

Received 9 June 2023, accepted 28 June 2023, date of publication 4 July 2023, date of current version 10 July 2023.

Digital Object Identifier 10.1109/ACCESS.2023.3292153

## RESEARCH ARTICLE

# Data-Driven Methodology for Knowledge Graph Generation Within the Tourism Domain

ALESSANDRO CHESSA<sup>1</sup>, GIANNI FENU<sup>2</sup>, ENRICO MOTTA<sup>3</sup>, FRANCESCO OSBORNE<sup>3,4</sup>,  
DIEGO REFORGIATO RECUPERO<sup>2</sup>, ANGELO SALATINO<sup>3</sup>, AND LUCA SECCHI<sup>1,2</sup>

<sup>1</sup>Linkalab, 09122 Cagliari, Italy

<sup>2</sup>Department of Mathematics and Computer Science, University of Cagliari, 09124 Cagliari, Italy

<sup>3</sup>Knowledge Media Institute, The Open University, MK7 6AA Milton Keynes, U.K.

<sup>4</sup>Department of Business and Law, University of Milano-Bicocca, 20126 Milan, Italy

Corresponding author: Diego Reforgiato Recupero (diego.reforgiato@unica.it)

**ABSTRACT** The tourism and hospitality sectors have become increasingly important in the last few years and the companies operating in this field are constantly challenged with providing new innovative services. At the same time, (big-) data has become the “new oil” of this century and Knowledge Graphs are emerging as the most natural way to collect, refine, and structure this heterogeneous information. In this paper, we present a methodology for semi-automatic generating a Tourism Knowledge Graph (TKG), which can be used for supporting a variety of intelligent services in this space, and a new ontology for modelling this domain, the Tourism Analytics Ontology (TAO). Our approach processes and integrates data from Booking.com, Airbnb, DBpedia, and GeoNames. Due to its modular structure, it can be easily extended to include new data sources or to apply new enrichment and refinement functions. We report a comprehensive evaluation of the functional, logical, and structural dimensions of TKG and TAO.

**INDEX TERMS** Knowledge graphs, ontology design, tourism ontology, web science, web mining, tourism, hospitality.

## I. INTRODUCTION

We are currently living in the age of big data, and the sheer volume of new data being generated is making the World Wide Web shifting from a web of content to a web of data. This gives all practitioners the opportunity to build more innovative and functional web services. Semantic Web and Linked Data technologies aim to represent the web itself through a large global graph that can be queried using standard protocols and languages [1]. The World Wide Web Consortium (W3C) has developed and promoted different standards, like RDF/S, OWL and SPARQL, that are now widely adopted to create knowledge bases that represent data as knowledge graphs (KGs). A knowledge graph is a graph of data whose nodes represent entities of interest and whose edges represent relations between these entities [2]. We are referring here to RDF knowledge graphs although other approaches, like labeled property graphs (LPG), are possible.

The associate editor coordinating the review of this manuscript and approving it for publication was Senthil Kumar<sup>1</sup>.

A few examples of knowledge graphs publicly available are DBpedia [1], YAGO (Yet Another Great Ontology) [3] or WikiData [4]. Knowledge graphs can store data and metadata using a common structure and are often used in application scenarios that involve extracting and integrating information from multiple, and possibly heterogeneous, sources. Typically the data in the knowledge graph are modelled according to a domain ontology, which gives meaning to the represented information and supports inferring new knowledge. The field of tourism is a natural domain of application of these technologies since stakeholders in this space need to integrate data from several heterogeneous sources in order to generate a multifaceted characterisation of tourist destinations and all relevant actors [5], [6], [7].

A tourist destination can be thought of as the place or area which is central in the decision of a tourist to take the trip<sup>1</sup> and is usually characterised according to two aspects: supply

<sup>1</sup>The World Tourism Organization (UNWTO) defines in its glossary a *destination* as “the place visited that is central to the decision to take the trip”. See <https://www.unwto.org/glossary-tourism-terms>

and demand. The supply side is based on the willingness and ability of producers to create goods and services to take them to market. Understanding the supply side of tourism includes all aspects related to tourism offerings and attractions (e.g., accommodations, events, points of interest, restaurants, and so forth). On the other hand, demand refers to how much (quantity) of a product or service is desired by buyers. Understanding which factors influence the demand side of tourism includes all aspects related to tourists' choices and opinions or their characteristics (e.g., socio-demographic, classification, provenance).

This information is crucial for informing business and marketing decisions as well as supporting a variety of software and services in this space, such as search engines and recommendation systems [8], [9].

The creation of KGs in this domain is a time-consuming and costly process, even with the help of mapping languages such as RML [10], [11], [12]. Indeed, it is still a challenge to automatically generate KGs from multiple semi-structured and textual sources (e.g., descriptions of specific accommodations, reviews, etc.) in order to describe the many facets of this domain, such as the different kinds of accommodations and amenities. Therefore, many KGs in this space are no longer maintained [6], [7] or cannot be easily extended to other tourist destinations [11]. In addition, the relevant ontologies, such as Accommodation Ontology,<sup>2</sup> Schema.org,<sup>3</sup> and Hontology [13] are to some degree incompatible with each other (as discussed in section III-C1) and do not offer a fine-grained representation of some crucial entities (e.g., amenities).

In this paper, we illustrate a general, reproducible, and easily extendable methodology for KG generation and the resulting framework for semi-automatically creating a *Tourism Knowledge Graph (TKG)*, which integrates information from Booking.com, Airbnb.com, DBpedia, and GeoNames. This advanced characterisation of tourism can be used to enable the quantitative analyses of a tourist destination and support several intelligent services. In order to model this data, we developed the Tourism Analytics Ontology (TAO), which offers a more granular characterisation of tourist locations, lodging facilities, and amenities than previous solutions and can be easily reused by similar initiatives.

We showcase our solution by applying it to touristic locations in Sardinia and London, producing over 10M triples describing almost 36K lodging facilities and 898K reviews. The resulting knowledge graph is available via a SPARQL end-point. The TAO ontology is available online.<sup>4</sup> Finally, for the sake of reproducibility, we share the code base for our knowledge graph generation pipeline, for engineering TAO, and the evaluation tests.<sup>5</sup>

To summarise, the contributions of this manuscript are the following:

- a general data-driven methodology for the semi-automatically generation of knowledge graph that we applied to the tourism domain;
- an open-source pipeline for generating a tourism knowledge graph from (semi-) structured and unstructured data;
- the new Tourism Analytics Ontology (TAO);
- an open-source program to produce the Tourism Analytics Ontology (TAO) using code and data;
- an instance of the tourism knowledge graph (TKG) with data relative to two Tourist Destinations (Greater London and Sardinia island in Italy);
- an evaluation assessing functional, logical, and structural dimensions of TAO and TKG.

The remainder of this paper is organised as follows. Section II describes related works about different knowledge graphs within the tourism domain and methodologies for their creation. Section III explains the methodology adopted to guide the knowledge graph creation, detailing the first three iterative phases related to the use cases refinement and ontology design. Section IV describes the other three phases of the adopted methodology related to the creation of the proposed knowledge graph. Section V presents the evaluation, and finally, Section VI ends the paper with conclusions and future directions of work.

## II. RELATED WORK

In this section, we will review the literature on the two main themes concerning this work: i) methodologies for ontology and knowledge graph creation and ii) knowledge graphs within the tourism domains.

### A. ONTOLOGY AND KNOWLEDGE GRAPH CREATION

Creating, maintaining, and further developing knowledge graphs requires the adoption of a number of ontology engineering methodologies (OEMs). Kotis et al. [14] classified such methodologies into three categories: collaborative, non-collaborative, and custom. A collaborative OEM is clearly and systematically defined and involves knowledge engineers, knowledge workers as well as domain experts in all the phases of ontology creation. A non-collaborative OEM does not focus on the collaboration of stakeholders although it still clearly defines phases, tasks, and workflows in a systematic and formal way. A custom OEM does not necessarily define phases, tasks, and workflows in a formal and systematic way; however, it looks for the involvement of communities of practice and the use of tools for the development of ontologies in an agile, decentralized, and most of the time collaborative manner.

There are plenty of works in literature dealing with the creation of knowledge graphs and their methodologies within different domains and constraints [15], [16], [17], [18]. The challenges to be faced depend on such constraints that need to be satisfied by the developers. For example, when

<sup>2</sup><http://ontologies.sti-innsbruck.at/acco/ns.html>

<sup>3</sup><https://schema.org/docs/hotels.html>

<sup>4</sup>See <http://purl.org/tao/ns>

<sup>5</sup>See <https://github.com/linkalab/tkg>

knowledge graphs need to be built starting from a complex database schema, there are difficulties (especially related to its dimension) that must be addressed (i.e., how to efficiently read tables, which columns to consider, how to map linked tables, and so on). In this direction, Sequeda et al. [15] presented a novel and unique pay-as-you-go approach to overcome the difficulties of understanding complex database schemas, providing also a use case from a large company. Tamašauskaitė and Groth [16] presented a systematic review of the process for knowledge graph creation. The review methodology aimed at collecting the various steps describing such a process and these include: identification of the data, construction of the knowledge graph ontology, extraction of knowledge, analysis of the extracted knowledge, creation of the knowledge graph and maintenance. The last step is the one that tends to provide periodical updates and edits to the current knowledge graph. In this review, the authors provide suggestions, best practices, and tools supporting the creation and maintenance of knowledge graphs.

In this paper, we present a data-driven methodology that encompasses the semi-automatic generation of the knowledge graph exploiting several off-the-shelf tools and the engineering of a supporting domain ontology using a collaborative OEM (as defined in [14]). The methodology is applied to generate a Knowledge Graph within the tourism domain.

## B. KNOWLEDGE GRAPHS WITHIN THE TOURISM DOMAIN

In previous years, various attempts have been made to build knowledge bases in several domains, including tourism, using information extracted from websites and social media.

For instance, the 3city platform [5] was built during Expo Milano 2015 to create comprehensive knowledge bases, containing descriptions of events and activities, places and sights, transportation facilities, and social activities collected from numerous, local and global data providers, including hyper-local sources. Using the sample platform, in 2016-2017 new knowledge bases have been created for the cities of London, Madeira, and Singapore, as well as for the entire French Cote d'Azur area. The project now seems no longer maintained and no source code was released to recreate the infrastructure. Although a SPARQL endpoint remains active it only allows the user to export data only in HTML and not as RDF.

The Tourpedia platform which was meant to be the DBpedia of tourism was developed within the OpeNER Project [6]. OpeNER (Open Polarity Enhanced Name Entity Recognition) was a project funded under the 7th Framework Program of the European Commission whose main objective was to implement a pipeline to process natural language. The project is no longer maintained although anyone can run the proposed pipeline to view categories, places information, and create and manage events and tour plans for users. Also, on the main website, it is still possible to run the web demo application, showing the sentiment about places through an interactive map. Some datasets are still available

for download although other tools, including the SPARQL endpoint, are no longer working.

DBtravel [7] is a tourism-oriented knowledge graph generated from the collaborative travel site Wikitravel that takes advantage of the recommended guidelines for contributors provided by Wikitravel and extracts the named entities available in Wikitravel Spanish version<sup>6</sup> by using an NLP pipeline. As for the previous two projects, the knowledge graph and the source code used to produce it are no longer maintained nor available online.

Other projects demonstrate that semantic technologies and knowledge graphs can be successfully applied to tourism when information is extracted from curated proprietary data sources. In the case of *La Rioja Turismo* Knowledge Graph, Alonso-Maturana et al. [11] retrieve and integrate information referring to attractions, accommodation, tourism routes, activities, events, restaurants, and wineries from heterogeneous and diverse management systems. This approach is focused on the La Rioja Turismo ecosystem but cannot be easily extended to other tourist destinations.

In the case of the Tyrolean Tourism Knowledge Graph [19], data based on Schema.org annotations are collected from destination management organisations (DMOs) and their IT service providers. In this case, the knowledge graph creation is based on the availability of coherent Schema.org annotations in the source websites, which was possible thanks to the cooperation of Tyrolean DMOs. Once again, this scenario is not always applicable because it requires a central organisation to coordinate the different stakeholders.

Another proposed approach was to collect, enrich, and publish Linked Open Data for the Municipality of Catania, a city in Southern Italy, in the context of the project PRISMA, "Platform Interoperable cloud for Smart-Government"<sup>7</sup> [12], [20], [21], [22]. In this case, Consoli and his colleagues presented the collected city data, described the process and issues to create a semantic data model for emergency vehicle routing and geo-linked data, and discussed a developed prototype. In particular, they described the employed procedures, ontology design patterns, and tools used for ensuring semantic interoperability during the transformation process. Although the project is flexible and can be generalized, the authors did not maintain the resulting knowledge graph.

Other state-of-the-art solutions include the generation of a knowledge graph of tourism in the Chinese language [10], [23]. The authors constructed such knowledge graphs by extracting knowledge from the existing encyclopedia knowledge graph and unstructured web pages in the Chinese language. Besides the fact that this knowledge graph is focused on the Chinese language, the authors focused on semi-structured knowledge extraction and deep learning algorithms to extract high-level entities and relations from

<sup>6</sup><https://wikitravel.org/es/Portada>

<sup>7</sup><http://www.ponsmartcities-prisma.it/>

unstructured travel notes. The project is no longer maintained and did not provide a SPARQL endpoint.

It is still a big challenge to automatically generate a knowledge graph about tourism that integrates the most important data sources in this field and can be easily extended to other touristic locations. We also lack a single ontology<sup>8</sup> that would offer a fine-grained description of touristic lodging (e.g., Hotel), accommodations (e.g., family room), amenities (e.g., swimming pool), locations (e.g., amusement park), and destinations (e.g., London). The work presented in this paper proposes to address this gap by introducing the Tourism Analytics Ontology (TAO), which offers a granular characterisation of accommodations, tourist locations and destinations, and a general, reproducible, and easily extendable pipeline to integrate relevant data sources and generate a knowledge graph for the tourism domain. Last but not least, we implemented a SPARQL endpoint and plan to periodically update our knowledge graph by using the proposed pipeline. Differently from the approaches discussed above, our proposal can be easily reused and extended to different tourist destinations since we release the full source code, allowing other users to generate new KGs from several data sources that offer worldwide coverage.

### III. METHODOLOGY PHASES FOR KNOWLEDGE GRAPH DESIGN

Our approach for KG construction is aligned with the general methodology analysed in [16] and is organised into six macro phases that can be iteratively repeated to refine the resulting KG. Specifically, the first three phases are the core of a **data-driven design process** that leverages the knowledge embedded in the data sources for guiding the use case refinement and ontology engineering. The last three phases drive the actual implementation of the knowledge graph, its publishing and validation. Figure 1 describes the different phases.

The *first phase* is focused on the definition of the use cases that the knowledge graph should support, that is to say, what are the desired outcomes a user or an application should be able to produce from it. Because our process is driven by what we can find in the data, this is a preliminary definition that is subject to further refinements and that should be revised multiple times until all use cases are positively supported by the KG.

The *second phase* is about understanding how the data at our disposal can support the use cases, but it is also about extracting knowledge from the data to support the ontology definition. On the one hand, the data is used to adapt the use cases to the actual information we have access to, thus extending the scope for some use cases or reducing it for others. For example, if we do not find in the data any information about the total number of rooms for a hotel, we cannot support any use case about the available accommodation capacity for a tourist destination

unless we find new data sources. On the other hand, the data is analysed to guide the ontology design. As an example, the accommodations offered on Airbnb have specific types, like shared rooms, which are peculiar to a sharing economy approach. They may also include amenities we seldom find in other forms of hospitality like hotel rooms. This information incorporates knowledge about the hospitality services for tourism that we can use in the process of ontology design and engineering together with the building of the knowledge graph itself.

The *third phase* focuses on the creation of an ontology to model lodging, tourist destinations, and locations that support all the use cases defined in the first phase and incorporate the domain knowledge distilled in the second phase.

The *fourth phase* is about transforming the data extracted from the data sources in order to prepare it to be used for triple creation in the following phase. During this process, various data wrangling techniques are applied to semi-structured data, whereas natural language processing is applied to unstructured texts (e.g., language detection, named entity extraction, and entity linking).

The *fifth phase* is concerned with triple creation using the data prepared in the previous phase. The triple creation is performed using RDF Mapping Language (RML) in order to include in the knowledge graph also the transformation process metadata.

Finally, the *sixth phase* focuses on the publication of the knowledge graph in a triple store and its validation with respect to the use cases.

The proposed methodology is general and it can be applied whenever it is deemed necessary to design a KG and its supporting ontology with a bottom-up approach. It is well suited to address the need to model a KG to support applications that are based on existing data sources that pose practical constraints to the design and implementation process.

In the following subsections, we describe in detail the first three phases related to use case refinement and ontology creation. In Section IV we will then describe the final three phases, related to the creation of the knowledge graph.

#### A. DEFINE THE USE CASES

We start with a first general definition of some use cases that we want to cover when building the KG, also considering what data sources could be used to support them. We should also define which kind of applications we would need to implement on top of the KG to support the use cases. This analysis can give us a more general scenario of how the KG would be used. This, in turn, is useful to understand to what extent the data sources can support the scenario and guide the design process on how the KG should be structured. In fact, this phase is intertwined with the second phase (i.e., *Find and study information sources*), discussed in Section III-B, because we need to consider the information we can extract from the web to support the selected use cases. It is also related to the third phase (i.e., *Define the*

<sup>8</sup>We analyse other ontologies in Section III-C1.





**FIGURE 1.** Tourism Knowledge Graph creation phases.

ontology) in Section III-C, because we can have different design approaches regarding the KG depending on what kind of methods and applications it should support (e.g., whether or not we want to apply reasoning techniques on the KG). In order to generate a KG that can be used to support the analysis of tourist destinations with respect to the supply and demand side, we have identified, together with the domain experts and stakeholders, the following use cases:

- (UC1) Support the identification of the topics of interest discussed by tourists in their reviews;
- (UC2) Support the identification of the topics of interest presented in the descriptions of lodging facilities<sup>9</sup> and accommodation<sup>10</sup> offers;
- (UC3) Support the recognition and linking of tourism entities in the KG for different applications revolving in the domain of social media, news, and blogs;
- (UC4) Support sentiment analysis [24], [25] applications about tourists toward lodging facilities and destinations;
- (UC5) Support the classification of tourist destinations on the basis of what they offer and on the basis of tourist opinions.

We also identified a number of applications that can leverage the KG to produce better results (see [26] for a

<sup>9</sup>Lodging facilities mean any hotel, motel, motor inn, lodge, and inn or other quarters that provide temporary sleeping facilities open to the public. See <https://www.lawinsider.com/dictionary/lodging-facilities>)

<sup>10</sup>An accommodation is a place that can accommodate human beings, e.g., a hotel room, a camping pitch, or a meeting room. An accommodation is always part of a lodging facility (e.g., a hotel room is part of a Hotel).

comprehensive overview of applications based on knowledge graphs). In turn, each one of the following applications can be used to better support one or more use cases:

- 1) **automatic reasoning**<sup>11</sup> and **graph learning**<sup>12</sup> on the KG allows for the entailment of new triples thus enriching the explicit knowledge other applications can work on; for this reason, it is indirectly related to all use cases;
- 2) **named entity recognition (NER) and entity linking (EL)** of tourist locations and lodging facilities using the KG have an immediate positive impact on use cases 3 and 5.
- 3) **relation extraction (RE) in a closed setting** for the tourism industry can be used to support a better understanding of the relations between users and touristic entities thus improving use cases 4 and 5.
- 4) **tourism-related Topic Modelling** (cluster words/phrases frequently co-occurring together in the tourism context) for texts and documents are written in natural language can be used to support use cases 1 and 2.
- 5) **tourism-related Topic Labelling** (for clusters of words identified as abstract topics, extract a single term or phrase that best characterises the topic) can also be used to support use cases 1 and 2.
- 6) **Text Classification** of documents concerning tourism topics can support use cases 1, 2, and 5.

<sup>11</sup>Leveraging Description Logic and OWL.

<sup>12</sup>Using Graph Neural Networks or similar techniques.

