

The euBusinessGraph Ontology: a Lightweight Ontology for Harmonizing Basic Company Information

Dumitru Roman ^{a,*}, Vladimir Alexiev ^b, Javier Paniagua ^c, Brian Elvesæter ^a,
Bjørn Marius von Zernichow ^a, Ahmet Soylu ^a, Boyan Simeonov ^b and Chris Taggart ^d

^a SINTEF AS, Norway

E-mail: {firstname.lastname}@sintef.no

^b Ontotext, Bulgaria

E-mail: {firstname.lastname}@ontotext.com

^c SpazioDati, Italy

E-mail: paniagua@spaziodati.eu

^d OpenCorporates, UK

E-mail: chris.taggart@opencorporates.com

Abstract. Company data, ranging from basic company information such as company name(s) and incorporation date to complex balance sheets and personal data about directors and shareholders, are the foundation that many data value chains depend upon in various sectors (e.g., business information, marketing and sales, etc.). Company data becomes a valuable asset when data is collected and integrated from a variety of sources, both authoritative (e.g., national business registers) and non-authoritative (e.g., company websites). Company data integration is however a difficult task primarily due to the heterogeneity and complexity of company data, and the lack of generally agreed upon semantic descriptions of the concepts in this domain. In this article, we introduce the euBusinessGraph ontology as a lightweight mechanism for harmonising company data for the purpose of aggregating, linking, provisioning and analysing basic company data. The article provides an overview of the related work, ontology scope, ontology development process, explanations of core concepts and relationships, and the implementation of the ontology. Furthermore, we present scenarios where the ontology was used, among others, for publishing company data (business knowledge graph) and for comparing data from various company data providers. The euBusinessGraph ontology serves as an asset not only for enabling various tasks related to company data but also on which various extensions can be built upon.

Keywords: Company data, Open data, Linked data, Ontology, Business knowledge graph

1. Introduction

Corporate information, including basic company information (e.g., name(s), incorporation data, registered addresses, ownership and related entities, etc.), financials (e.g., balance sheets, ratings, etc.) as well as contextual data (e.g., cadastral data on corporate properties, geo data, personal data about directors and shareholders, public tenders data, etc.) are the foundation that many data value chains depend upon in different sectors. The most evident examples of sectors are the business information sector, the

*Corresponding author. E-mail: dumitru.roman@sintef.no.

marketing and sales sector and the business publishing industry. At the same time, the use of company data is extremely significant in many other business sectors and societal activities including transparency and accountability [1].

Recently, a number of initiatives have been established to harmonise and increase the interoperability of corporate and financial data across national borders, including public initiatives such as the Global Legal Entity Identification System—GLEIS¹, Bloomberg’s open FIGI system for securities², as well as long-established proprietary initiatives such as the Dun & Bradstreet DUNS number³. Other notable initiatives include the European Business Register (EBR)⁴, which aims to federate several national business registers in order to offer a unique point of access, and BREX⁵, which “wraps” the EBR, extends its country coverage and offers a pricing model to access the underlying data. Additionally there are established and widespread adopted standardisation systems in the area of company financials (e.g., official deposited and public balance sheets data, which is in most cases exchanged in the XBRL format⁶). However, due to various reasons including technical, operational and organizational limitations, the systems and data sources mentioned above are mostly fragmented across borders, limited in scope and size⁷, and siloed within specific business communities with limited accessibility from outside their originating sectors. For example, register exchanges only offer access to official national registry data, not linked to any other contextual datasets (i.e., there is no obvious way of following a link from a company’s registered data to a tender it has won in another country), nor among themselves across countries (which means that there is no “machine-readable” and easy way to follow, for example, a shareholding relationship from an individual to companies in two different countries).

As a result, collecting and aggregating information about a business entity from several public sources (official and non-official ones, such as public tender registries, press mentions of companies and related entities, cadastral records, etc.), and especially across borders and languages is a tedious and very expensive task which renders many potential business models non-feasible. As a step in addressing this challenge, governments and other public bodies are increasingly publishing open data about firmographics and contextual databases, which reference companies. For example, the UK, Norway, France, and Denmark make the public records about companies available as open data, and other countries have different degrees of openness for their company registries⁸. Examples of contextual databases include the EU TED (Tenders Electronic Daily) public procurement notices⁹, gazette notices, Horizon 2020 project data¹⁰, and Structural Funds¹¹. Unfortunately, firmographics datasets are not yet fully harmonised and interoperable because data differs widely in semantics from one source to the other, and due to data formats ranging from UK’s five star Linked Data [2] to poorly accessible and poorly documented ones. Furthermore, contextual databases are not linked to the company registries and they still use different

¹<https://www.gleif.org>

²https://en.wikipedia.org/wiki/Financial_Instrument_Global_Identifier

³<http://www.dnb.com/duns-number.html>

⁴<http://www.ebr.org>

⁵<https://brex.io>

⁶<https://www.xbrl.org>

⁷Less than 1.6M companies worldwide were assigned a Legal Entity Identifier (LEI) number as of December 2019 (<https://search.gleif.org>) and these are only used in financial transactions of certain kind.

