---

[11] https://www.data.org/

# Data Resource Repositories



Figure: Open Data Trentino[12]

# Data Resource Repositories



Figure: Geonames Home[13]

---

[13]https://www.geonames.org/

# Data Resource Repositories



Figure: Open Street Map Home[14]

---

[14]https://www.openstreetmap.org/

# Data Resource Repositories



Figure: DBpedia Home[15]

---

[15] https://wiki.dbpedia.org/

# Data Resource Repositories



Figure: Data Hub Home[16]

---

[16] https://datahub.io/

# Contents

# Data reuse - Producer & Consumer

- The Purpose changes between producer and consumer users.

- Different kinds of resources are considered if the Purpose aims to produce **new structured resources** (suitable to be reused by KGE processes), or to exploit as much as possible, **already existing KG-based resources**.

# Reuse for the Data Producer



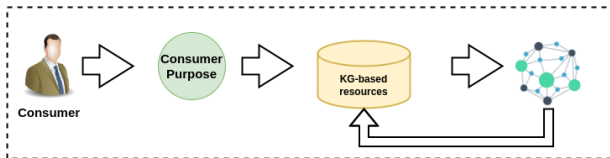- The producer collects resources from catalogs, and/or produces from scratch the resources she needs, with the objective of **produce KG-based version of such resources**, which can be exploited for specific purposes.

- The producer has to deal with the data **heterogeneity** present in the existing catalogs resources:
    - Low quality resources.
    - Noisy resources.
    - Resources not described by metadata.
    - Missing values.
    - Syntactic heterogeneity (see next lectures)

# Reuse for the Data Consumer



- The consumer collects KG-based resources with the objective of **produce KG-based new resources by composition**, supporting her specific purpose.

- The consumer has to deal mostly with the data semantic heterogeneity (see next lectures).

- The KG-based resources already have some useful features like:

  - More quality.
  - Structured resources (Knowledge + Data layer).
  - Mandatory metadata description.
  - Syntactically aligned (see Syntactic heterogeneity lecture)

**Fausto Giunchiglia**

**The Reuse Problem**