- 7) **Semantic Annotation** of documents about tourism with entities, classes, and topics based on the KG can be used to support all the use cases by improving user interfaces and user interactions with the textual data.

It is important to note that, the actual feasibility of a use case can be confirmed only when the knowledge graph is built and one or more of the supporting applications are implemented. This validation phase is out of the scope of the present work, which focuses on the design and construction of the knowledge graph.

## B. FIND AND STUDY INFORMATION SOURCES

To support the use cases described in Section III-A we need to identify a minimum set of information sources we need throughout the construction of a *core* version of the Tourist Knowledge Graph. After this core Knowledge Graph is created, new information sources could be added by applying the same process described in this work. This is because knowledge graphs have a flexible schema which makes them easily extendable.

Observing the use cases, we can see that we need information sources about:

- lodging facilities and the accommodation they offer;
- user reviews and opinions;
- tourist locations (i.e., points of interest for a tourist such as a train station or a beach);
- tourist destinations such as London or the Costa Smeralda (i.e., the place visited that is central to the decision to take the trip);

The first set of information sources adequately covering the listed items consists of:

- Booking.com, a digital travel company specialised in hotels, B&Bs, and other types of hospitality; from its website we can collect information about accommodations and related offers but also users' opinions expressed as reviews.
- Airbnb, an American company that connects hosts, offering their accommodation spaces (e.g., apartments, rooms, etc.), and travelers, looking for a place to stay; it adopts a peer-to-peer model that originates from sharing economy and represents a new emerging reality in the tourism and accommodation market; its website is a source of information about accommodations and related offers but also users' opinions expressed as reviews.
- DBpedia, an open knowledge graph built with structured content extracted from the information created in various Wikimedia projects (e.g., Wikipedia). Specifically, we link entities in TKG to the DBpedia entities of selected classes (e.g., `DBpedia:Places` or `DBpedia:Food`).
- GeoNames, a geographical database exposed through APIs and as RDFs documents. We connect entities in TKG with GeoNames entities representing places.

It is worth noticing that, although there are many other websites and applications for tourism and hospitality,

Booking.com and Airbnb are market leaders and together cover both the traditional accommodation industry and the emerging sharing economy. A similar consideration could be made for DBpedia and GeoNames when we consider places (DBpedia and GeoNames) or general topics related to tourism (DBpedia).

For the present work, we build upon the results of an industrial project about Tourism 4.0 called *Data Lake Turismo* developed by Linkalab s.r.l.,<sup>13</sup> which was the evolution of a previous research project promoted by the Digital Innovation Hub of Sardinia and Fondazione Banco di Sardegna. The project aimed at creating a digital platform for tourism data analysis. One of the main components of this platform was a data lake for collecting, transforming, and analysing data in this sector. However, the project lacked a semantic layer that could support and enhance the data analysis, which is the starting point and motivation of the present work. Through this infrastructure, we have access to data assets related to lodging facilities, user reviews, and opinions; and we enrich them with DBpedia and Geonames. The data source selection influences both the use case and the ontology definition phases. Although it could be possible to add new data sources to the mix from the beginning, it has a cost and should be postponed wherever possible, because our objective is to complete the construction of a core version of the knowledge graph before expanding its coverage. On the other hand, we should always select data sources that incorporate a rich and well-established model of the business sector (tourism in our case) in the data itself. This is important to support the ontology design with a data-driven analysis process.

## 1) SOURCE DATA EXPLORATION

The first step of this phase is to understand what kind of data we can use. We should examine the documentation but we also need to perform an exploratory data analysis on the files and tables accessible in the source data lake in order to have a complete grasp of its contents. This analysis is focused on the following resources available in the data lake:

- data about hospitality:
  - information related to lodging facilities (e.g., hotels, b&b's, resorts) and their characteristics (e.g., name, address, type, hospitality features);
  - information related to accommodations offered by a lodging facility (e.g., hotel room, b&b room, apartment).
  - rent offers for accommodation (e.g., price, number of people, etc.).
- data about user reviews (e.g., user, date, rating, text).

Data is extracted from the data lake in tables with nested structures and needs to be "flattened" to be used by the downstream tasks. This is due to the way the data lake stores information in a redundant and not normalised way.

The result of the exploratory analysis has shown:

<sup>13</sup>Linkalab s.r.l. is an Italian company specialised in data science and data engineering. Home page <https://www.linkalab.it/>

- how data is organised in fields and sub-structures;
- that structured and unstructured data (i.e., texts) is available;
- that texts can be in many different languages and it is not always specified in which one;
- that structured data fields can contain numbers, Boolean values, time/date values, or categorical values;
- that data is not always typed and can be represented internally as strings;
- that categorical data is not related to a lookup table or taxonomy;
- that in some cases there are no unique IDs that can be used to identify a resource.

This analysis led us to define some fundamental data pre-processing steps to be executed before building the knowledge graph and the ontology:

- **data preparation:** in this step, we extracted the data from the source data lake via SQL queries; next, we stored it on a local file system to be prepared (cleaned, flattened, combined) so that it can be used for downstream tasks.
- **data enrichment:** in this step, we augmented the data using various techniques; specifically, we applied NLP techniques to identify the language of the text (e.g., English, Italian, French, and so on), because downstream tasks depend on it to work properly.

We also found that the data lake source should be integrated with data about attractions and points of interest from other sources. To support this need we identified DBpedia and GeoNames as the most appropriate data sources for the following reasons: i) both sources are stable and constantly maintained, with a vast supporting community; ii) both sources cover the identified destinations (and many others) in depth; iii) both sources are exposed as linked open data and APIs.

### C. DEFINE THE DOMAIN ONTOLOGY

To support the identified use cases and the related applications we want to generate a KG that includes all the relevant entities and their relations. We thus need to define a domain ontology that can model them. To guide the ontology design, we defined a set of functional and non-functional requirements in collaboration with domain experts from Linkalab. We also expressed the same requirements in a more operational form using competency questions, i.e., queries expressed in natural language [27], [28]. In Appendix A, we describe in detail the resulting competency questions as well as both functional and non-functional requirements.

Competency Questions (CQ) are useful as they: i) can be easily understood by non-technical people; ii) can guide the ontology engineering process working as a practical reference of what should be implemented; iii) can be easily tested during the validation process.

We adopted a data-driven design process and followed two complementary approaches when defining the competency questions: i) top-down, by developing new questions with

a domain expert and then checking whether they could be answered with our data; and ii) bottom-up, by deriving them from the information available in the source data. Because CQs express all functional requirements in other terms, at the end of this process we could verify that the ontology would successfully model the data in the knowledge graph, which in turn would satisfy the use cases and support the related applications. At the end of this process, we identified the main aspects to model within our domain ontology: i) lodging facilities (buildings), ii) accommodations within lodging facilities, iii) amenities offered to tourists, iv) tourist destinations and locations, and v) user reviews. Next, we analysed several state-of-the-art ontologies covering the tourism domain (detailed in Section III-C1), but none of them fully satisfy our requirements. Therefore, we designed and implemented a new ontology, the Tourism Analytics Ontology (TAO), leveraging existing ontologies (e.g., Schema.org, Hontology).

We devote the following subsections to describing: i) the ontologies which we used as a starting point; and ii) the final version of TAO and our design choices.

#### 1) REUSE OF EXISTING ONTOLOGIES

We analysed several tourism ontologies to assess if they could be reused to support our use cases.

We identified three main families of ontologies:

- 1) ontologies based on Open-Travel or other heavyweight industrial standards, typically focused on information exchange among tourism organisations (e.g., the Harmonise Ontology [29]).
- 2) ontologies produced by researchers to support specific tasks, such as question answering (e.g., QALL-ME Ontology [30]) and information retrieval (e.g., GET-ESS [31]) as well as ontologies that combine or build on them (e.g., cDOTT [32], Hontology [13]).
- 3) ontologies based on Schema.org [33] and GoodRelations [34], such as the Accommodation Ontology.

Based on the functional and non-functional requirements, we then selected three of them: (i) Accommodation Ontology, (ii) the Schema.org markup for hotels, and (iii) Hontology. The latter is currently not available as OWL serialisation at any specific URI and does not seem to be maintained anymore. TAO also reuses other two ontologies: (iv) GeoNames,<sup>14</sup> which is used to specify the geographic locations, and (v) the DBpedia ontology,<sup>15</sup> which is used for further characterising locations and food types (e.g., pizza, sushi). Next, we will describe the selected ontologies and vocabularies and how they have been reused in TAO.

**Accommodation Ontology** (prefix `acco:`) is an extension of GoodRelations (prefix `gr:`) focused on describing accommodation offers from an e-commerce perspective. It provides additional vocabulary elements for describing hotel rooms, hotels, camping sites, and other forms of

<sup>14</sup><https://www.geonames.org/ontology/documentation.html>

<sup>15</sup><https://www.dbpedia.org/resources/ontology/>

accommodations as well as their features. However, it does not make a distinction between the lodging facility (e.g., a hotel as a whole), and the individual accommodations on a lease (e.g., the hotel rooms), because all lodging facility types and accommodation types are sub-classes of the same class (`acco:Accommodation`). The Accommodation Ontology does not define specific types of amenities (called accommodation features) but “provides a consolidated conceptual model for encoding proprietary feature information”. So instead of defining classes for room and hotel features, the ontology provides the generic class `acco:AccommodationFeature` that can hold feature information in varying degrees of formality. A leasing offer is modelled using the `GoodRelations` relation `gr:Offering` specifying that the offering is a `gr:LeaseOut` using the property `gr:hasBusinessFunction`. Unfortunately, the Accommodation ontology does not cover several concepts that are required for our use cases, including 1) tourist destinations (e.g., London), 2) tourist locations (e.g., beach, church, subway station), 3) tourist reviews.

**Schema.org markup for hotels** (prefix `schema:`), incorporates and extends many Accommodation Ontology [35] concepts. Schema.org models hospitality according to three main classes:

- 1) A **lodging business**, (e.g., a hotel, hostel, resort, or a camping site): essentially it represents both the lodging facility, which is the place that houses the actual units of the establishment (e.g., hotel rooms) and the business organisation governing it. The lodging business can encompass multiple buildings but is in most cases a coherent place.
- 2) An **accommodation**, i.e., the relevant units of the establishment (e.g., hotel rooms, suites, apartments, meeting rooms, camping pitches, etc.). These are the actual objects that are offered for rental.
- 3) An **offer** to let a hotel room, or other forms of accommodations, for a particular price and a given type of usage (e.g., occupancy), typically further constrained by booking requirements and other terms and conditions.

In this case, we have a clear distinction between lodging business and accommodation because we have two distinct classes: `schema:Accommodation` and `schema:LodgingBusiness`. Unfortunately, Schema.org is not intended to be used as an OWL ontology because its data model is very generic and derived from RDF Schema.<sup>16</sup> The main purpose of Schema.org is to enable sharing of structured data on the Internet whereas OWL is based on formal semantics that enables reasoning on the knowledge graph. In addition, the `schema:LodgingBusiness` class cannot be used in conjunction with `GoodRelations` ontology without introducing logical contradictions. Specifically, Schema.org defines `schema:LodgingBusiness` as a subclass of `schema:LocalBusiness` which is a subclass of both `schema:Organisation` and

`schema:Place`. On the other hand, `GoodRelations` states that `schema:Organisation` and `schema:Place` are disjoint. We reused Schema.org in TAO by importing and extending a few classes and properties, including `schema:PostalAddress`, `schema:UserReview`, `schema:address`, `schema:subjectOf`. We also selected appropriate Schema.org types that describe places to enrich TAO tourism location classes using `rdfs:seeAlso` to establish a mapping with them.<sup>17</sup>

**Hontology** (prefix `ho:`) is a multilingual ontology for the accommodation sector (H stands for hotel, hostel, and hostel). It is a freely available domain-specific ontology in four languages: English, Portuguese, Spanish, and French [13], [36]. It was partially aligned with QALL-ME and Schema.org and described several useful concepts in this domain such as Facilities (a.k.a. amenities), Services, Staff, and Points Of Interest. The ontology is not published as linked data but can be downloaded and used in a local environment. Its latest version dates back to 2012 and therefore it is not aligned with the most recent extensions of Schema.org. In addition, since it is not based on `GoodRelations`, it does not fulfill our non-functional requirements. We re-implemented within TAO some of its classes describing location amenities, such as `ho:Balance`, `ho:AirConditioning`, `ho:Ballroom`, and `ho:BeautySalon`.

**DBpedia** Ontology (prefix `dbpedia:`) is a shallow, cross-domain ontology, which has been manually created based on the most commonly used infoboxes within Wikipedia.<sup>18</sup> The ontology currently covers 685 classes which form a subsumption hierarchy and are described by 2,795 different properties. We used some of the classes from this ontology to enrich TAO tourist location types (subclasses of `tao:TouristLocation`) also mapped to GeoNames geographic features.

**GeoNames** Ontology (prefix `gn:`) provides elements of description for geographical features, in particular those defined in the `geonames.org` database. It has three key ontology classes: `Feature` (a set of all geospatial instances in GeoNames like cities and countries), `Class` (a set of all feature schemes defined in GeoNames), and `Code` (a set of abbreviation feature codes in different feature schemes). GeoNames `Feature` is used for describing concrete geospatial entities (UK, Washington, Colosseum, etc.), whereas GeoNames `Class` and `Code` are used for representing meta-information about features. All feature instances are uniquely identified by URI in GeoNames.

We used GeoNames `gn:Feature` class to model classes that are also places (e.g., lodging facilities, tourist locations) and to express their geographic relations using `gn:parentFeature`. We also used GeoNames to enrich TAO tourist location types with specific codes, for example, `tao:Park` was associated to the `gn:L.PRK` code.

<sup>17</sup>In this respect we can consider TAO ontology an external extension of Schema.org as described in the page <https://schema.org/docs/extension.html>

<sup>18</sup>As defined in the DBpedia ontology page <http://web.archive.org/web/20210416134559/http://wikidata.dbpedia.org/services-resources/ontology>

<sup>16</sup>See <https://schema.org/docs/datamodel.html>



## 2) THE TOURISM ANALYTICS ONTOLOGY

In this section, we describe the new Tourism Analytics Ontology (TAO) and discuss our design choices. We aimed at developing an ontology i) for which all the requirements listed in Appendix are fulfilled, ii) that would be able to integrate all relevant information from the data sources, and iii) that would be fully compatible with the Accommodation Ontology, GoodRelations, and Schema.org. Specifically, the Accommodation Ontology is explicitly imported using `owl:imports`, GoodRelations is imported indirectly through Accommodation Ontology and Schema.org is partially included by reusing specific classes and properties or making explicit mappings to it.

The new ontology has the following characteristics:

- 1) introduces the `LodgingFacility` class which represents any hotel, motel, inn, or other quarters that provide temporary sleeping facilities open to the public<sup>19</sup>;
- 2) distinguishes between lodging facilities and specific accommodations within lodging facilities;
- 3) includes an extended hierarchy<sup>20</sup> of lodging facilities types (e.g., hotel, house, resort);
- 4) includes an extended hierarchy of the amenities (e.g., oven, parking garage, baby monitor) offered by lodging facilities;
- 5) includes an extended hierarchy of geographic features relevant to tourism (based on Schema.org) and enriched with GeoNames feature taxonomy (leveraging the GeoNames mapping<sup>21</sup> data-set);
- 6) uses Schema.org to model tourist reviews;
- 7) uses Schema.org to model Tourist Destinations and Tourist Locations;
- 8) can be easily extended to model other kinds of entities relevant to tourism in the future (e.g., events or restaurants).

Figure 2 illustrates the schema of the TAO ontology where the reader can identify the reused classes of the existing ontologies, mentioned above. We will refer to TAO using the `tao:` prefix from now onward. The central classes are `tao:LodgingFacility` and `tao:Accommodation` which are respectively used to model lodging facilities and their accommodations. The `tao:LodgingFacility` class is related to the lodging business concept used in Schema.org, but only refers to the physical place where the accommodations within the facility are located (e.g., a hotel is considered as the building that contains rooms). In this way, there is a clear distinction with the business organisation that governs or owns the lodging facility and no inconsistencies are generated by GoodRelations disjunction between `schema:Place` and `schema:Organization` classes, as discussed in Section III-C1. A facility location

is described according to its latitude and longitude literal properties and also using the `schema:PostalAddress` class, which favors very detailed specification of the address. To complete the facility description we have literal properties for its name (`schema:name`) and a relevant web page (`schema:mainEntityOfPage`). We can use the object property `tao:aggregateRating`<sup>22</sup> to associate a lodging facility to an overall rating, modelled with a node of type `tao:NormAggregateRating`<sup>23</sup> annotated using the data property `tao:normRatingValue` to specify a float value between 0 and 1. A lodging facility can also be associated, through the property `schema:subjectOf`, with a textual description modelled using the `tao:LodgingDescription` class.<sup>24</sup> Finally, lodging facilities can be connected, using the `schema:review` property, to one or more user reviews, modelled using the `schema:UserReview` class. Each review is characterised by the date of creation and associated, using the `schema:reviewRating` property, with a rating (vote) modelled with a `tao:NormRating` class,<sup>25</sup> that can be used to specify the normalised rating in a specific review. The facility description and the reviews can mention every kind of entity, including those defined in other knowledge graphs (DBpedia and GeoNames) using the `schema:mentions` property. This information will be typically extracted from the text of descriptions and reviews with various entity linking techniques, such as DBpedia Spotlight [37] (the solution we have chosen for DBpedia), Mordecai [38] (the solution we have chosen for GeoNames), OpenTapioca [39], or Falcon [40]. Entity linking is the task of linking a portion of texts with their corresponding entities in a knowledge graph [41]. These approaches can be used to identify a variety of entities defined in the external knowledge graphs, such as “Eiffel Tower” or “Paris”.

The `tao:Accommodation` class, analogously to `schema:Accommodation`, represents the actual relevant units of the lodging facility that are offered for rental. It is formally distinct<sup>26</sup> from the physical place where the accommodations are located, which is modelled with the `tao:LodgingFacility` class instead. TAO uses the `tao:includes` object property to define the relation between a lodging facility and one of its accommodations. In order for the TAO ontology to maintain a certain degree of compatibility with the Accommodation Ontology, and potentially reuse semantic entities and annotations expressed using it, we defined the `tao:Accommodation` class as a subclass of `acco:Accommodation`. In this way, if a node in the KG is a member of `tao:Accommodation` it is also a

<sup>19</sup>Definition from Law Insider, see <https://www.lawinsider.com/dictionary/lodging-facilities>

<sup>20</sup>We use the term “hierarchy” to define a subsumption hierarchy of concepts such as the one used by DBpedia, that is a hierarchy of classes connected with `rdfs:subClassOf` property.

<sup>21</sup>[https://www.geonames.org/ontology/mappings\\_v3.01.rdf](https://www.geonames.org/ontology/mappings_v3.01.rdf)

<sup>22</sup>`tao:aggregateRating` is defined as a subproperty of `schema:aggregateRating` (relation not shown in Figure 2).

<sup>23</sup>`tao:NormAggregateRating` is defined as a subclass of `schema:AggregateRating` (relation not shown in Figure 2).

<sup>24</sup>`tao:LodgingDescription` is a subclass of `schema:CreativeWork`. (relation not shown in Figure 2).

<sup>25</sup>`tao:NormRating` is defined as a subclass of `schema:Rating` (relation not shown in Figure 2).

<sup>26</sup>Using `owl:disjointWith` property.