⁸<https://index.okfn.org/dataset/companies> and <http://registries.opencorporates.com>

⁹<https://ted.europa.eu>

¹⁰<https://data.europa.eu/euodp/en/data/dataset/cordisH2020projects>

¹¹<https://cohesiondata.ec.europa.eu>

identifier systems or, in some cases, no identifiers at all. Private businesses are also producers of valuable company-related data, which is seldom linked to the public sources mentioned above. For example, media publishers often reference businesses and legal entities by name (hence ambiguously) even within their digital publications (with the exception of traded company tickers, which are sometimes used by specialised financial publishers), because there isn't any widespread markup schema to annotate a digital reference to a company, nor a standardised way of accessing its information once it is unambiguously identified. As a result, it is extremely expensive, time consuming and error prone to find, interpret and reconcile these data from private sector sources. One of the immediate consequences is that the business information sector is very cost-inefficient in itself, which is reflected in a lack of transparency and efficiency of the markets. Nevertheless, the most relevant consequence in this context is that these inefficiencies severely harm digital innovation across sectors, which is often introduced by small and agile actors (e.g., startups, civil society organizations) who lack the capacity to invest time and resources in overcoming these problems.

In this article, we follow the established approach for harmonizing and integrating data based on ontologies (e.g., [3, 4]). In particular, we develop an ontology—the euBusinessGraph ontology—for harmonising and integrating basic company information. The ontology is meant to be used as a key mechanism for aggregating, linking, provisioning and analysing company-related data. The article provides an overview of the related work, ontology scope, ontology development process, explanations of core concepts and relationships, implementation of the ontology, and examples of scenarios where the ontology was used, among others, for publishing company data (business knowledge graph) and for comparing data from various company data providers.

The remainder of the article is organised as follows. Section 2 provides an overview of related work and ontologies relevant to company-related data. Section 3 describes the euBusinessGraph ontology development process, covering the scope, requirements, and the development approach. Section 4 gives an overview of the core concepts and relations in the euBusinessGraph ontology, together with details about the realization of the ontology. Section 5 provides examples of the usage of the ontology. Finally, Section 6 concludes this article and outlines possible future work.

2. Related Work

Several ontologies and data models were developed in the literature and have relevance to capturing the structure and complexity of company-related data. In what follows, we look specifically at works dealing with *basic* information about companies, covering organizational structures of companies, economical classifications of companies, company identification schemes, and locations of companies. This includes actual ontologies and vocabularies, and also several initiatives and data models relevant in the development of the euBusinessGraph ontology for basic company information.

The ontologies and vocabularies discussed in this section either insufficiently cover basic company information or are too complex due to many ontological commitments. Nevertheless, as we shall see below, relevant ontologies and data models were partly re-used and/or provided inspiration in the development of the euBusinessGraph ontology.

2.1. Organizational Structure

The W3C Organization ontology (ORG) [5] is a W3C recommendation since 2014. It aims to capture information about organizations and organizational structures, including governmental organizations. It

primarily captures organizational structure (e.g., sub-organizations and classification), reporting structure (e.g., roles and posts), location information (e.g., sites and buildings), and organizational history (e.g., merger and renaming). ORG is highly generic and designed as a core ontology, capturing general concepts and encouraging extensions for specific domains. It has been reused by other ontologies such as PPROC [6] in the procurement domain. The W3C Registered Organization Vocabulary (RegOrg)¹² is a profile of the W3C Organization ontology for describing organizations that have gained legal entity status through a formal registration process, typically in a national or regional register.

The e-Government Core Vocabularies [7] were developed in order to provide a minimum level of semantic interoperability for e-Government systems developed under the ISA program of the European Commission¹³. They include basic concepts about legal entities, locations, persons, public services, public organizations, and criterion to become eligible for public services and procurement. The Core Public Organization Vocabulary (CPOV) and the Core Business Vocabulary (CBV) are the most relevant vocabularies in our context. The CBV is published by W3C as a part of public working draft named RegOrg since 2013.

The Popolo Project defines data interchange formats and data models in the context of the Open Government initiative¹⁴. A set of concepts and relations are provided for capturing persons and organizations and the relationships between them (e.g., membership properties). A vocabulary for describing organizations is also provided. This vocabulary reuses terms from the ORG ontology and adds some new ones (e.g., other name, area, and contact detail).

The Application Profile of the Organization Ontology (ORG-AP-OP) was developed by the Publications Office of the European Union and supports its Whoiswho service¹⁵. It provides actual contact information for staff working at the European Institutions. It is concerned with people and the roles they play in the actual institutions. Similarly, in 2015, the ISA Programme of the EC initiated the development the Core Public Service Vocabulary and its Application Profile (CPSV-AP) [8]. However, it defines a number of terms closely related to CPOV, such as the administrative level, the type of organization, and its home page.