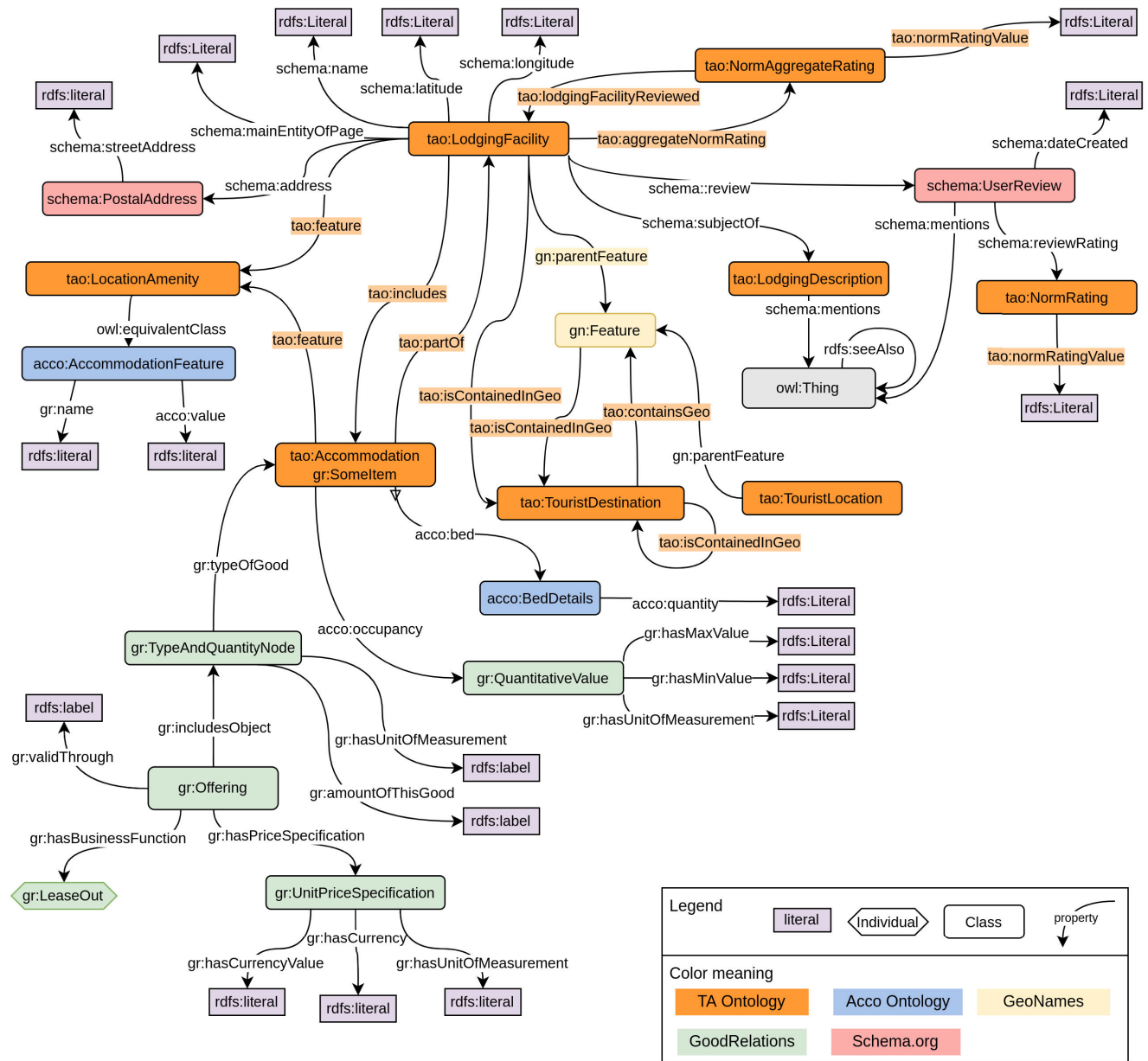


FIGURE 2. TAO ontology schema. In this schema, each arrow represents a semantic relationship, starting from its domain and ending in its range.

member of `acco:Accommodation`, and all the properties defined in the Accomodation ontology for accommodations are still valid. On the contrary, not all the nodes that are members of `acco:Accommodation` are also members of `tao:Accommodation`.

Following GoodRelations best practices, a lease out offering a `tao:Accommodation` individual is modelled using a combination of GoodRelations classes to define the offering price, type, and quantity:

- the individual is also defined by type `gr:SomeItem`<sup>27</sup>;

<sup>27</sup>Besides being of type `tao:Accommodation`.

- the offering itself is modelled with a node of type `gr:Offering`, which has an end of validity expressed with the `gr:validThrough` data property and which is characterised with a specific business function using `gr:hasBusinessFunction` to specify that is a `gr:LeaseOut`<sup>28</sup>;
- the offering includes the accommodation indirectly through a `gr:TypeAndQuantityNode` node using the `gr:includesObject` property and can define

<sup>28</sup>An individual of type `gr:BusinessFunction` defined in the GoodRelations ontology.

its price through a `gr:UnitPriceSpecification` node;

- a `gr:TypeAndQuantityNode` node is used to specify which `tao:Accommodation` node is offered (through the `gr:typeOfGood` relation), the amount of the good included (using `gr:amountOfThisGood` data property) and the unit of measure for the amount included (using `gr:hasUnitOfMeasurement` data property);
- a `gr:UnitPriceSpecification` node is used to specify the price (using `gr:hasCurrencyValue` data property), the currency (using `gr:hasCurrency` data property), and what you are getting for the price (using `gr:hasUnitOfMeasurement`) i.e., a DAY in the accommodation.

The occupancy accommodation is modelled by using the `acco:occupancy` property whose value is a `gr:QuantitativeValue` object, which uses the `gr:hasUnitOfMeasurement` to specify “C62” literal (used by `GoodRelations` to indicate “one piece” of something, in this case, a person<sup>29</sup>) as well as the `gr:hasMinValue` and `gr:hasMaxvalue` relations to define the minimum and maximum number of allowed persons. To model an amenity offered by a lodging facility as a whole or as part of a specific accommodation TAO uses the `tao:LocationAmenity` class, which is defined as an equivalent class of `acco:AccommodationFeature` for compatibility with the Accommodation Ontology. It also uses the `tao:feature` property to associate a lodging facility or an accommodation with one or more amenities.

A tourist location (e.g., London’s Big Ben or the city of Alghero) is a point or area of interest from a tourist point of view and is modelled with a `tao:TouristLocation` class, which is a subclass of both `schema:Place` and `gn:Feature`. A tourist destination (e.g., Sardinia) is defined as a place that is central to the decision to take the trip and is modelled with a `tao:TouristDestination` class, which is declared as `owl:equivalentClass` of `schema:TouristDestination` and as a subclass of `gn:Feature`. Tourist locations and lodging facilities can be included in a tourist destination using the property `tao:isContainedInGeo`.

For instance, if a tourist destination includes the City of London, all `tao:LodgingFacility` individuals in the City of London (according to `gn:parentFeature` property) are also considered within the same destination. This is because the TAO ontology includes an axiom that defines a chain of properties that state that if  $X$  `gn:parentFeature`  $Y$  and  $Y$  `tao:isContainedInGeo`  $Z$ , then  $X$  `tao:isContainedInGeo`  $Z$ , which can be expressed in functional-style syntax as:

```
SubObjectPropertyOf( ObjectProperty
Chain( gn:parentFeature tao:isContained
InGeo ) tao:isContainedInGeo ).
```

TAO includes also several subsumption hierarchies describing the relationships of relevant classes, including:

- 1) the *lodging hierarchy* with 35 types of lodging facilities (e.g., `tao:Hotel`, `tao:Apartment`, `tao:House`) across 4 levels;
- 2) the *accommodation hierarchy* with 17 types of accommodations (e.g., `Room`, `EntireApartment`, `Suite`) across 4 levels;
- 3) the *location amenity hierarchy* with 343 types of amenities (e.g., `Wifi`, `Minigolf`, `Dryer`) across 5 levels;
- 4) the *tourist location hierarchy* with 146 types of tourist locations (e.g., `City`, `Museum`, `Mountain`) across 5 levels;

Appendix B describes these four hierarchies in more detail.

### 3) TAO ENRICHMENT

The TAO ontology was produced using a programmatic approach instead of manual editing. Specifically, we developed a building process in Python using the `owlready2` [42] library. Compared with other approaches based on templates (OPPL [43], OTTR [44]) or on other languages (like Tawny-OWL [45] which is based on Clojure) we preferred the use of a full programming language like Python which is also very well suited to data manipulation and data transformation. This choice also allowed us to apply well-known software engineering tools and practices and automate some aspects of the ontology building process (e.g., creation of axioms), to version the code instead of just the final ontology, to reduce human errors, and to easily produce inline documentation about the ontology creation process. We also release an open-source version of the Python code that builds the TAO ontology as a Jupyter Notebook.<sup>30</sup>

The TAO ontology has to be able to model information derived from typical data sources in the tourism domain, such as Booking.com and Airbnb, which provide (semi)structured data as key/value properties and unstructured data as text regarding lodging facilities, accommodations, amenities, and user reviews. Therefore, we developed a human-in-the-loop strategy, reported in Figure 3, to produce new versions of TAO by continuously enriching the ontology with new types of `tao:LodgingFacility`, `tao:Accommodation` and `tao:LocationAmenity` or new labels for existing types which are derived from the source data. This solution allows us to keep the ontology updated and well-aligned with the actual data.

We start with the basic version of the ontology (orange bullet 1 in the figure), set up external imports, and define classes, properties, and axioms (bullet 2). To further enrich TAO, our ontology engineers in collaboration with domain experts, analyse several analytics about the most

<sup>29</sup><http://www.heppnetz.de/ontologies/goodrelations/v1#UnitPriceSpecification>

<sup>30</sup>See [https://github.com/linkalab/tkg/tree/main/tao\\_modelling](https://github.com/linkalab/tkg/tree/main/tao_modelling)

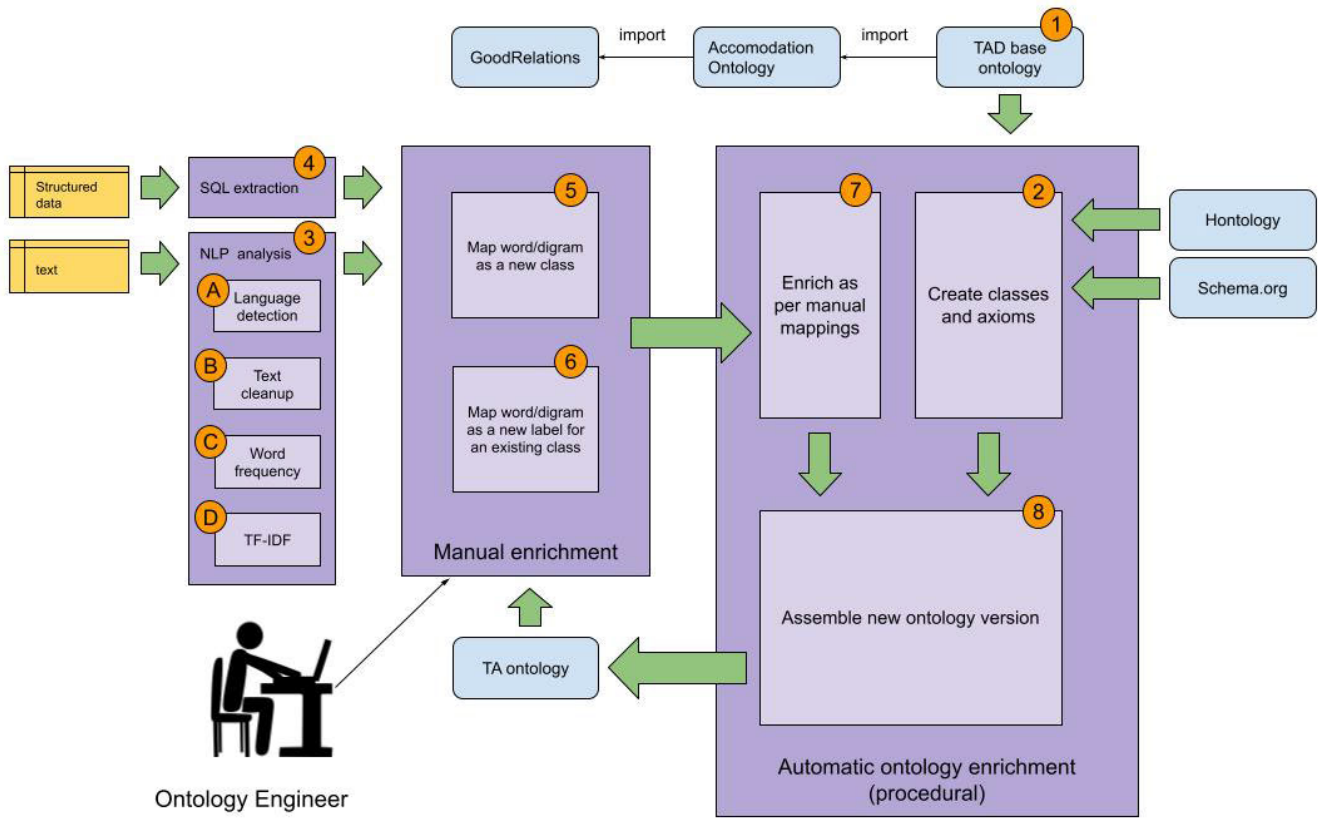


FIGURE 3. Ontology enrichment workflow.

frequent terms associated with facilities, accommodations, and amenities. Then, they use them to create new relevant classes in the ontology (bullet 5) or add additional labels to an existing class (bullet 6). For example, the mini-golf amenity class was identified in the amenities list extracted from Booking.com, while the holiday home lodging facility alternative label “holiday house” was extracted from Airbnb texts. The analytics are produced by two automatic pipelines (3 and 4). The first one processes the unstructured text, extracting and ranking frequent uni-grams and bi-grams from the text descriptions of lodging facilities or user reviews. To achieve this, we relied on Spacy Python library to perform the following sub-tasks: 1) identify language to filter English text only (bullet A), 2) clean the text from special characters (bullet B), 3) perform text frequency analysis (bullet C), and 4) perform TF-IDF analysis (bullet D). The second processes structured data, extracting a list of all possible values for categorical fields that refer to accommodation types, accommodation features, or types of lodging facilities.

Finally, the ontology engineers produce a mapping file that is used (bullet 7) to create new classes, and sub-class relations (using the `rdfs:subClassOf` property) or add labels to existing classes (using the `skos:altLabel` property). We also track the provenance of these changes using the `dc:source` property for classes and the `rdfs:comment`

property for labels. The final process (bullet 8) produces a new version of the TAO ontology.

In Appendix E, we report the code snippets for the iterative extension of the TAO ontology.

#### IV. METHODOLOGY PHASES FOR KNOWLEDGE GRAPH GENERATION

In the following, we will describe the last three phases for the construction of the Knowledge Graph (depicted in Figure 1).

##### A. TRANSFORM THE DATA

The transformation of data is the fourth phase in our approach to building our Tourism Knowledge Graph. Specifically, this phase consists of transforming the information extracted from the data sources into a set of tables, which will be used in the next phase (described in Section IV-B) to produce the actual knowledge graph triples. We devote this section to describing the data transformation process and the technologies for implementing it. Depending on the source data structure and the desired output, we can apply different transformation steps organised as data pipelines. A data pipeline is a series of computational steps organised as a direct acyclic graph where the output of one step becomes the input of one or more downstream steps.



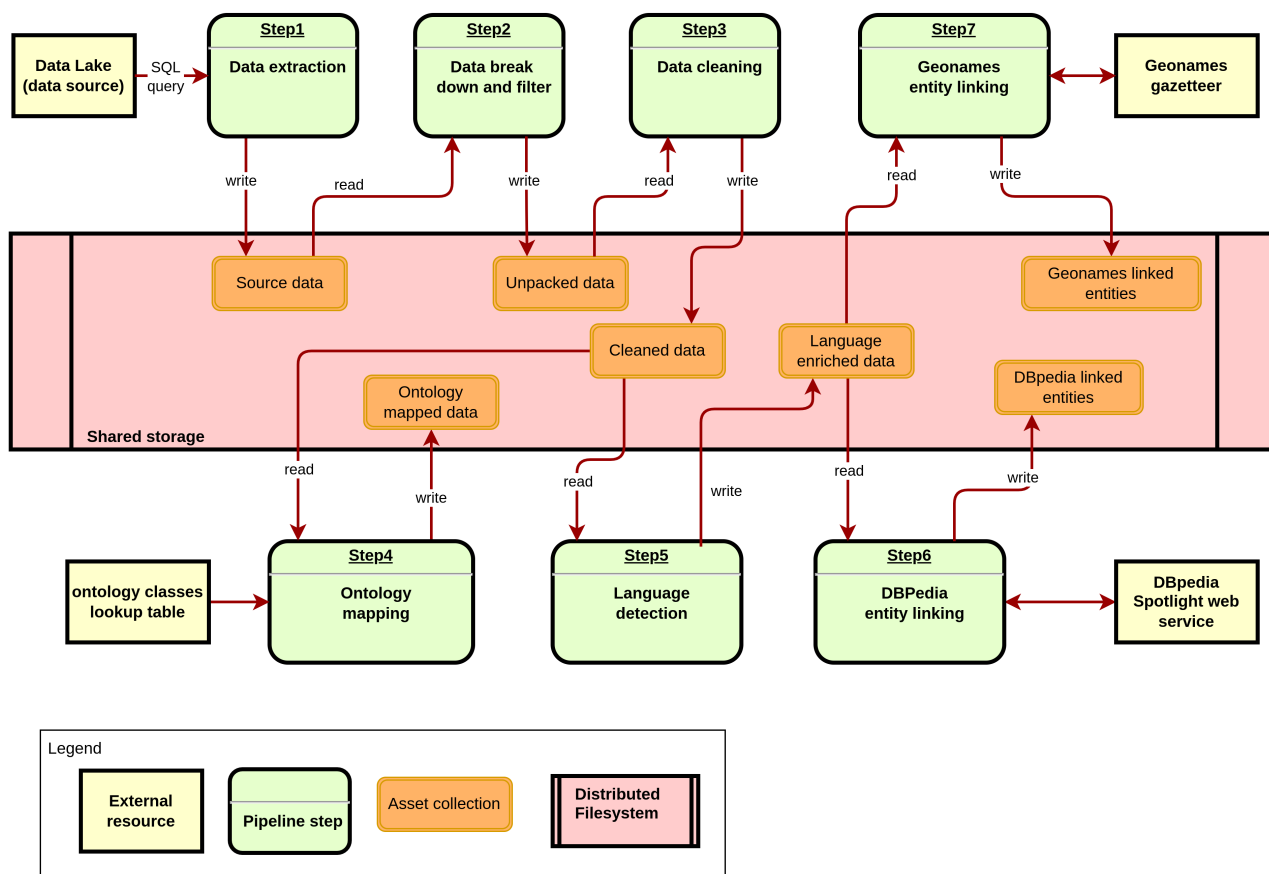


FIGURE 4. High level data transformation workflow diagram.

Figure 4 depicts the complete data transformation workflow. Each step can materialise its output (henceforth referred to as *asset*), saving it as a file or storing it in a database application. From the diagram, we can observe four types of components:

- 1) external resources that are used during the pipeline execution (yellow boxes) representing
  - a) tables in the data lake,
  - b) files mapping text strings to TAO ontology classes,
  - c) DBpedia Spotlight public web service,
  - d) GeoNames gazetteer exposed as an Elasticsearch endpoint;
- 2) pipeline execution steps (green boxes);
- 3) collections of data assets (files) produced by the execution steps (orange boxes);
- 4) a distributed file system that stores all the data assets produced and consumed by one or more processing steps (pink box).

At a high level, the workflow consists of 7 steps. The first 3 steps are executed on both structured (key/values) and unstructured (text) data:

- 1) **Data extraction:** acquires the source data and produces the *Source data assets collection*;

- 2) **Data break down and filter:** rearranges the data structure and filters out unnecessary data; works in combination with the Data cleaning step and materialises the *Unpacked data assets collection*;
- 3) **Data cleaning:** reads from the *Unpacked data assets collection*; corrects or removes corrupt, duplicated or inaccurate data; produces the *Cleaned data assets collection*;

The cleaned data is processed differently depending on if it is structured or unstructured. For structured data, the final step is:

- 4) **Ontology mapping:** uses heuristic rules to identify what ontology class should be used to model each entity described in the data; it produces the *Ontology mapped data assets collection*;

For unstructured data, our objective is to enrich TKG with links from lodging descriptions and user reviews to semantic entities in DBpedia and GeoNames. In this way, TKG would be connected to external knowledge graphs revealing what tourists and business owners are considering important and worth noting. To perform this enrichment we perform entity linking, in three more steps:

- 5) **Language detection**: identifies the language used in texts to process only English text; produces the *Language enriched data assets collection*;
- 6) **DBpedia entity linking**: descriptions and reviews texts are processed to recognise and link DBpedia entities; produces the *DBpedia linked entities data assets collection*;
- 7) **GeoNames entity linking**: descriptions and reviews texts are processed to recognise and link GeoNames entities; produces the *GeoNames linked entities data assets collection*.

In Appendix C, we describe each processing step as well as the employed technological architecture.

## B. TRIPLES CREATION

This section presents the fifth phase for the creation of the Tourism Knowledge Graph, shown in Fig. 1, which deals with the creation of the RDF triples. For this, we leveraged the RDF Mapping Language (RML) [46], to build data pipelines for producing RDF triples<sup>31</sup> from text files, and subsequently save them in a serialised format. The RML language is a declarative language used to define how Linked Data is generated from corresponding data sources, using annotations provided through vocabulary terms. RML can use also files as data sources, which is very useful for our scenario. An RML transformation requires the following elements<sup>32</sup>:

- 1) an RML processor that performs the actual transformation;
- 2) an input to the RML mapping which is called input data source;
- 3) an RML mapping, that defines the rules of conversion from any input (structured) data to RDF.

These rules define how to convert an input record (row, XML element, or JSON object) to one or more RDF triples. They are independent of the process of executing the conversion, thus decoupling the implementation from the rules themselves.

In our implementation, we used RMLMapper [47] which is an open-source RML processor developed in Java.<sup>33</sup> We designed different mappings to handle the different sources, i.e., Booking.com and Airbnb. The output of the RML processor is a set of files containing the RDF triples serialisation in n-quads.<sup>34</sup> To improve the development, debugging and maintenance of RML triple maps we adopted YARRRML [48], a human-readable text-based representation for declarative generation rules.<sup>35</sup> In the following paragraphs, we will examine an example of how a Lodging Facility and all the other related entities can be expressed in TKG by a set of triples created through the process described

<sup>31</sup><https://www.w3.org/TR/2004/REC-rdf-concepts-20040210/#dfn-rdf-triple>

<sup>32</sup>See <https://rml.io/specs/rml/>

<sup>33</sup><https://github.com/RMLio/rmlmapper-java>

<sup>34</sup><https://www.w3.org/TR/n-quads/>

<sup>35</sup><https://rml.io/yarrml/>

above. We will represent triples in a graphical form to better understand the knowledge graph structure.

## 1) HIGH-LEVEL TOURISM KNOWLEDGE GRAPH TRIPLES STRUCTURE

The triples creation process for describing accommodation offers follows the Accommodation Ontology prescriptions and is compliant with GoodRelations and Schema.org best practices. Figure 5 shows an example of TKG structure at a high level. We can observe a lodging facility (:lodging\_1) that is the subject of a descriptive text (:lodging\_description\_1), that has one review (:review\_1), and contains one accommodation (:accommodation\_1). A description is a special kind of creative work (modelled using the tao:LodgingDescription class) that can mention one or more real entities like places or food. In the example, the description mentions the Big Ben tower (through the schema:mentions property). Also, reviews are considered creative works in Schema.org<sup>36</sup> and are thus related to other real-world entities using schema:mentions property. There is an offer (:offer\_1) to lease out an accommodation that is contained in the lodging facility; :offer\_1 is related to the offered accommodation (:accommodation\_1) utilizing (:quantity\_1) node whose properties define what is offered using the property gr:hasUnitOfMeasurement (e.g., DAY) and in what quantity using the property gr:amountOfThisGood (e.g., 2).

In Appendix D we describe in more detail the structure of triples representing lodging facilities, accommodations, offers, and user reviews in the Tourism Knowledge Graph.

## C. KNOWLEDGE GRAPH PUBLISHING AND VALIDATION

In this section, we present the triple store publishing TKG, discuss how to identify the different resources in the knowledge graph, and finally how we encoded the provenance. The TKG validation is covered in Section V-A where the capability of TKG and TAO to address the use cases and offer a useful representation of the tourism domain is considered.

For publishing the knowledge graph we relied on Ontotext GraphDB. The knowledge graph itself is a collection of multiple RDF graphs. Each RDF graph has an associated URI which defines its graph name. For both Booking.com and Airbnb we created two kinds of named graphs:

- 1) hospitality named graph that contains all the triples created using data assets produced at the end of the ontology mapping step processing semi-structured data extracted from a specific source (e.g., Airbnb) for a certain tourist destination (i.e., London or Sardinia);
- 2) linked entities named graph that contains all the triples created using data assets produced at the end

<sup>36</sup>See <https://schema.org/UserReview>



- 3) source:
  - a) bkg: is used to identify the source Booking.com;
  - b) air: is used to identify the source Airbnb;
- 4) enrichment:
  - a) internal: is used for all the RDF assets that are produced with no entity linking during the transformation phase;
  - b) dbpedia\_el: on assets that are enriched with Entity Linking against DBpedia;
  - c) geonames\_el: on assets that are enriched with Entity Linking against GeoNames.

As an example, the named graph which is a collection of triples about London hospitality, produced from Booking.com (semi-)structured data (with no entity linking) would have the following URI: <http://tourism.kg.linkalab-cloud.com/ng/london/bkg/internal>.

The use of named graphs implemented as described simplifies the distinction of resources related to a specific tourist destination because we can use the named graphs in SPARQL queries and identify subsets of data through Implicit Graphs using Triple Pattern Fragments<sup>38</sup> (TPF) [49], [50]. This distinction is also useful to express provenance metadata at the named graph level as described in Section IV-C1.

Concerning identifying a resource in the knowledge graph, we use URIs that explicitly contain the external source (e.g., Booking, Airbnb), and the type of resource. The resource URI is structured as follows: `base_url/resource_type/source/unique_id`

In particular:

- 1) base\_url: <http://tourism.kg.linkalab-cloud.com/>
- 2) resource\_type:
  - a) lf: is used to identify Lodging Facility entities;
  - b) ac: is used to identify Accommodation entities;
  - c) of: is used to identify Offering entities;
  - d) rv: is used to identify User Reviews entities.
- 3) source:
  - a) bkg: the resource is derived from Booking.com;
  - b) air: the resource is derived from Airbnb.
- 4) unique\_id: is an identifier produced by the data transformation phase which is unique for the data source.

As an example, the following URI identifies a lodging facility derived from Airbnb: <http://tourism.kg.linkalab-cloud.com/lf/air/30840569>.

The Tourism Analytics ontology is published as an RDF/XML file at the following URI: <http://purl.org/tao/ns>.<sup>39</sup> To access a specific class or property the hash URI approach is adopted<sup>40</sup> (e.g., <http://purl.org/tao/ns#LodgingFacility> is the URI for LodgingFacility class).

## 1) PROVENANCE AND DATASET METADATA

In a dedicated named graph, we loaded also the metadata triples describing the other named graphs and their provenance: <http://tourism.kg.linkalab-cloud.com/ng/meta/prov>.

<sup>38</sup><http://linkeddatafragments.org/>

<sup>39</sup>This is a redirect to <http://schema.linkalab-cloud.com/tao.rdf>

<sup>40</sup>See <https://www.w3.org/TR/cooluris/#hashuri> for an in-depth explanation

A named graph can be referenced using Quad Pattern Fragments<sup>41</sup> with a URI with the following structure: `base_url?graph=graph_name` where we have:

- 1) base\_url: <http://tourism.ldf.linkalab-cloud.com/graph>
- 2) graph\_name: is the URI associated with the named graph as its name

As an example, the named graph containing the triples about London hospitality produced from Booking.com (semi-)structured data (with no entity linking) would be referenced as:

<http://tourism.ldf.linkalab-cloud.com/graph?graph=http://tourism.kg.linkalab-cloud.com/ng/london/bkg/internal>.

To express the provenance information we used the W3C PROV provenance model. This allows us to track the lineage of data assets produced during the data transformation and triple creation phases following a similar approach as that described in [47] and implemented in the RMLMapper tool. With respect to what is proposed in [47], we applied the Implicit Graphs approach to capture metadata at the Named Graph detail level of granularity, thus generating a minimum number of additional RDF triples for provenance. In this case, the metadata generation time is negligible compared to the overall triple generation time similar to what can be experimented with using the RMLMapper metadata generation feature with a similar configuration.

In PROV we have three main classes:

- `prov:Entity` - a physical, digital, conceptual, or other things with some fixed aspects; entities may be real or imaginary;
- `prov:Activity` - something that occurs over a while and acts upon or with entities; it may include consuming, processing, transforming, modifying, relocating, using, or generating entities;
- `prov:Agent` - something that bears some form of responsibility for an activity taking place, for the existence of an entity, or for another agent's activity.

Figure 6 shows a high-level provenance schema describing how provenance metadata for a specific named graph is modelled. Specifically, we can recognise the following PROV entities:

- 1) `source` - represents the web source for our data (e.g., Booking.com);
- 2) `dataLakeTablesFromSource` - represents the tables exposed by the data lake containing the data extracted from the source;
- 3) `assetsFromSource` - represents all the assets created during the transformation phase which are used to produce the RDF triples for a specific named graph;
- 4) `rmlMapForSource` - represents the RML map document used to produce the RDF triples for a specific named graph;
- 5) `rdfDatasetFromSource` - represents the RDF graph (serialised as one or more files) that is produced from the source using specific `assetsFromSource` and `rmlMapForSource` entities;

<sup>41</sup><https://linkeddatafragments.org/specification/quad-pattern-fragments/>



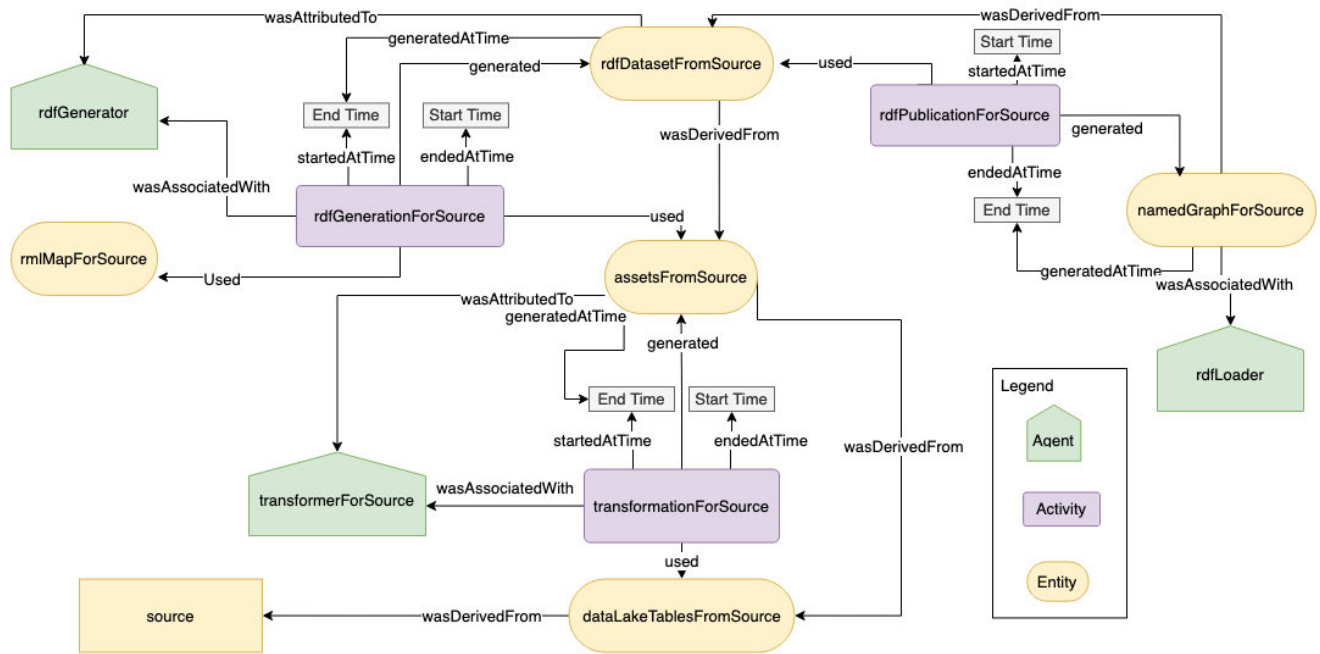


FIGURE 6. A high-level schema of provenance metadata for a named graph in the Tourism Knowledge Graph.

6) `namedGraphForSource` - represents the published named graph.

Moreover, in the same schema, we can identify the PROV Activities involved in the production of a specific named graph:

- 1) `transformationForSource` - performed to prepare/enrich the data for the triple creation;
- 2) `rdfGenerationForSource` - performed to produce the triples;
- 3) `rdfPublicationForSource` - performed to load the triples in the triple store as named graphs.

Finally, we can identify in the schema the following PROV Agents:

- 1) `transformerForSource` - represents the entire transformation pipeline described in Section IV-A;
- 2) `rdfGenerator` - represents the RML processor software (RMLMapper in our case);
- 3) `rdfLoader` - represents the agent that loads the RDF graph in the triple store.

The proposed PROV schema can be easily adapted to specify a particular named graph provenance information and can track: (i) when all triples in the named graph are created/updated, (ii) what assets are used to generate the triples, (iii) what RML mapping document was used to generate them. The same can be specified for all the assets produced by the transformation pipeline. The agent entities are also useful to track the software version used to produce each named graph. It is worth noting that, although RMLMapper software is capable of producing provenance metadata, we decided to use the described provenance schema and a custom metadata generator because

we wanted to cover all the pipelines (data transformation, RDF generation, and RDF publication) using a common approach and leveraging our orchestration service (Dagster) as described in Appendix C-H.

## V. EVALUATION

We evaluated TKG and TAO<sup>42</sup> according to functional, logical, and structural dimensions as suggested by previous works [51], [52]. The *functional dimension* refers to the capability of addressing the requirements and offering a useful representation of the tourism domain while the *logical dimension* is about the ability to be successfully processed by a reasoner and produce sound new knowledge.

We evaluated both functional and logical dimensions by defining and running a set of tests. We implemented the test cases as RDF files modelled with the TestCase OWL meta-model (prefix `test:`), following Blomqvist et al. [53]. Each test case specifies its inputs, conditions for the execution, the actual testing procedure, and the expected results. All the resulting RDF files are available at <https://github.com/linkalab/tkg/tree/main/validation>. Finally, the analysis of the *structural dimension* aims at assessing the topological properties of TKG and TAO, which is also compared with other ontologies (i.e., Hontology and Acco).

<sup>42</sup>It is worth noting that TAO has been also verified using OOPS! (<https://oops.linkeddata.es/>) to find and correct common pitfalls. We manually inspected the results of the tool and, after excluding problems regarding other ontologies or related to incorrect results, we identified and fixed 47 missing annotations, 3 missing domain and range specifications in object properties, 1 wrong equivalent class definition, and 3 inverse relationships not explicitly declared

All these analyses provide useful insights on design choices and can be used to iteratively refine the knowledge graph. We detail them in the following three subsections.

### A. FUNCTIONAL DIMENSIONS

To verify that the functional requirements are satisfied, we followed the **CQ (Competency Question)<sup>43</sup> verification** approach proposed by Carriero et al. [52]. Specifically, this approach aims at testing whether the competency questions can be answered by running SPARQL queries on the KG. To this purpose, we defined 12 test cases by translating the competency questions, defined in Appendix A-B, into SPARQL queries. The input data were selected from the knowledge graph to test each specific functionality. We used this process to drive the creation and refining of TAO, identifying missing classes or properties and adding them to the ontology. We also used it for verifying that TKG can answer in a meaningful way to all competency questions.

The execution of each test case consists of performing the relative SPARQL query against TKG end point. Queries were manually executed and the results were checked against the expected values. Some CQs required the execution of federated queries to access triples from DBpedia and GeoNames. To this end, we used the `SERVICE` keyword to access Ontotext FactForge SPARQL endpoint,<sup>44</sup> which exposes both of them.

All the 12 competency question tests ran successfully. The following example (Listing V-A) shows a federated SPARQL query that aims to answer “*What are the apartments with wi-fi near at least 2 parks?*”<sup>45</sup>

```

PREFIX gdb-geo: <http://www.ontotext.com/
  owlim/geo#>
PREFIX dbo: <http://dbpedia.org/ontology/>
PREFIX gn: <http://www.geonames.org/ontology
  #>
PREFIX tao: <http://purl.org/tao/ns#>
PREFIX acco: <http://purl.org/acco/ns#>
PREFIX schema: <http://schema.org/>
PREFIX onto: <http://www.ontotext.com/>

SELECT ?lodge (SAMPLE(?name) AS ?apartment)
  (COUNT(?park) AS ?num_parks_nearby)
FROM onto:explicit ## use only explicit
  statement without any inference
WHERE {
  { SELECT DISTINCT ?lodge ?name ?lat ?
    long WHERE {
      ?lodge a tao:Apartment; schema:
        latitude ?lat;
      schema:longitude ?long; schema:
        name ?name; tao:feature ?b.
      ?b a tao:Wi-FiZone. } }
  SERVICE <http://factforge.net/
    repositories/ff-news> {

```

<sup>43</sup>A Competency Question is a query expressed in natural language as described in Section III-C.

<sup>44</sup>See <http://factforge.net/>

<sup>45</sup>In this case a park is considered near the apartment if it is within a distance of 1 km.

```

?park gdb-geo:nearby(?lat ?long
  "1km"); gn:featureCode gn:L.
  PRK.
}
}
GROUP BY ?lodge HAVING ( ?num_parks_nearby >
  1)
ORDER BY DESC(?num_parks_nearby)
LIMIT 3

```

The query returns the following results.

Lodge	Apartment	Near_parks
<a href="http://tourism.kg.linkalab-cloud.com/lf/bkg/9bd5bef8f50e0e03">http://tourism.kg.linkalab-cloud.com/lf/bkg/9bd5bef8f50e0e03</a>	"1 Bedroom Luxury Apartment Chancery Lane"	"3" ^ ^xsd:integer
<a href="http://tourism.kg.linkalab-cloud.com/lf/air/42701380">http://tourism.kg.linkalab-cloud.com/lf/air/42701380</a>	"2 bedroom apartment with 50 inch TV"	"3" ^ ^xsd:integer
<a href="http://tourism.kg.linkalab-cloud.com/lf/bkg/51e2e2d011d57200">http://tourism.kg.linkalab-cloud.com/lf/bkg/51e2e2d011d57200</a>	"3 Bedroom Palatial Apartment Chancery Lane"	"3" ^ ^xsd:integer

All competency question test cases are available at [https://github.com/linkalab/tkg/tree/main/validation/competency\\_questions](https://github.com/linkalab/tkg/tree/main/validation/competency_questions)<sup>46</sup>

### B. LOGICAL DIMENSIONS

To assess the logical dimension, we first ran a reasoner on TAO and checked for any inconsistency. We then assessed the full TKG according to two strategies suggested in Carriero et al. [52]:

- 1) **inference verification**, which checks if the inference over the KG produces the expected results (as an example, if a `tao:HotelRoom` accommodation is part of a generic `tao:LodgingFacility` we can infer that the latter is a Hotel);
- 2) **error provocation**, which aims to provoke an inconsistency error by injecting data that violates the requirements (as an example, an instance of a lodging facility can not be defined of type `tao:Hotel` and `tao:BedAndBreakfast` at the same time).

We thus formulate the relevant test cases to assess what inferences can be performed and what types of errors may be produced by the reasoner. In this case, we can no longer rely on CQs but we have to examine the ontology structure and consider how classes and properties are defined by axioms.

In the following subsection, we will describe more in detail how we conducted these two tests.

#### 1) INFERENCE VERIFICATION

For evaluating this dimension, we modelled 15 test cases as OWL files using the TestCase OWL meta-model.

<sup>46</sup>We suggest using Protégé for opening the competency questions test cases files.

These files are identified by a unique IRI and contain only the ABox, relying<sup>47</sup> on the TBox of the TAO ontology and the TestCase metamodel.<sup>48</sup> The ABox contains a set of individuals necessary to execute the test and obtain the expected results. All inference verification test cases and the related data sets are available at [https://github.com/linkalab/tkg/tree/main/validation/inference\\_verification](https://github.com/linkalab/tkg/tree/main/validation/inference_verification).

These tests are useful to understand if the ontology can be successfully used to extend the knowledge graph with reasoning e.g., using inverse properties definitions to materialize backlinks,<sup>49</sup> using a chain of object properties to infer new relationships,<sup>50</sup> inferring the type of an entity from its properties.<sup>51</sup> For example, let us consider a LodgingFacility individual (named Hotel Splendor) which is related to Greater London, a second-level administrative division defined in GeoNames,<sup>52</sup> through the ObjectProperty `gn:parentADM2`. Let us also suppose that there exists a TouristDestination individual called GreatLondonDestination which includes (via the `tao:containsGeo` property) Greater London. Then, the reasoner should infer that Hotel Splendor is also part of GreatLondonDestination. It is worth noting that the creation of inference verification tests has been used also during the ontology engineering process for guiding the introduction and refinement of new axioms in TAO.

We performed the final evaluation by loading the test files in Protégé and running the Pellet reasoner.<sup>53</sup>

All 15 test cases yielded the expected results.

## 2) ERROR PROVOCATION

This test aims at understanding how the knowledge graph (TKG) reacts to the injection of inconsistent data. As an example, since an entity cannot be at the same time a `tao:Hotel` and a `tao:BedAndBreakfast`, we can validate the ontology with regards to this requirement by injecting an individual which is defined as belonging to both classes. The test is successful if the reasoner finds an inconsistency because the appropriate disjointness axiom is defined in the ontology.

We followed the same strategy used in the inference verification tests described above. In addition, for some

<sup>47</sup>Using `owl:imports`.

<sup>48</sup><http://www.ontologydesignpatterns.org/schemas/testannotationschema.owl>

<sup>49</sup>As an example if an Accommodation is `tao:partOf` a lodging facility the inverse relation `tao:includes` can be added to the knowledge graph.

<sup>50</sup>A TouristDestination can be expressed as the composition of other geographic features (using `gn:parentFeature`) so that all lodging facilities contained in those features become also part of the TouristDestination itself.

<sup>51</sup>A lodging facility can be inferred to be of type LowRatedFacility if its normalised rating value is less or equal to a certain value.

<sup>52</sup>See Greater London <http://www.geonames.org/2648110/greater-london.html>

<sup>53</sup>We used the Pellet reasoner, see the Protégé plug-in <https://github.com/stardog-union/pellet/tree/master/protege/plugin>

tests, we developed also a SHACL file defining further constraints.<sup>54</sup>

We implemented 12 test cases for error provocation, testing the identification of wrong patterns in the knowledge graph such as the inclusion of hotel rooms as accommodations in a lodging facility that is not a hotel, the inclusion of accommodation to multiple disjoint lodging facilities, the presence of isolated nodes like a location amenity not connected to any accommodation or lodging facility.<sup>55</sup> We loaded each test file within Protégé, and then we ran both reasoner and the SHACL rules engine.<sup>56</sup> A test is successful if the injected inconsistencies are detected by the reasoner and/or the SHACL validator.

We used this same error provocation technique to test the correct creation of triples during the triple creation process (see section IV-B) and to refine axioms and constraints in TAO.

All error provocation test cases and the related data sets are available at [https://github.com/linkalab/tkg/tree/main/validation/error\\_provocation](https://github.com/linkalab/tkg/tree/main/validation/error_provocation).

## C. STRUCTURAL DIMENSION

We assessed the structural dimension of TAO and TKG by computing different metrics for assessing ontologies and KG that have been defined and used in the literature [51] and [52]. In particular, we followed a similar approach to Carriero et al. [52], which considered both base and topological metrics. Base metrics are used to assess the following quantitative aspects:

- *number of axioms* - the total number of axioms defined for classes, properties, datatype definitions, assertions, and annotations;
- *number of logical axioms* - the number of axioms that affect the logical meaning of an ontology;
- *number of classes* - the total number of classes defined in the ontology;
- *number of object properties* - the total number of object properties defined in the ontology;
- *number of datatype properties* - the total number of datatype properties defined in the ontology;
- *number of annotation assertions* - the total number of annotations in the ontology;
- *DL expressivity* - the description logic expressivity of the ontology.

On the other hand, topological metrics are useful to understand ontology richness, width/depth, inheritance structure, cohesion, and multi-hierarchical degree.

In particular, we adopted the following metrics<sup>57</sup>:

<sup>54</sup>In some tests we use SHACL language to test for integrity constraints that are not limited by the Open World Assumption (OWA).

<sup>55</sup>This case requires the use of SHACL rules because of the open world assumption in OWL.

<sup>56</sup>Using SHACL4Protege Constraint Validator, see <https://github.com/fekaputra/shacl-plugin>

<sup>57</sup>For a more detailed description of the following metrics see: <https://ontometrics.informatik.uni-rostock.de/wiki/index.php/OntoMetrics>

- *Inheritance Richness (IR)* - measures the average number of sub-classes per class. Low values indicate a vertical (deep) ontology whereas high values indicate a horizontal (shallow) ontology.
- *Relationship Richness* - measures the ratio of the number of non-inheritance relationships divided by the number of relationships of all kinds. Values are normalised to one, where 0 indicates that only inheritance relations exist in the ontology and 1 that no inheritance relations are present.
- *Axiom Class Ratio* - measures the ratio of the number of axioms divided by the number of classes. A scarcely axiomatised ontology has a low value of this metric (near zero); higher values are an indication of a better axiomatisation, but very high values can state an excessive axiomatisation.
- *Class/property ratio* - measures the ratio of the number of classes divided by the number of relations. Low values (i.e.,  $\sim 0$ ) are found in ontologies with many properties connecting a few concepts. On the contrary, high values indicate that the ontology has many classes connected by few properties.
- *NoR* - number of root classes (a class which is not a subclass of other classes). The interpretation of NoR depends on the total number of classes. We expose (i) the ordinal values of NoR and (ii) the ratios between NoR and the number of classes between parenthesis.
- *NoL* - number of leaf classes (all classes that have no sub-classes). The interpretation of NoL depends on the total number of classes. We expose (i) the ordinal values of NoL and (ii) the ratios between NoL and the number of classes between parenthesis.
- *NoC* - number of external classes<sup>58</sup> defined by [54]. A low value of NoC can indicate that the ontology is semantically independent; a high value can indicate that the ontology depends on concepts defined in other ontologies. The interpretation of NoC depends also on the number of classes in ontology. We expose (i) the ordinal values of NoC and (ii) the ratios between NoC and the number of classes between parenthesis.
- *ADIT-LN* (Average depth of inheritance tree of leaf nodes) - is the average depth of the graph constructed considering classes as nodes and `subClassOf` properties as arcs.
- *Max breadth* - the maximal value of breadth computed on the graph constructed as for the *ADIT-LN* metric. The value of *Max breadth* should be considered concerning the number of classes in ontology.
- *Average breadth* - the average breadth computed on the graph constructed as for the *ADIT-LN* metric.
- *Max depth* - the maximal depth obtained by traversing the graph constructed as for the *ADIT-LN* metric. The

<sup>58</sup>A class is considered external when it is defined in a different ontology. This metric has been calculated using Protégé.

TABLE 1. Base metrics.

Metric name	TAO	Hontology	Acco
# axioms	3960	1453	344
# logical axioms	1237	448	111
# classes	588	284	31
# object properties	19	8	21
# datatype properties	3	31	14
# annotation assertions	2074	682	161
DL expressivity	SROIQ(D)	ALCHQ(D)	ALUH(D)

TABLE 2. Number of classes by tourism aspect.

Metric name	TAO	Hontology	Acco
Lodging facility types	35	19 <sup>1</sup>	5 <sup>5</sup>
Accommodation types	17	54 <sup>2</sup>	4 <sup>6</sup>
Amenities types	343	93 <sup>3</sup>	not provided <sup>7</sup>
Tourist location types	146	22 <sup>4</sup>	not provided <sup>8</sup>

1. ho:Accommodation sub classes

2. ho:Room sub-classes

3. ho:Facility sub-classes types

4. union of ho:PointOfInterest and ho:Location sub-classes

5. Only selected sub-classes of acco:Accommodation

6. Only selected sub-classes of acco:Accommodation

7. Class acco:AccommodationFeature can hold feature information using acco:value and gr:name data properties to create custom sub-classes.

8. Acco does not model tourist locations.

value of *Max depth* should be considered concerning the number of classes in ontology.

- *Tangledness* - is the degree of multi-hierarchical classes (which are classes with more than one super-class). It is related to the multi-hierarchical nodes of the graph constructed for the *ADIT-LN* metric. A value of 0 indicates no tangledness; a value of 1 indicates that each class has multiple super-classes.

Tables 1, 2, and 3 report the base and topological metrics measured on TAO, Hontology, and the Accommodation Ontology (Acco). It should be noted that when analysing TAO we considered only the classes and properties defined in this ontology and not the ones imported from other ontologies (i.e., the Accommodation Ontology, GoodRelations). This was done to allow a fair comparison with the Accommodation Ontology, which we import.

All metrics were calculated using OntoMetrics<sup>59</sup> web tool.

Table 1 shows that TAO is significantly larger than Hontology and Accommodation Ontology in terms of a number of classes, axioms, logical axioms,<sup>60</sup> and annotation assertions. The additional classes mostly describe different types of lodging facilities (35 classes), accommodations (17 classes), amenities (343 classes), and tourist locations

<sup>59</sup>See <https://ontometrics.informatik.uni-rostock.de/ontologymetrics/index.jsp>

<sup>60</sup>Logical axioms affect the logical meaning of an ontology. See [https://ontometrics.informatik.uni-rostock.de/wiki/index.php/Base\\_Metrics#Logical\\_Axiom](https://ontometrics.informatik.uni-rostock.de/wiki/index.php/Base_Metrics#Logical_Axiom). On the other hand, non-logical axioms, like entity declarations or annotations, do not affect the consequences of an OWL 2 ontology. See [https://www.w3.org/TR/owl2-syntax/#Entity\\_Declarations\\_and\\_Typing](https://www.w3.org/TR/owl2-syntax/#Entity_Declarations_and_Typing)



TABLE 3. Topological metrics.

Metric name	TAO	Hontology	Acco
Inheritance Richness	1.177	0.961	0.742
Relationship Richness	0.413	0.321	0.477
Axiom Class Ratio	6.735	5.116	11.097
Class/property ratio	0.499	0.706	0.705
NoR	14 (0.02)	17 (0.06)	13 (0.42)
NoL	494 (0.84)	247 (0.87)	23 (0.74)
NoC	19 (0.03)	0 (0.00)	2 (0.06)
ADIT-LN	3.612	2.725	2.439
Max depth	6	5	3
Average breadth	6.578	7.375	5.077
Max breadth	54	29	13
Tangledness	0.179	0.018	0.097

(146 classes). Table 2 shows a comparison of the mentioned classes. TAO has more types of lodging facilities, amenities, and tourist locations with respect to Hontology and Acco. Hontology has apparently more accommodation types, but this number may be due to the fact that these types actually combine room types with amenities (e.g., `ho:FamilyRoomWithBalcony`) or the number of beds (e.g., `ho:SingleRoom`, `ho:10BedFemaleDorm`). In this case, we preferred to avoid the addition of specific sub-classes but instead, we used amenities (e.g., `acco:Terrace`) and bed details specifications (using `acco:BedDetails`) to better characterize accommodations.

In terms of properties, TAO introduces only a few new ones, since it reuses most of them from Acco (4), GoodRelations (15), Schema.org (11), and GeoNames (1) as discussed in Section III-C2.

Finally, in terms of expressivity, TAO is similar to Hontology because they share ALCQU features and Acco because they share ALCU features; TAO does not have the H feature because it does not express role hierarchies (SubPropertyOf) as Hontology and Acco; however, TAO has the IS features, indicating the presence of inverse and transitive roles (relations), that the other two ontologies do not have.

The indicators in Tables 1, 2 and 3 can be used to assess and compare TAO, Hontology, and the Accommodation Ontology according to their *transparency*, *flexibility*, and *cognitive ergonomics* [51]. *Transparency* has been defined as “the property of an ontology to be analysed in detail, with a rich formalisation of conceptual choices and motivation”. *Flexibility* is related to how easy is to change and evolve the ontology with limited side effects. Finally, *cognitive ergonomics* is the ability of an ontology to be “easily understood, manipulated, and exploited by final users”. In the following, we discuss the main indicators of these properties.

TAO performs well according several indicators of transparency [51] as it offers:

- a relative *high number of axioms per class* (6.578). This is higher than Hontology, but lower than

Accommodation Ontology, mostly due to the much lower number of classes in the latter;

- a *small coupling* with external ontologies (0.03), similar to Hontology (0) and the Accommodation Ontology (0.06). This is computed as the number of external classes defined in other ontologies (NoC) normalized by the total number of classes. Low coupling allows users to inspect and understand an ontology.
- a *strong cohesion* (i.e., relatedness among classes) due to the low depth of the class hierarchy (ADIT-LN = 3.612), the small number of root classes (NoR = 14), and the high number of leaf classes (NoL = 494);
- a *high inheritance richness* (1.177), which accounts for a more vertical structure, reflecting a more comprehensive coverage of the tourism domain. This is higher than both Hontology (0.961) and Accommodation Ontology (0.742).

The combination of *low coupling* and *strong cohesion* are also indicators of *flexibility* [51].

Finally, TAO exhibits several indicators that are typically associated with a good *cognitive ergonomics*, such as:

- a relatively *low class/property ratio* (0.499), also smaller than Hontology (0.706) and Accommodation Ontology (0.705);
- a *sub-class tree with low depth and breadth* as indicated by ADIT-LN (3.612), max depth (6), and average breadth (6.578);
- a relatively *low tangledness* (0.179 in a range from 0 to 1) that suggests that the inheritance tree has low complexity.

Table 4 reports some statistics about the current prototype of TKG, which includes over 10M triples describing 35K facilities and almost 898K reviews.

Figure 7 shows the distribution of individuals in terms of classes. The most frequent classes are (i) `tao:NormRating` and `schema:UserReview` which are used for reviews; (ii) `acco:AccommodationFeature`<sup>61</sup> that is used as a generic class for amenities together with a specific class from `tao` (e.g., `tao:Kitchen`, `tao:Television`); (iii) the classes used to model an offer such as `gr:Offering`, `gr:TypeAndQuantityNode`, and `gr:UnitPriceSpecification`; (iv) `tao:Accommodation`, `gr:QuantitativeValue`, `gr:SomeItems`, and `acco:BedDetails` are the classes used to model an accommodation; (v) `tao:LodgingDescription`, `tao:LodgingFacility` (and its subclasses), `schema:PostalAddress`, and `tao:NormAggregateRating` that are used to model the lodging facilities. The other classes in the diagram are sub-classes of `tao:LocationAmenity`, `tao:Accommodation` or `tao:LodgingFacility`, which are used to specify precisely their type.

<sup>61</sup>`tao:LocationAmenity` is defined as an equivalent class to `acco:AccommodationFeature`.

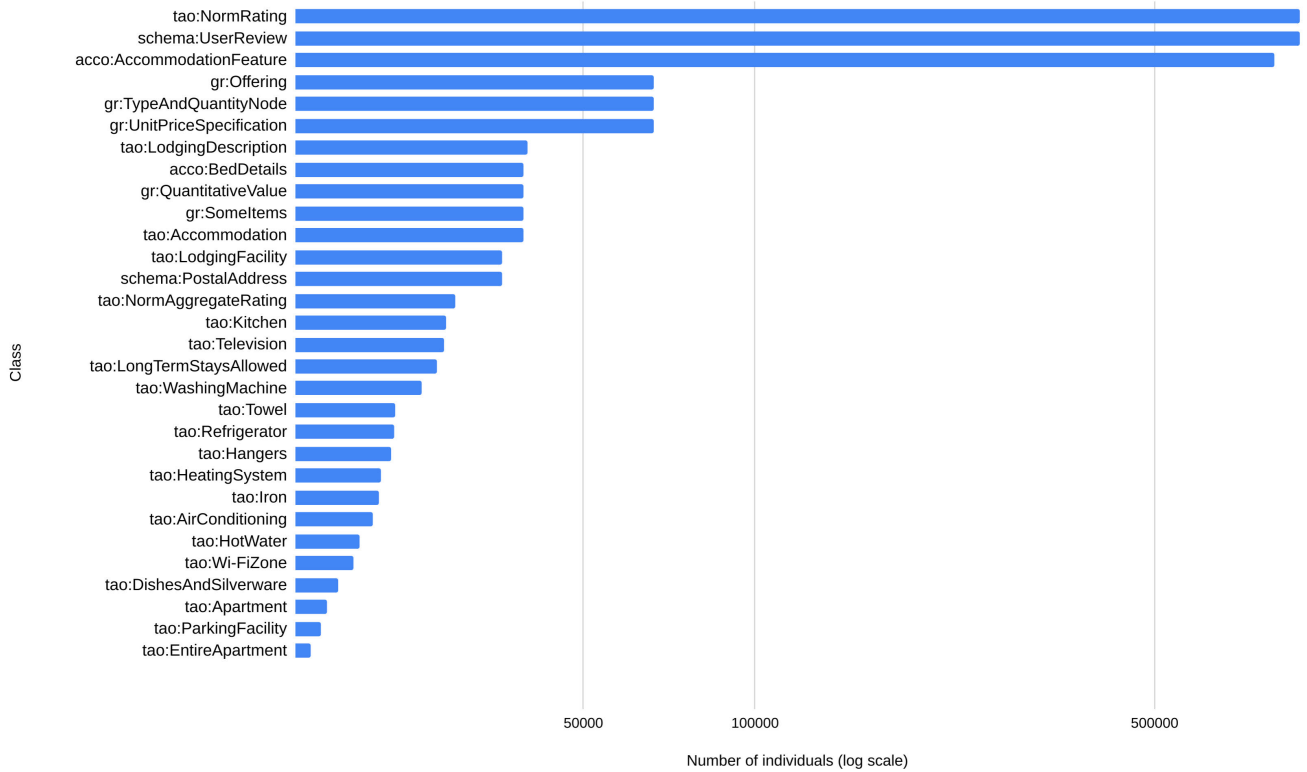


FIGURE 7. Top 30 classes by the number of individuals in the knowledge graph.

TABLE 4. Knowledge graph metrics.

Metric	Value
Number of triples	10,917,081
Number of distinct relations	146
Number of links to DBpedia entities	210,245
Number of unique DBpedia entities linked	3,851
Number of links to GeoNames entities	142,043
Number of unique GeoNames entities linked	3,487
Number of Airbnb reviews entities	358,005
Number of Booking.com reviews entities	539,834
Number of Airbnb LodgingFacility entities	29,870
Number of Booking.com LodgingFacility entities	6,126

## VI. CONCLUSION

In this paper, we presented a framework for the semi-automatic construction of a Tourism Knowledge Graph (TKG) and introduced a new ontology for modelling this domain: the Tourism Analytics Ontology (TAO). We have evaluated TKG and TAO according to functional, logical, and structural dimensions.

The evaluation suggests that TAO is i) larger than the alternatives (Hontology and Accommodation Ontology) in terms of the number of classes and axioms and ii) also offers higher transparency, flexibility, and cognitive ergonomics.

In future work, we aim to pursue three main pathways. First, we are working on developing NLP solutions to improve the extraction of entities from text, such as

descriptions and reviews, so to further enrich the representation of lodging facilities. This step includes the extraction of data from other sources related to several other touristic destinations. Solutions such as Entity Fishing<sup>62</sup> or Open Information Extraction<sup>63</sup>

Second, we want to develop a more scalable solution for integrating data about millions of facilities and users. To achieve such a goal, we will rely on big data frameworks such as Apache Spark and Elasticsearch running in a cluster of machines on cloud computing facilities and we will implement a dedicated DBpedia Spotlight web service to speedup the entity linking process. Third, we want to develop a range of intelligent services based on TKG, including an entity-linking application for automatically annotating accommodations according to reviews and a conversational agent able to answer questions regarding the tourism sector. Knowledge Graph completion will provide a means to predict relations between entities of the knowledge graph and will be performed by leveraging Knowledge Graph Embedding models (e.g., TransE [55], RotatE [56], ComplexE [57]) or methods based on Graph Neural Networks [58], path-based features [59] and Few-Shot Learning [60].

Transversally to them, we want to extend TAO ontology in order to model other aspects related to tourism, starting

<sup>62</sup><https://github.com/kermitt2/entity-fishing>

<sup>63</sup><https://openie.allenai.org/> can be leveraged for named entity extraction, including entity detection, name resolution, and named entity recognition

with events and restaurants. We also plan to explore other APIs relevant to tourism such as Google Hotel API,<sup>64</sup> Google Places API<sup>65</sup> or TripAdvisor<sup>66</sup> we are working on automatising as much as possible the pipeline we have used intending to create knowledge graphs with related ontologies in any domain and sources.

## APPENDIX A REQUIREMENTS AND COMPETENCY QUESTIONS

In this section we detail the functional and non-functional requirements identified during the definition of the domain ontology (TAO), described in Section III-C, as well as the relevant use cases. We also describe the Competency Questions and how they express functional requirements in a more operative form. Finally, we examine the information available in the data sources that supported the formulation of the CQs.

### A. REQUIREMENTS

To successfully be used to model a knowledge graph that can support the use cases identified in Section III-A, we envisaged that the ontology would need to fulfill the following functional requirements (FR):

- FR 1** model lodging facilities and define a hierarchy<sup>67</sup> of their types (e.g., hotels, hostels, apartments),
- FR 2** model accommodations and define a hierarchy of their types (e.g., room, entire apartment, suite);
- FR 3** model amenities offered to tourists and define a hierarchy of their types (e.g., disable access, parking garage, baby monitor);
- FR 4** model tourist locations (e.g., waterfall, beach, museum, park) and define a hierarchy of their types;
- FR 5** model the relations among entities (e.g., geographic relations, mentions, composition/inclusion);
- FR 6** model tourist reviews;
- FR 7** model tourist destinations (e.g., Sardinia, London), which is the place that is central to the trip.

Functional requirements for the ontology are mapped to the knowledge graph's use cases as described in Table 5. As an example, we can see that since "KG should support the identification of the topics of interest discussed by tourists in their reviews" the ontology should model user reviews (FR6) and concepts typically related to what tourists speak about as lodging facilities (FR1), accommodations (FR2), amenities (FR3), and tourist locations (FR4).

When considering non-functional requirements (NFR), the ontology should support reasoning and be based on widely adopted technical and market standards. In particular:

- NFR 1** should be defined in OWL<sup>68</sup>;
- NFR 2** should be based on two *de-facto* standards to model business data:
  - Schema.org,<sup>69</sup> which is a set of vocabularies developed through a collaborative effort for structuring data on the web. It was originally founded by Google, Microsoft, Yahoo, and Yandex.
  - GoodRelations, which is a lightweight ontology for exchanging e-commerce information, namely data about products, offers, points of sale, prices, terms, and conditions, on the Web.
- NFR 3** should be easy to extend in order to cover other use cases in the tourism domain.

### B. COMPETENCY QUESTIONS

Based on the functional requirements we defined the following 12 competency questions:

- CQ 1** Which are the first n (e.g., 10) lodging facilities of a specific type (e.g, hotels) with more than m (e.g, 1,000) reviews and the lowest mean value of users' review scores?
- CQ 2** Find three apartments with a specific amenity (e.g.,Wi-Fi), within a specific distance Km (e.g. 2Km) from at least a specific number (e.g. 2) of tourist attraction (e.g., Parks).
- CQ 3** Which Tourist Destinations have the highest percentage of high-priced Lodging Facilities (at least one offer for accommodation for two persons with a nightly price two times over the mean price)?
- CQ 4** What are the n (e.g, 10) tourist locations cited most by hotel descriptions that also offer a specific amenity (e.g., day Spa) in a specific tourist destination?
- CQ 5** What are the most cited Tourist Locations in all Lodging Facility descriptions within a certain tourist destination?
- CQ 6** What are the Tourist Locations cited most in positive user reviews?
- CQ 7** What are the n (e.g, 10) cheapest apartments that offer at least m (e.g., 2) beds and a specific amenity (e.g., secured parking) and are within a certain distance (e.g., 10km) from a certain type of tourist attraction (e.g, airport)?
- CQ 8** Which type of Lodging Facility is more reviewed by tourists in a specific Tourist Destination?
- CQ 9** What are the top Tourist Destinations with respect to positive sentiment about food (i.e., percentage of Lodging Facilities with positive reviews that cite food)?

<sup>64</sup>See <https://developers.google.com/hotels>

<sup>65</sup>See <https://developers.google.com/maps/documentation/places/web-service>

<sup>66</sup>See <https://www.tripadvisor.com/developers> to reuse and extend TAO to model also their data in a unified way. To conclude

<sup>67</sup>For a description of hierarchies and their implementation in the TAO ontology see Appendix B.

<sup>68</sup>More specifically it should be based on OWL DL dialect which is designed to provide the maximum expressiveness possible while retaining computational completeness, decidability, and the availability of practical reasoning algorithms.

<sup>69</sup>See <https://schema.org/>

**TABLE 5. Mapping knowledge graph’s use cases with ontology’s functional requirements.**

Use Case	FR1	FR2	FR3	FR4	FR5	FR6	FR7
UC1 - KG should support the identification of the topics of interest discussed by tourists in their reviews	X	X	X	X		X	
UC2 - KG should support the identification of the topics of interest presented in the descriptions of lodging facilities and accommodation offers	X	X	X	X			
UC3 - KG should support the recognition and linking of tourism entities in the KG for different applications revolving in the domain of social media, news, and blogs	X	X	X	X			
UC4 - KG should support sentiment analysis applications about tourists toward lodging businesses and destinations	X	X	X	X	X		X
UC5 - KG should support the classification of tourist destinations on the basis of what they offer and on the basis of tourist opinions	X	X	X	X	X	X	X

**TABLE 6. Mapping competency questions with functional requirements.**

	FR1	FR2	FR3	FR4	FR5	FR6	FR7
CQ1	X					X	
CQ2	X	X	X	X	X		
CQ3	X	X			X		X
CQ4			X	X	X	X	X
CQ5	X			X	X		X
CQ6				X	X	X	
CQ7		X	X	X	X		
CQ8	X				X	X	X
CQ9	X				X	X	X
CQ10	X				X	X	
CQ11				X	X	X	
CQ12		X			X		X

**CQ 10** In which months do we have the highest number of user reviews for lodging facilities of a specific type (e.g, hotels)?

**CQ 11** What Tourist Locations can be found in a Tourist Destination?

**CQ 12** How many beds are offered on lease in a certain Tourist Destination?

As we can see, CQs can be more generic or specific depending on which aspect of the ontology we want to describe and eventually test, but all CQs are expressed in terms of questions that can be translated into SPARQL queries against the KG. This is why in some CQs we can use concrete examples (i.e., Wi-Fi) instead of more generic entity classes (i.e., “a location amenity”).

A given competency question usually includes information related to different functional requirements and vice-versa, a certain functional requirement is covered by different competency questions. We can see the mapping between CQs and functional requirements in Table 6.

**C. INFORMATION IN DATA SOURCES**

The formulation of the competency questions was also supported by the information available in the data sources. Here, we report a list of the most relevant information

available in the data sources (discussed in Section III-B) that drove the CQs formulation:

- 1) information about lodging facilities:
  - a) name(s)
  - b) position
  - c) geographic relations with administrative divisions
  - d) geographic relations with tourist destinations
  - e) type (e.g., Hotel, Resort, Motel, B&B, Holiday Accommodations)
  - f) type of accommodation offered (e.g., room, apartment, villa, bungalow, etc.)
  - g) amenities (e.g., sauna, parking, swimming pool, breakfast, air conditioning, etc.)
  - h) accommodation prices exposed on the web
  - i) user ratings
  - j) textual descriptions (to perform Named Entity Recognition, Entity Linking and Relation Extraction, etc.)
- 2) information about tourist locations:
  - a) name (in multiple languages)
  - b) position
  - c) geographic relations with administrative divisions
  - d) geographic relations with tourist destinations
- 3) information about tourist destinations:
  - a) name (in multiple languages)
  - b) position
  - c) geographic relations with administrative divisions
  - d) geographic relations with tourist locations
- 4) tourist reviews about lodging businesses and locations
  - a) user votes
  - b) tourist nationality and type of tourist (family, couple, etc.)
  - c) textual review (to perform Named Entity Recognition, Entity Linking and Relation Extraction, etc.)

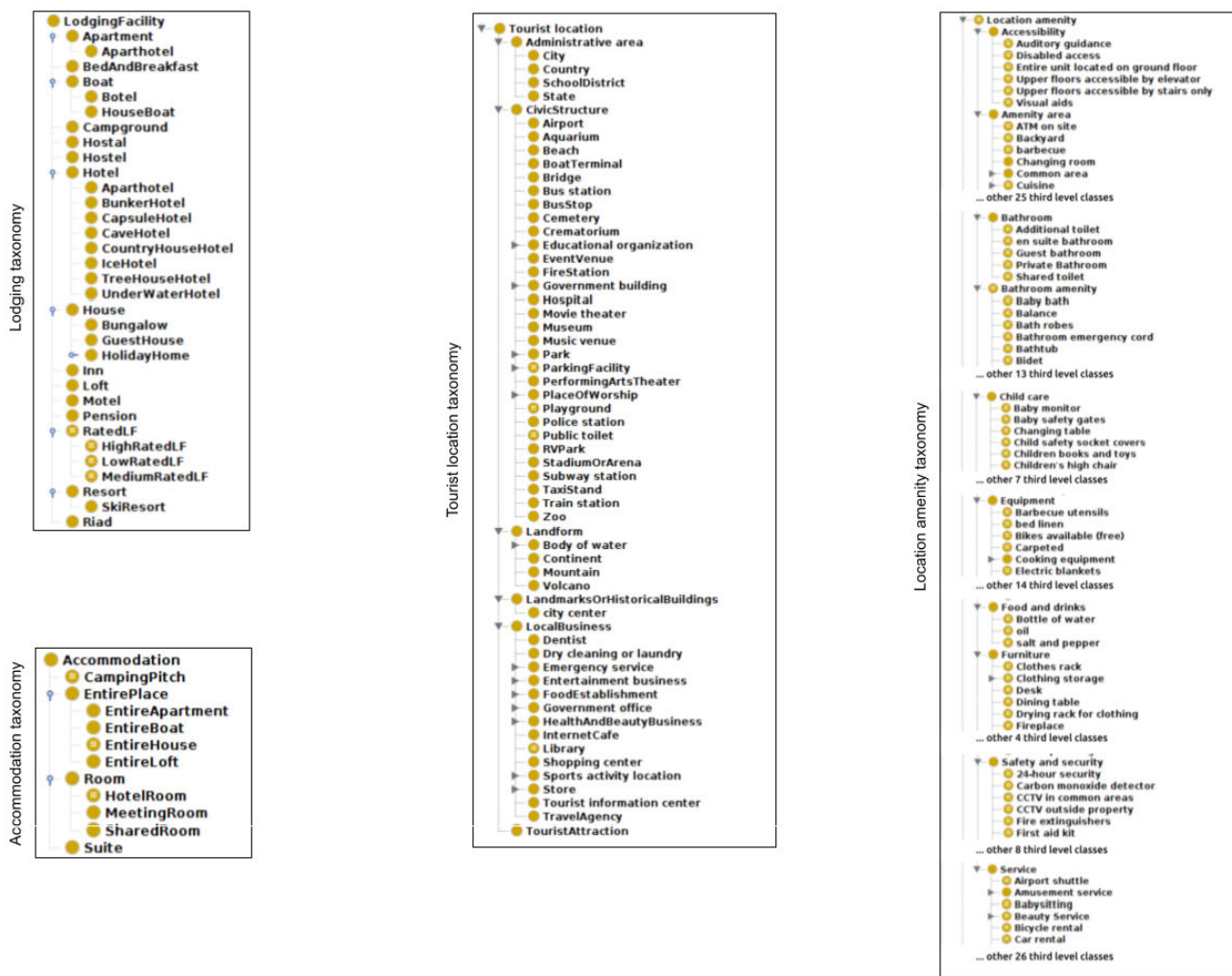
This list was employed during the process of ontology engineering, as it helps to define the set of entities and properties that should be modelled by the TAO ontology.

**APPENDIX B**

**CLASS HIERARCHIES IN TAO**

TAO includes several hierarchies of classes connected with rdfs:subClassOf property. This approach was chosen above





**FIGURE 8.** A tree representation of the four hierarchies included in the TAO ontology expanded to the third level (some class removed in the location amenity hierarchy for sake of clarity and space).

others (e.g., model taxonomies using SKOS<sup>70</sup>) because we wanted to be compatible with the Accommodation Ontology (where accommodations and amenities types are represented as sub-classes) and simplify the use of Schema.org where class hierarchies are also used. In particular, we have four hierarchies describing the relationships of relevant classes, including:

- 1) the *lodging hierarchy* with 35 types of lodging facilities (e.g., `tao:Hotel`, `tao:Apartment`, `tao:House`) across 4 levels;
- 2) the *accommodation hierarchy* with 17 types of accommodations (e.g., `Room`, `EntireApartment`, `Suite`) across 4 levels;
- 3) the *location amenity hierarchy* with 343 types of amenities (e.g., `Wifi`, `Minigolf`, `Dryer`) across 5 levels;

- 4) the *tourist location hierarchy* with 146 types of tourist locations (e.g., `City`, `Museum`, `Mountain`) across 5 levels;

Figure 8 reports the first three levels of each hierarchy. For each sub-class in a hierarchy we can have one or more of the following implementations:

- if a class is conceptually related to a similar class in other ontologies (e.g., DBpedia), this is modelled with the annotation property `rdfs:seeAlso`;
- if a class is derived from other ontologies, we track the provenance using the `dc:source` property to indicate the original class<sup>71</sup>;
- if a class extension<sup>72</sup> is the same as the extension of a class in other ontologies we link them

<sup>70</sup>See <https://www.w3.org/TR/2009/REC-skos-reference-20090818/>

<sup>71</sup>Note that `dc:` stands for Dublin Core.

<sup>72</sup>The set of individuals that are members of the class.

with the `owl:equivalentClass` property,<sup>73</sup> or the `rdfs:subClassOf` property if it is narrower<sup>74</sup>;

- for each class, we use `rdfs:label` to indicate the primary label and `skos:altLabel` to indicate alternate labels;
- disjoint axioms are added when appropriate to better support the reasoning.

### A. LODGING TAXONOMY

The first hierarchy describes the different types of lodging facilities and their sub-types like in the case of *Aparthotel*, which is a special case of a hotel. We also introduce a special case with `tao:RatedLF` and its sub-classes which are used to classify *Lodging facilities* according to their ratings (`tao:NormAggregateRating`). Specifically, `tao:NormAggregateRating` has 3 sub-classes: `tao:LowNormRating`, `tao:MediumNormRating` and `tao:HighNormRating`. These classes can be extended using a data property restriction<sup>75</sup> on `tao:normRatingValue` to implement an automatic classification of a *Lodging facility*.

A rated lodging facility is also part of `tao:RatedLF` (rated lodging facility) class<sup>76</sup> and it can also be inferred whether it is part of one of the following three sub-classes:

- is part of `tao:HighRatedLF` class if it is associated<sup>77</sup> with a `tao:HighNormRating` node;
- is part of `tao:MediumRatedLF` class if it is associated with a `tao:MediumNormRating` node;
- is part of `tao:LowRatedLF` class if it is associated with a `tao:LowNormRating` node;

### B. ACCOMMODATION HIERARCHY

When modelling *accommodations*, we distinguished two general offerings: (i) entire place (i.e., *EntirePlace*), and (ii) room (i.e., *Room*). For these, we also defined sub-classes (e.g., *EntireHouse* for *EntirePlace*, *HotelRoom* for *Room*). In addition, we modelled two special cases (i.e., *CampingPitch* and *Suite*), which are not covered by the general cases. When appropriate, we used equivalence axioms to add useful constraints as in the case of *HotelRoom* which must be part of one *Hotel*. Moreover, to support high compatibility between TAO and the Accommodation Ontology, we defined the accommodation classes of TAO as subclasses of the Accommodation Ontology ones (e.g., `tao:CampingPitch` is a subclass of `acco:CampingPitch`).

<sup>73</sup>It is the case of `tao:TouristDestination` which is declared to be `owl:equivalentClass` of `schema:TouristDestination`.

<sup>74</sup>It is the case of `tao:EntireApartment` which is declared to be `rdfs:subClassOf` of `acco:Apartment` because in the Accommodation ontology `acco:Apartment` can refer to an apartment as a lodging facility or as an actual accommodation offered on lease.

<sup>75</sup>See OWL2 specifications [https://www.w3.org/TR/2012/REC-owl2-syntax-20121211/#Data\\_Property\\_Restrictions](https://www.w3.org/TR/2012/REC-owl2-syntax-20121211/#Data_Property_Restrictions)

<sup>76</sup>Because this class is defined using an existential quantification on the object property `tao:aggregateNormRating` that has some `tao:NormAggregateRating`.

<sup>77</sup>Using `tao:aggregateNormRating` object property.

### C. LOCATION AMENITY HIERARCHY

In the case of *location amenities*, we added equivalence axioms to support a certain degree of mapping with how specific accommodation features could be more probably defined using the Accommodation ontology approach.<sup>78</sup> To this end, each sub-class in this hierarchy is also declared as `owl:equivalentClass` to an anonymous class defined in accordance to Accommodation Ontology prescriptions.<sup>79</sup> Thus we define each anonymous class as a subclass of `acco:AccommodationFeature` and as an `owl:intersectionOf` of `owl:Restriction` based on `gr:name` and `acco:value` data properties from *GoodRelations*. An example is given below in Turtle:

```
tao:AirportShuttle rdf:type owl:Class;
owl:equivalentClass [
  rdf:type owl:Class
  owl:intersectionOf (
    acco:AccommodationFeature
    [
      rdf:type owl:Restriction;
      owl:onProperty acco:value;
      owl:hasValue "yes"@en
    ]
    [
      rdf:type owl:Restriction;
      owl:onProperty gr:name;
      owl:hasValue "AirportShuttle"@en
    ]
  )
];
```

In this way, a reasoner can map to the appropriate `tao:LocationAmenity` sub-class an accommodation feature defined using `acco:value` and `gr:name` as prescribed in the Accommodation ontology specifications.

### D. TOURIST LOCATION HIERARCHY

*Tourist locations* are modelled, whenever possible, according to their respective *GeoNames* feature codes. This is done by declaring them as `owl:equivalentClass` to an anonymous class which is a restriction on the property `gn:featureCode` that must have an appropriate value from the *GeoName* feature codes list.<sup>80</sup> An example is given below:

```
tao:Zoo rdf:type owl:Class;
owl:equivalentClass [
  rdf:type owl:Restriction;
  owl:onProperty <http://www.geonames.org/ontology#featureCode>;
  owl:hasValue <http://www.geonames.org/ontology#S.ZOO>
];
rdfs:subClassOf <http://www.geonames.org/ontology#Feature>;
rdfs:label "Zoo"@en.
```

<sup>78</sup>Because there is not a defined taxonomy but a textual label is used to define a specific feature we can only try to guess the label most probably used.

<sup>79</sup>It is defined as “a structured value representing the feature of an accommodation as a property-value pair of varying degrees of formality”; see <http://ontologies.sti-innsbruck.at/acco/ns.html#AccommodationFeature>

<sup>80</sup>See <https://www.geonames.org/export/codes.html>

**APPENDIX C  
TRANSFORM THE DATA**

In this Appendix, we describe the processing steps for transforming the data shown in Figure 4.

**A. DATA EXTRACTION**

As the first step, we extracted the relevant data from the source data lake. The extraction process is performed using a SQL big data engine.<sup>81</sup> During this process, the data is also combined and arranged to be more easily processed in the following steps (e.g., unique ids are calculated, and nested columns are exploded). This produces the *Source data* assets collection which consists of:

- 1) *hospitality\_supply\_assets*: containing information about lodging facilities, accommodation, and offers.
- 2) *hospitality\_demand\_assets*: containing information about user reviews.

**B. DATA BREAK DOWN AND FILTER**

This second step organizes and structures the information produced in the previous step. Specifically, we need to:

- 1) break down the information so that we have a distinct asset for each semantic entity we want to model as triples (e.g., lodging facility, accommodation, offer, review);
- 2) apply a flat structure to the data, because some columns contain complex data structures as arrays or key/value structures;
- 3) separate text blobs from the other data preserving their relation to the semantic entity they refer to (e.g., the lodging facility description, the review content).

We can obtain the right structure using specific data pipelines that produce multiple assets out of a single one, flattening the data and filtering out unnecessary columns. This produces an unpacked version of the assets for each source:

- 1) *hospitality\_unpacked\_supply\_assets*: containing unpacked information about lodging facilities, accommodation, and offers.
- 2) *hospitality\_unpacked\_demand\_assets*: containing unpacked information about user reviews.

**C. DATA CLEANING**

Here we correct or remove corrupt or inaccurate records from the assets produced in the previous step. In particular, we need to drop duplicated records, remove special characters, normalize categorical fields, normalize date and numeric fields.

From *hospitality\_unpacked\_supply\_assets*, the Data Cleaning step produces:

- 1) *lodging\_assets* - containing all structured data relative to lodging facility entities (i.e., entities of type `tao:LodgingFacility`); for each lodging facility a unique ID is produced;

- 2) *lodging\_description\_assets* - containing all descriptions relative to a lodging facility (used to perform Named Entity Extraction and Linking);
- 3) *accommodation\_assets* - containing all structured data relative to accommodation entities (i.e., entities of type `tao:Accommodation`) in a lodging facility; for each accommodation, a unique ID is produced;
- 4) *offers\_assets* - containing all structured data relative to accommodation offers (i.e., entities of type `gr:Offering` that will be modelled as prescribed by the Accommodation Ontology); for each offer, a unique ID is produced;
- 5) *amenities\_assets* - containing all accommodation features (a.k.a. amenities) that are related to a lodging facility and/or to accommodation.

Instead, from *hospitality\_unpacked\_demand\_assets*, the Data Cleaning produces:

- 1) *reviews\_assets* - containing all structured data relative to user reviews about a lodging facility; for each review, a unique ID is produced;
- 2) *reviews\_content\_assets* - containing all text content for user reviews about a lodging facility (used to perform Named Entity Extraction and Linking);

**D. ONTOLOGY MAPPINGS**

At this stage, we identify and map the classes of the structured data to transform them into triples.

For instance, if a lodging business is represented as a record like:

Key	Value
hotel_id	9f40f613d308cf80
name	Chelsea BnB
structure_type	Bed and breakfast

after the ontology mapping step, a new field *lf\_class* (lodging facility class) is added with the “BedAndBreakfast” class name:

Key	Value
hotel_id	9f40f613d308cf80
name	Chelsea BnB
structure_type	Bed and breakfast
lf_class	BedAndBreakfast

Structured data include categorical columns that refer to concepts in the TAO ontology. In particular, there are three hierarchies in the ontology (See Appendix B for details) that we have to reconcile with categorical columns in the data:

- 1) *lodging facility types*: for each lodging table record we have a text field that contains the name of the lodging facility type; this field can be used to associate the correct `tao:LodgingFacility` subclass to the individual lodging facility the record is about;
- 2) *accommodation types*: for each accommodation table record we have a text field that contains the name of the

<sup>81</sup>Amazon Athena, see <https://aws.amazon.com/en/athena>

accommodation facility type; this field can be used to associate the correct `tao:Accommodation` subclass to the individual accommodation the record is about;

- 3) accommodation features (amenities) types: for each amenity table record we have an accommodation feature associated with a specific lodging facility (via an external key ID that refers to the lodging table). This field can be used to associate the correct `tao:LocationAmenity` subclass to the individual amenity the record is about.

To perform the reconciliation we use a heuristic process based on rules that can identify the most appropriate class to use to model an entity. The heuristic process uses lookup tables extracted from the ontology where we have each class associated with each of its labels. In this way, we leverage the ontology enrichment we already described in Section III-C3. The reconciliation is thus performed by adding the correct class name in a new column of the data table so that it can be used during the triple-creation phase. The ontology mapping step produces new types of assets that are part of the *Ontology mapped data* asset collection:

- 1) `classified_lodging_assets`;
- 2) `classified_accommodation_assets`;
- 3) `classified_amenities_assets`.

These assets will be fed into the triple creation process.

### E. LANGUAGE DETECTION

This step applies a language detection algorithm [61] to the text contained in the lodging description and reviews content tables. The detected language is used to enrich *lodging\_description\_assets* and *reviews\_content\_assets* with a new language column so that subsequent steps can process only English texts. The enriched assets are part of the *Language enriched data* asset collection.

### F. DBpedia ENTITY LINKING

To perform the Entity Linking task against DBpedia we have applied DBpedia Spotlight [62], [63] APIs<sup>82</sup> to the English text contained in the lodging description and reviews content tables. DBpedia Spotlight identifies and annotates entities based on the following pipeline process:

- Spotting: identifies possible entity mentions (surface forms) from the original input text.
- Candidate selection: selects the DBpedia resources that are candidate meanings for each surface form.
- Disambiguation: determines which candidate is the most likely resource for each surface form.
- Filtering: adjusts the annotation task based on the user requirements.

For the filtering step, we restricted the annotation scope to the following type of entities: `DBpedia:Activity`, `DBpedia:Food`, `DBpedia:Holiday`, `DBpedia:MeanOfTransportation`, `DBpedia:Place`, `Schema:Event`, `Schema:Place`. The result of the

DBpedia entity linking process produces two new types of assets which are part of the *DBpedia linked entities* asset collection:

- 1) `lodging_dbpedia_linked_assets` - containing a record for each DBpedia entity linked to a lodging facility identified by its unique ID;
- 2) `review_dbpedia_linked_assets` - containing a record for each DBpedia entity linked to a user review identified by its unique ID.

We used these assets in the triple-creation process.

### G. GeoNames ENTITY LINKING

This step performs an Entity Linking task against GeoNames so that places named in the lodging descriptions or the reviews are linked to the GeoNames corresponding entities.

To this end, we employed an open-source software called Mordecai<sup>83</sup> [38], a full-text geoparsing system that extracts place names from the text, resolves them to their correct entries in a gazetteer, and returns structured geographic information for the resolved place name. Mordecai is based on a language-agnostic architecture that uses word2vec [64] for inferring the correct country for a set of locations in a piece of text. As a gazetteer, it uses a custom-built Elasticsearch database populated with GeoNames data. Mordecai is integrated within the Spacy library.<sup>84</sup> Analogously to what is described in Appendix C-F for DBpedia, we used Mordecai to process all English text contained in the lodging description and review content tables. The result of the GeoNames entity linking process produces two new types of assets which are part of the *GeoNames linked entities* asset collection:

- 1) `lodging_geonames_linked_assets` - containing a record for each GeoNames entity linked to a lodging facility identified by its unique ID;
- 2) `review_geonames_linked_assets` - containing a record for each GeoNames entity linked to a user review identified by its unique ID.

We used these assets in the triple-creation process.

### H. IMPLEMENTATION STRATEGY

To support the data transformation described in the previous sections, we identified the following requirements for our technological architecture:

- Data-driven,
- Flexible and easily extensible,
- Scalable in a distributed computing environment,
- Easily manageable,
- Easily instrumented for lineage (a.k.a. provenance) metadata collection.

Following the requirements, the data computation is organised using the pipeline approach already described. This approach is optimal to create a distributed computation if the intermediate and final materializations are stored on a distributed file system. This is the same approach adopted by Apache Spark and other big data frameworks.

<sup>83</sup><https://github.com/openeventdata/mordecai>

<sup>84</sup>Only Spacy v2.x is supported at the moment.

<sup>82</sup><https://www.dbpedia.org/resources/spotlight/>



To manage the execution of a set of data pipelines, we used Dagster, an open-source orchestrator service. Dagster can be deployed on a single machine or a distributed environment like Kubernetes or AWS Elastic Container Service clusters. Thanks to this flexibility we started using a single machine to simplify the deployment process, without losing the opportunity to switch to a distributed architecture in the future. Dagster can also expose metadata about the execution of each pipeline and the produced assets, enabling our system to generate provenance information for the Knowledge Graph. The data transformation code is developed using Python Pandas library. We released the pipelines built on Dagster as an open-source resource for the paper.<sup>85</sup>

### I. PERFORMANCE ON A SINGLE SERVER

We used a single node with CPU AMD Ryzen™7 5800H, 32GB of RAM, 1TB SSD, and Ubuntu 20.04. With this setup the data transformation over the booking.com and Airbnb data was about 8 hours and 45 minutes, where the entity linking process took 7h 14m, language detection 1h and 26m, leaving all the other data extraction and transformation steps only 14 minutes of execution time. This is because entity linking is performed invoking DBpedia Spotlight public endpoints so that we could only apply a limited concurrency on the requests to the external web service to avoid server-side errors. This is the main limitation to scalability for the present implementation because the other data processing steps are very fast being executed using a big data query engine for the extraction (Amazon Athena) or using Python pandas with all data loaded in RAM. If a higher entity linking speed is needed it is possible to create a self-managed setup for DBpedia Spotlight as described in their website.<sup>86</sup> Regarding language detection, it can be optimised in a single-node setup using a multithreading approach similar to what has been implemented for entity linking and can also scale horizontally on multiple nodes because it only requires local CPU time.

We reduced the used disk space using Parquet files for tabular data. The total storage space was 3GB which can be reduced to 1.6GB if all triple files are compressed. To support storage scalability a distributed filesystem could be used as suggested in Appendix C.8.

## APPENDIX D TRIPLE STRUCTURE DETAILS FOR TKG

In this Appendix, we describe the structure of triples representing lodging facilities, accommodations, offers, and user reviews in the Tourism Knowledge Graph. We refer to Figure 5 in the following sections.

### A. LODGING FACILITY ENTITIES TRIPLE STRUCTURE

In Figure 5, we can steer our focus to observe triples modelling a lodging facility, which includes:

- 1) an address entity (`:address_1`), modelled as a `schema:PostalAddress` class that gives us great flexibility to define the facility position;
- 2) one or more accommodation features entities that are associated with the lodging facility using the `tao:feature` property; in our example, we have the node `:amenity_1` of type `tao:Parking`.<sup>87</sup>
- 3) an aggregated rating entity (`:agg_rating_1` in our example) that is used to model the overall user rating for the lodging facility (which is related to the ratings expressed by the single users' reviews) that specifies the vote in a normalised range from 0 to 1.

### B. ACCOMMODATION ENTITIES TRIPLE STRUCTURE

Accommodation is always related to a lodging facility, in compliance with the Accommodation ontology, and it includes:

- 1) its maximum and minimum occupancy capacity, using a `gr:QuantitativeValue` node (`:capacity_1` in our example);
- 2) its provision of beds, using an `acco:BedDetails` node (`:beds_1` in our example);
- 3) the type of accommodation<sup>88</sup> (using one of the TAO ontology classes like `tao:Room`).

### C. OFFER ENTITIES TRIPLE STRUCTURE

We describe a commercial offer for leasing out an accommodation leveraging GoodRelations. As shown in Figure 5 an offer can be expressed in terms of:

- 1) a node (`:quantity_1`) of type `gr:TypeAndQuantityNode` used to specify the number of days it is offered using `gr:amountOfThisGood` and `gr:hasUnitOfMeasurement` properties;
- 2) a node (`:price_spec_1`) of type `gr:UnitPriceSpecification` used to specify the price and currency for each day using the `gr:hasUnitOfMeasurement`, `gr:hasCurrency` and `gr:hasCurrencyValue` properties.

### D. USER REVIEWS TRIPLE STRUCTURE

A user review of the lodging facility is represented in TKG by two entities:

- 1) a node (`:review_1`) of type `schema:UserReview` with a `schema:dateCreated` property used to specify the review creation date;
- 2) a node (`:review_rating_1`) of type `tao:NormRating` that is used to specify the actual rating normalised to 1 (using `tao:normRatingValue`) property.

<sup>87</sup>In general the class of the amenity should be the most appropriate TAO ontology class among all the subclasses of `tao:LocationAmenity` as detected during the Ontology mapping step described in Appendix C-D.

<sup>88</sup>As detected during the Ontology mapping step described in Appendix C-D.

<sup>85</sup>See [https://github.com/linkalab/tkg/tree/main/kg\\_pipelines](https://github.com/linkalab/tkg/tree/main/kg_pipelines)

<sup>86</sup>See [http://dev.dbpedia.org/Dbpedia\\_Spotlight](http://dev.dbpedia.org/Dbpedia_Spotlight)

## APPENDIX E

## TAO EXTENSION

Here we report an example of the Python code we implemented on top of owlready2 for extending the TAO ontology with new classes:

```

from owlready2 import *
world = World()
tao_ontology = world.get_ontology("./ontologies
/tao_base.rdf").load()
tao = tao_ontology.get_namespace("http://purl.
org/tao/ns#")
with tao:
    class TouristLocation(schema.Place, gn.
    Feature):
        label = [locstr("Tourist_location",
        lang = "en")]
        comment = """A location is a point or
        area of interest from a tourist
        point of view, which a particular
        product or service is available, e.
        g. a museum, a beach, a bus stop, a
        gas station, or a ticket booth.
        The difference to gr:BusinessEntity
        is that the gr:BusinessEntity is
        the legal entity (e.g. a person or
        corporation) making the offer,
        while tao:Location is the store,
        office, or place. A chain
        restaurant will e.g. have one legal
        entity but multiple restaurant
        locations. Locations are
        characterized by an address or
        geographical position and a set of
        opening hour specifications for
        various days of the week."""
        altLabel = [locstr("Point_of_interest",
        lang = "en"), locstr("Area_of_
        interest", lang = "en"), locstr("
        Location", lang = "en")]
        seeAlso = gr.Location
    class TouristDestination(gn.Feature):
        label = [locstr("Tourist_destination",
        lang = "en")]
        comment = """A tourist destination. A
        TouristDestination is defined as a
        Place that contains, or is
        colocated with, one or more
        TouristLocation and LodgingFacility
        , often linked by a similar theme
        or interest to a particular tourist
        audience. The [UNWIO](http://www2.
        unwto.org/) defines Destination (
        main destination of a tourism trip)
        as the place visited that is
        central to the decision to take the
        trip."""
        equivalent_to = [schema.
        TouristDestination]
tao_ontology.save(file = "output_ontology/
tao_new.rdf", format = "rdxml")

```

**Listing 1.** Python snippet to extend the TAO ontology with new classes.

In the following code snippet we show an example of how we can process a CSV file that describes new classes to be integrated into the ontology. All data from the CSV file are loaded in a pandas dataframe and processed by a custom function (process\_entity function) that uses owlready2 to

```

from owlready2 import *
import pandas as pd
world = World()
tao_ontology = world.get_ontology("./ontologies
/tao_base.rdf").load()
tao = tao_ontology.get_namespace("http://purl.
org/tao/ns#")
df = pd.read_csv("./enrichment/
booking_facilities.csv")
df.apply(lambda r: process_entity(
    [tao_solo, acco], r['entity'], r['
    parent_class'], r['class'], r['type'], r
    ['is_amenity'],
    provenance = "Booking.com_features_lists_
    extraction.",
    comment_text = "Enriched_Booking.com_
    features_lists_extraction"), axis=1)
tao_ontology.save(file = "output_ontology/
tao_new.rdf", format = "rdxml")

```

**Listing 2.** Python snippet to extend the TAO ontology with new classes.

handle OWL class creation or modification. For more detail see the full source code.

## ACKNOWLEDGMENT

The authors thank Linkalab s.r.l., for supporting this research article by providing access to the Data Lake Turismo cloud infrastructure.

## REFERENCES

- [1] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P. N. Mendes, S. Hellmann, M. Morsey, P. van Kleef, S. Auer, and C. Bizer, "DBpedia—A large-scale, multilingual knowledge base extracted from Wikipedia," *Semantic Web*, vol. 6, no. 2, pp. 167–195, 2015.
- [2] D. Fensel, U. Şimşek, K. Angele, E. Huaman, E. Kärle, O. Panasiuk, I. Toma, J. Umbrich, and A. Wahler, *Knowledge Graphs: Methodology, Tools and Selected Use Cases*. Springer, 2020. [Online]. Available: <https://books.google.co.uk/books?id=1qnNDwAAQBAJ>
- [3] F. M. Suchanek, G. Kasneci, and G. Weikum, "YAGO: A core of semantic knowledge," in *Proc. 16th Int. Conf. World Wide Web (WWW)*, May 2007, pp. 697–706.
- [4] F. Erxleben, M. Günther, M. Krötzsch, J. Mendez, and D. Vrandečić, "Introducing Wikidata to the linked data web," in *The Semantic Web—ISWC 2014*, P. Mika, T. Tudorache, A. Bernstein, C. Welty, C. Knoblock, D. Vrandečić, P. Groth, N. Noy, K. Janowicz, and C. Goble, Eds. Cham, Switzerland: Springer, 2014, pp. 50–65.
- [5] R. Troncy, G. Rizzo, A. Jameson, O. Corcho, J. Plu, E. Palumbo, J. C. B. Hermida, A. Spirescu, K. D. Kuhn, C. Barbu, M. Rossi, I. Celino, R. Agarwal, C. Scanu, M. Valla, and T. Haaker, "3cixty: Building comprehensive knowledge bases for city exploration," *J. Web Semantics*, vols. 46–47, pp. 2–13, Oct. 2017.
- [6] D. Gazzè, A. L. Duca, A. Marchetti, and M. Tesconi, "An overview of the tourpedia linked dataset with a focus on relations discovery among places," in *Proc. ACM Int. Conf. Ser.*, Sep. 2015, pp. 157–160.
- [7] P. Calleja, F. Priyatna, N. Mihindukulasooriya, and M. Rico, "DBtravel: A tourism-oriented semantic graph," in *Current Trends in Web Engineering (Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11153, C. Pautasso, F. Sánchez-Figueroa, K. Systä, and J. M. M. Rodríguez, Eds. Cham, Switzerland: Springer, 2018, pp. 206–212.
- [8] M. Tenemaza, J. Limaico, and S. Luján-Mora, "Tourism recommender system based on natural language classifier," in *Advances in Artificial Intelligence, Software and Systems Engineering*, T. Z. Ahrum, W. Karwowski, and J. Kalra, Eds. Cham, Switzerland: Springer, 2021, pp. 230–235.

- [9] L. R. Ropesh and T. Bomatpalli, "A survey of travel recommender system," *Int. J. Comput. Sci. Eng.*, vol. 7, no. 3, pp. 356–362, Mar. 2019.
- [10] W. Zhang, H. Cao, F. Hao, L. Yang, M. Ahmad, and Y. Li, "The Chinese knowledge graph on domain-tourism," in *Proc. Int. Conf. Multimedia Ubiquitous Eng.*, in Lecture Notes in Electrical Engineering, vol. 590, Nov. 2020, pp. 20–27.
- [11] R. Alonso-Maturana, E. Alvarado-Cortes, S. López-Sola, M. O. Martínez-Losa, and P. Hermoso-González, "La Rioja turismo: The construction and exploitation of a queryable tourism knowledge graph," in *Proc. Int. Conf. Web Eng.*, in Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 11153, 2018, pp. 213–220.
- [12] S. Consoli, V. Presutti, D. R. Recupero, A. G. Nuzzolese, S. Peroni, M. Mongiovì, and A. Gangemi, "Producing linked data for smart cities: The case of catania," *Big Data Res.*, vol. 7, pp. 1–15, Mar. 2017.
- [13] M. S. Chaves and C. Trojahn, "Towards a multilingual ontology for ontology-driven content mining in social web sites," in *Proc. CEUR Workshop*, vol. 687, 2010, pp. 1–10.
- [14] K. I. Kotis, G. A. Vouros, and D. Spiliotopoulos, "Ontology engineering methodologies for the evolution of living and reused ontologies: Status, trends, findings and recommendations," *Knowl. Eng. Rev.*, vol. 35, p. e4, 2020. [Online]. Available: <https://www.cambridge.org/core/journals/knowledge-engineering-review/article/ontology-engineering-methodologies-for-the-evolution-of-living-and-reused-ontologies-status-trends-findings-and-recommendations/7A2D8D844EE0369C24967E156910AB50>
- [15] J. F. Sequeda, W. J. Briggs, D. P. Miranker, and W. P. Heideman, "A pay-as-you-go methodology to design and build enterprise knowledge graphs from relational databases," in *The Semantic Web—ISWC 2019*, C. Ghidini, O. Hartig, M. Maleshkova, V. Svátek, I. Cruz, A. Hogan, J. Song, M. Lefrançois, and F. Gandon, Eds. Cham, Switzerland: Springer, 2019, pp. 526–545.
- [16] G. Tamašauskaitė and P. Groth, "Defining a knowledge graph development process through a systematic review," *ACM Trans. Softw. Eng. Methodol.*, vol. 32, pp. 1–40, Feb. 2022.
- [17] D. Dessì, F. Osborne, D. R. Recupero, D. Buscaldi, and E. Motta, "Generating knowledge graphs by employing natural language processing and machine learning techniques within the scholarly domain," *Future Gener. Comput. Syst.*, vol. 116, pp. 253–264, Mar. 2021.
- [18] D. Dessì, F. Osborne, D. R. Recupero, D. Buscaldi, E. Motta, and H. Sack, "AI-KG: An automatically generated knowledge graph of artificial intelligence," in *Proc. 19th Int. Semantic Web Conf.*, in Lecture Notes in Computer Science, Athens, Greece, vol. 12507, J. Z. Pan, V. A. M. Tamma, C. d'Amato, K. Janowicz, B. Fu, A. Polleres, O. Seneviratne, and L. Kagal, Eds. Cham, Switzerland: Springer, Nov. 2020, pp. 127–143.
- [19] E. Kärle, U. Şimşek, O. Panasiuk, and D. Fensel, "Building an ecosystem for the tyrolean tourism knowledge graph," 2018, pp. 260–267, *arXiv:1805.05744*.
- [20] S. Consoli, A. Gangemi, A. G. Nuzzolese, S. Peroni, D. R. Recupero, and D. Spampinato, "Setting the course of emergency vehicle routing using geolinked open data for the municipality of Catania," in *The Semantic Web: ESWC 2014 Satellite Events*, in Lecture Notes in Computer Science, Anissaras, Greece, vol. 8798, V. Presutti, E. Blomqvist, R. Troncy, H. Sack, I. Papadakis, and A. Tordai, Eds. Cham, Switzerland: Springer, May 2014, pp. 42–53.
- [21] S. Consoli, A. Gangemi, A. G. Nuzzolese, S. Peroni, V. Presutti, D. R. Recupero, and D. Spampinato, "Geolinked open data for the municipality of Catania," in *Proc. 4th Int. Conf. Web Intell., Mining Semantics (WIMS)*, Thessaloniki, Greece, R. Akerkar, N. Bassiliades, J. Davies, and V. Ermolayev, Eds., Jun. 2014, pp. 58:1–58:8.
- [22] S. Consoli, A. Gangemi, A. G. Nuzzolese, S. Peroni, V. Presutti, D. R. Recupero, and D. Spampinato, "Towards emergency vehicle routing using geolinked open data: The case study of the municipality of Catania," in *Proc. 11st Workshop Semantic Sentiment Anal. (SSA), Workshop Social Media Linked Data Emergency Response (SMILE) Co-Located With 11th Eur. Semantic Web Conf. (ESWC)*, Crete, Greece, vol. 1329, A. Gangemi, H. Alani, M. Nissim, E. Cambria, D. R. Recupero, V. Lanfranchi, T. Kauppinen, Eds., May 2014, pp. 31–42.
- [23] D. Xiao, N. Wang, J. Yu, C. Zhang, and J. Wu, "A practice of tourism knowledge graph construction based on heterogeneous information," in *Proc. China Nat. Conf. Chin. Comput. Linguistics*, in Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 12522, 2020, pp. 159–173.
- [24] M. Atzeni and D. R. Recupero, "Multi-domain sentiment analysis with mimicked and polarized word embeddings for human–robot interaction," *Future Gener. Comput. Syst.*, vol. 110, pp. 984–999, Sep. 2020.
- [25] A. Dridi and D. R. Recupero, "Leveraging semantics for sentiment polarity detection in social media," *Int. J. Mach. Learn. Cybern.*, vol. 10, no. 8, pp. 2045–2055, Aug. 2019.
- [26] J. L. Martínez-Rodríguez, A. Hogan, and I. Lopez-Arevalo, "Information extraction meets the semantic web: A survey," *Semantic Web*, vol. 11, no. 2, pp. 255–335, Feb. 2020.
- [27] M. Grüninger, M. S. Fox, and M. Gruninger, "Methodology for the design and evaluation of ontologies," in *Proc. Int. Joint Conf. Artif. Intell. (IJCAI), Workshop Basic Ontol. Issues Knowl. Sharing*, 1995, pp. 1–10.
- [28] F. N. Noy and L. D. McGuinness, "Ontology development 101: A guide to creating your first ontology," Stanford Knowl. Syst. Lab., Stanford Univ., Stanford, CA, USA, Tech. Rep., 2001, p. 25.
- [29] O. Fodor and H. Werthner, "Harmonise: A step toward an interoperable E-tourism marketplace," *Int. J. Electron. Commerce*, vol. 9, no. 2, pp. 11–39, Jan. 2005.
- [30] S. Ou, V. Pekar, C. Orasan, C. Spurk, and M. Negri, "Development and alignment of a domain-specific ontology for question answering," in *Proc. 6th Int. Conf. Lang. Resour. Eval. (LREC)*, 2008, pp. 2221–2228.
- [31] S. Staab, C. Braun, I. Bruder, A. Düsterhöft, A. Heuer, M. Klettke, G. Neumann, B. Prager, J. Pretzel, H. P. Schnurr, R. Studer, H. Uszkoreit, and B. Wrenger, "GETESS—Searching the web exploiting German texts," in *Proc. Int. Workshop Cooperat. Inf. Agents*, in Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 1652, 1999, pp. 113–124.
- [32] R. Barta, C. Feilmayr, B. Pröll, C. Grün, and H. Werthner, "Covering the semantic space of tourism: An approach based on modularized ontologies," in *Proc. ACM Int. Conf. Ser.*, 2009, pp. 1–8.
- [33] R. Guha, D. Brickley, and S. Macbeth, "Schema.org: Evolution of structured data on the web," *Commun. ACM*, vol. 59, no. 2, pp. 44–51, Jan. 2016.
- [34] M. Hepp, "GoodRelations: An ontology for describing products and services offers on the web," in *Proc. Int. Conf. Knowl. Eng. Knowl. Manag.*, in Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 5268, 2008, pp. 329–346.
- [35] E. Kärle, U. Simsek, Z. Akbar, M. Hepp, and D. Fensel, "Extending the schema.org vocabulary for more expressive accommodation annotations," in *Information and Communication Technologies in Tourism 2017*, R. Schegg and B. Stangl, Eds. Cham, Switzerland: Springer, 2017, pp. 31–41.
- [36] M. S. Chaves, L. Freitas, and R. Vieira, "Hontology: A multilingual ontology for the accommodation sector in the tourism industry," in *Proc. Int. Conf. Knowl. Eng. Ontol. Develop. (KEOD)*, 2012, pp. 149–154.
- [37] P. N. Mendes, M. Jakob, A. García-Silva, and C. Bizer, "DBpedia spotlight: Shedding light on the web of documents," in *Proc. 7th Int. Conf. Semantic Syst.*, Sep. 2011, pp. 1–8.
- [38] A. Halterman, "Mordecari: Full text geoparsing and event geocoding," *J. Open Source Softw.*, vol. 2, no. 9, p. 91, Jan. 2017.
- [39] A. Delpuech, "OpenTapioca: Lightweight entity linking for Wikidata," 2019, *arXiv:1904.09131*.
- [40] A. Sakor, K. Singh, A. Patel, and M.-E. Vidal, "Falcon 2.0: An entity and relation linking tool over Wikidata," in *Proc. 29th ACM Int. Conf. Inf. Knowl. Manag.*, Oct. 2020, pp. 3141–3148.
- [41] C. Möller, J. Lehmann, and R. Usbeck, "Survey on English entity linking on Wikidata," 2021, *arXiv:2112.01989*.
- [42] J.-B. Lamy, "Owlready: Ontology-oriented programming in Python with automatic classification and high level constructs for biomedical ontologies," *Artif. Intell. Med.*, vol. 80, pp. 11–28, Jul. 2017.
- [43] L. Iannone, A. Rector, and R. Stevens, "Embedding knowledge patterns into OWL," in *The Semantic Web: Research and Applications*, L. Aroyo, P. Traverso, F. Ciravegna, P. Cimiano, T. Heath, E. Hyvönen, R. Mizoguchi, E. Oren, M. Sabou, and E. Simperl, Eds. Berlin, Germany: Springer, 2009, pp. 218–232.



- [44] G. M. Skjæveland, H. Forssell, W. J. Klüwer, P. D. Lupp, E. Thorstensen, and A. Waaler, "Pattern-based ontology design and instantiation with reasonable ontology templates," in *Proc. WOP@ISWC*, 2017, pp. 1–15.
- [45] W. P. Lord, "The semantic web takes wing: Programming ontologies with tawny-OWL," 2013, *arXiv:1303.0213*.
- [46] A. Dimou, M. V. Sande, P. Colpaert, R. Verborgh, E. Mannens, and R. Van De Walle, "RML: A generic language for integrated RDF mappings of heterogeneous data," in *Proc. CEUR Workshop*, vol. 1184, 2014, pp. 1–5.
- [47] A. Dimou, T. Nies, R. Verborgh, E. Mannens, and R. D. Walle, "Automated metadata generation for linked data generation and publishing workflows," in *Proc. 9th Workshop Linked Data Web*, vol. 1593, 2016, pp. 1–10.
- [48] P. Heyvaert, B. de Meester, A. Dimou, and R. Verborgh, "Declarative rules for linked data generation at your fingertips!" in *Proc. Eur. Semantic Web Conf.*, in Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 11155, Aug. 2018, pp. 213–217.
- [49] R. Verborgh, M. V. Sande, P. Colpaert, S. Coppens, E. Mannens, and R. D. Walle, "Web-scale querying through linked data fragments," in *Proc. CEUR Workshop*, vol. 1184, 2014.
- [50] R. Verborgh, O. Hartig, B. D. Meester, G. Haesendonck, L. D. Vocht, M. V. Sande, R. Cyganiak, P. Colpaert, E. Mannens, and R. V. D. Walle, "Querying datasets on the web with high availability," in *The Semantic Web—ISWC 2014*, P. Mika, T. Tudorache, A. Bernstein, C. Welty, C. Knoblock, D. Vrandečić, P. Groth, N. Noy, K. Janowicz, and C. Goble, Eds. Cham, Switzerland: Springer, 2014, pp. 180–196.
- [51] A. Gangemi, C. Catenacci, M. Ciaramita, and J. Lehmann, "Modelling ontology evaluation and validation," in *Proc. ESWC*, 2006, p. 15.
- [52] V. A. Carriero, A. Gangemi, M. L. Mancinelli, A. G. Nuzzolese, V. Presutti, and C. Veninata, "Pattern-based design applied to cultural heritage knowledge graphs," *Semantic Web*, vol. 12, no. 2, pp. 313–357, Jan. 2021.
- [53] E. Blomqvist, A. S. Sepour, and V. Presutti, "Ontology testing—Methodology and tool," in *Proc. Int. Conf. Knowl. Eng. Knowl. Manag.*, in Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 7603, Nov. 2020, pp. 216–226.
- [54] A. M. Orme, H. Tao, and L. H. Etkorn, "Coupling metrics for ontology-based system," *IEEE Softw.*, vol. 23, no. 2, pp. 102–108, Mar. 2006.
- [55] A. Bordes, N. Usunier, A. Garcia-Durán, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," in *Proc. 26th Int. Conf. Neural Inf. Process. Syst. (NIPS)*, vol. 2. Red Hook, NY, USA: Curran Associates, 2013, pp. 2787–2795.
- [56] M. Nickel, V. Tresp, and H.-P. Kriegel, "Factorizing YAGO: Scalable machine learning for linked data," in *Proc. 21st Int. Conf. World Wide Web (WWW)*. New York, NY, USA: Association for Computing Machinery, 2012, pp. 271–280.
- [57] T. Trouillon, J. Welbl, S. Riedel, E. Gaussier, and G. Bouchard, "Complex embeddings for simple link prediction," in *Proc. 33rd Int. Conf. Mach. Learn. (ICML)*, vol. 48, 2016, pp. 2071–2080.
- [58] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, "Graph neural networks: A review of methods and applications," *AI Open*, vol. 1, pp. 57–81, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2666651021000012>
- [59] A. Borrego, D. Ayala, I. Hernández, C. R. Rivero, and D. Ruiz, "CAFE: Knowledge graph completion using neighborhood-aware features," *Eng. Appl. Artif. Intell.*, vol. 103, Aug. 2021, Art. no. 104302.
- [60] C. Zhang, H. Yao, C. Huang, M. Jiang, Z. Li, and N. V. Chawla, "Few-shot knowledge graph completion," 2019, *arXiv:1911.11298*.
- [61] N. Shuyo. (2010). *Language Detection Library for Java*. [Online]. Available: <http://code.google.com/p/language-detection/>
- [62] N. P. Mendes, M. Jakob, A. García-Silva, and C. Bizer, "DBpedia spotlight: Shedding light on the web of documents," in *Proc. 7th Int. Conf. Semantic Syst. (I-Semantics)*, New York, NY, USA, 2011, pp. 1–8.
- [63] J. Daiber, M. Jakob, C. Hokamp, and P. N. Mendes, "Improving efficiency and accuracy in multilingual entity extraction," in *Proc. 9th Int. Conf. Semantic Syst.*, Sep. 2013, pp. 121–124.
- [64] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013, *arXiv:1301.3781*.



**ALESSANDRO CHESSA** received the Ph.D. degree in theoretical physics. He was an Adjunct Professor with LUISS. He is currently the CEO of Linkalab, a complex systems computational laboratory, and the Scientific Advisor with Eni Datalab. He has been a Research Associate with Boston University working on econophysics. His scientific research interests include applying quantum mechanics to the World Wide Web, the study of the social graphs of new communities on the Internet, and the data-driven journalism. Recently, he is studying the impact of artificial intelligence on human creativity.



**GIANNI FENU** received the Laurea degree in engineering from the University of Cagliari, Italy, in 1985. He joined the University of Cagliari, in 1988, where he has been the Director of bioinformatics and innovation and informatics services masters and is currently a Full Professor of computer science with the Department of Mathematics and Computer Science. He teaches courses of computer networks and information systems for first level degree in computer science students, and network architecture for specialized degree courses in informatics students. From 2008 to 2015, he was the Coordinator of the course of studies of computer science and informatics. He was scientifically responsible for Smart Cities Project of E-learning Ilearnv MIUR-UE (10 ME and six partners) from 2014 to 2017. He is currently involved in two regional projects Natura 2000 (L.R. 7/2007) and in the European Research M-Commerce and Development. He is responsible for the Bilateral Agreement with the Universidad Nacional de Tucuman, Argentina. He is the Delegate of Rector of the ELIOS Project (MIUR and 1.2 ME). He is a Council Member of the UnitelSardegna Consortium (the University of Cagliari and the University of Sassari) and the Director of the E-Learning Center, University of Cagliari. He is the Rector's Delegate for information and communication technologies, the Delegate of Rector in GARR-CRUI, and the President of the Faculty of Science. He is the author of more than 100 papers in international journals and conference proceeding.

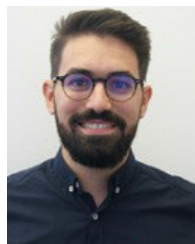


**ENRICO MOTTA** received the Laurea degree in computer science from the University of Pisa, Italy, and the Ph.D. degree in artificial intelligence from The Open University, U.K. From 2000 to 2007, he was the Director of the Knowledge Media Institute (KMi), The Open University, where he is currently a Professor in knowledge technologies. His research interests include the intersection of large-scale data integration and modeling, semantic and language technologies, intelligent systems, and human-computer interaction. Over the years, he has been leading KMi's contribution to numerous high-profile projects, receiving over €10.4M in external funding, since 2000, from a variety of institutional funding bodies and commercial organizations.





**FRANCESCO OSBORNE** is currently a Senior Research Fellow with the Knowledge Media Institute, The Open University, U.K., where he leads the Scholarly Data Mining Team. He is also an Assistant Professor with the University of Milano-Bicocca. His research covers artificial intelligence, information extraction, knowledge graphs, science of science, and semantic web. He has authored more than 100 peer-reviewed publications in top journals and conferences of these fields. He collaborates with major publishers, universities, and companies in the space of innovation for producing a variety of innovative services for supporting researchers, editors, and research politics makers. He released many well-adopted resources, such as the Computer Science Ontology and the Computer Science Knowledge Graph.



**ANGELO SALATINO** received the Ph.D. degree, studying methods for the early detection of research trends. He is currently a Research Associate with the Intelligence Systems and Data Science (ISDS) Group, Knowledge Media Institute (KM<sub>i</sub>), The Open University. In particular, his project aimed at identifying the emergence of new research topics at their embryonic stage. His research interests include semantic web, network science, and knowledge discovery technologies, with a focus on the structures and evolution of science.

**ANGELO SALATINO** received the Ph.D. degree, studying methods for the early detection of research trends. He is currently a Research Associate with the Intelligence Systems and Data Science (ISDS) Group, Knowledge Media Institute (KM<sub>i</sub>), The Open University. In particular, his project aimed at identifying the emergence of new research topics at their embryonic stage. His research interests include semantic web, network science, and knowledge discovery technologies,



**DIEGO REFORGIATO RECUPERO** received the Ph.D. degree in computer science from the University of Naples Federico II, Italy, in 2004. From 2005 to 2008, he was a Postdoctoral Researcher with the University of Maryland, College Park, MD, USA. He has been a Full Professor with the Department of Mathematics and Computer Science, University of Cagliari, Italy, since February 2022. He co-founded six companies within the ICT sector and is actively involved in European projects and research (with one of his companies he won more than 40 FP7 and H2020 projects). His current research interests include sentiment analysis, semantic web, natural language processing, human–robot interaction, financial technology, and smart grid. He is the author of more than 190 conference and journal papers in these research fields, with more than 2400 citations. He won different awards in his career (such as the Marie Curie International Reintegration Grant, the Marie Curie Innovative Training Network, the Best Researcher Award from the University of Catania, the Computer World Horizon Award, the Telecom Working Capital, the Startup Weekend, and the Best Paper Award).

**DIEGO REFORGIATO RECUPERO** received the Ph.D. degree in computer science from the University of Naples Federico II, Italy, in 2004. From 2005 to 2008, he was a Postdoctoral Researcher with the University of Maryland, College Park, MD, USA. He has been a Full Professor with the Department of Mathematics and Computer Science, University of Cagliari, Italy, since February 2022. He co-founded six companies within the ICT sector and is actively



**LUCA SECCHI** received the Laurea degree in electronic engineering from the University of Cagliari, Italy, where he is currently pursuing the Ph.D. degree. His research interests include knowledge graphs, natural language processing, big data, and semantic web. He has more than 20 years of experience as a Professional in the IT field both in the public and private sectors. He is one of the partners and the Chief Research and Development Officer of Linkalab, a private computational laboratory on complex systems.

**LUCA SECCHI** received the Laurea degree in electronic engineering from the University of Cagliari, Italy, where he is currently pursuing the Ph.D. degree. His research interests include knowledge graphs, natural language processing, big data, and semantic web. He has more than 20 years of experience as a Professional in the IT field both in the public and private sectors. He is one of the partners and the Chief Research and Development Officer of Linkalab, a private computational laboratory on complex systems.

...