The Schema.org initiative [9] is spearheaded by the big four search engines, Google, Yahoo, Bing and Yandex, and is a collaborative effort to create, maintain, and promote schemas for structured data on the Internet. It is highly reusable since it makes few ontological commitments in order to cater to a truly global audience of millions of Web sites. Schema.org considers schemas as a set of types arranged in a hierarchy and associated with a set of properties. The core vocabulary is currently composed of 614 types and 902 properties. The “Organization” concept is among one of the commonly used types (among with, e.g., person, product, event) and models businesses (e.g., type, contact, etc.) and marketing aspects (e.g., logo, social profile, etc.).

2.2. Financial and Economic

The Financial Industry Business Ontology (FIBO) [10] is a joint effort of the Enterprise Data Management Council (EDMC) and the Object Management Group (OMG), aiming to go beyond a mere dictionary and capture the semantics of the business domain from a financial perspective. FIBO formalizes entities such as companies, directors, ownership and control relations, business registers, monetary

¹²<https://www.w3.org/TR/vocab-regorg>

¹³<https://ec.europa.eu/isa2>

¹⁴<http://www.popoloproject.com/specs>

¹⁵<http://whoiswho.europa.eu>

amounts, debts, obligations, contracts, and financial instruments. It is composed of a large number of smaller ontologies, with a modular perspective, each of which models a specific financial area [11]. The result is a large and very complex set of ontologies for the financial industry consisting of 11 core domains and 49 modules made available in more than 400 ontology files.

There are a number of classification vocabularies to specify the kind of economic activity such as International Standard Industrial Classification of All Economic Activities (ISIC) [12], which is a United Nations industry classification system, and European Commission's NACE [13], which is preferred in the context of European interoperability. The Wikipedia Business Entities¹⁶ provides a world-wide list for the types of business entities including a translation to English and approximate equivalents in the company law of English-speaking countries.

2.3. Company Identification and Location

The Global Legal Entity Identifier Foundation (GLEI) established a registration structure to issue Legal Entity Identifiers (LEI) to legal entities participating in financial transactions. The LEI structure is standardized as ISO 17442 [14]. LEI includes two code lists that are relevant in the context of basic company information, that is registration authorities list including 651 national official registers with their descriptions such as authority code, jurisdiction, and website; and, entity legal form code resolving variant names for each valid legal form within a jurisdiction to a single code per legal form.

The Business Registers Interconnection System (BRIS) interconnects business registers across Europe and provides a single (though limited) company search form¹⁷. The list of legal forms, list of national registers, and the pan-European company identifier (which is formed by register and company identifiers) are relevant for capturing basic company information.

With respect to capturing various forms of locations for companies, several initiatives are relevant. Eurostat has established a unified hierarchy of regions across the EU, EFTA and Candidate Countries. It consists of a nomenclature of Territorial Units for Statistics (NUTS) [15] and Local Administrative Units (LAU)¹⁸. NUTS and LAU are important geographic resources since a significant amount of open data is available that can support address data mapping (e.g., from postal code to NUTS) and use cases (e.g., hierarchical facets, distance calculations, spatial inclusion); and, NUTS and LAU provide a uniform hierarchy, whereas the administrative hierarchy varies greatly in different countries.

The ISA Programme Location Core Vocabulary [16] aims at describing any place in terms of its name, address or geometry through a minimum set of classes and properties. It is closely integrated with the Business (i.e., RegOrg) and Person Core Vocabularies of the EU ISA Programme.

GeoVocab.org¹⁹ provides vocabularies for geospatial modelling. This includes vocabularies NeoGeo Geometry Ontology for describing geographical regions and NeoGeo Spatial Ontology for describing topological relations between features.

Finally, GeoNames²⁰ provides a free geographical database covering all countries and containing over eleven million place names. It includes data elements such as administrative regions and settlements, and physical places.

¹⁶https://en.wikipedia.org/wiki/List_of_legal_entity_types_by_country

¹⁷<https://e-justice.europa.eu>

¹⁸<https://ec.europa.eu/eurostat/web/nuts/local-administrative-units>

¹⁹<http://geovocab.org>

²⁰<http://www.geonames.org>

2.4. Other relevant initiatives

In addition to well known initiatives such as FOAF²¹, Dublin Core²² and DBPedia²³, there are other ontologies, vocabularies and initiatives that are relevant in the context of modelling basic company information, including:

- ADMS ontology [17] describes various interoperability assets, including XML schemas, generic data models, code lists, taxonomies, dictionaries, vocabularies. ADMS is relevant in our context since we aggregate free company datasets from various company data providers.
- Vocabulary of Interlinked Datasets (VoID) [18] provides terms and patterns for describing RDF datasets and could be used in a variety of situations such as data discovery, cataloging and archiving of datasets.
- Simple Knowledge Organization System (SKOS) [19] offers a vocabulary for expressing the basic structure and content of concept schemes. This is essential for example for company classification (e.g., type and status).
- The IANA language code registry²⁴ uses ISO 639-1, 639-2 and 639-3 language codes (2 and 3-letter codes) and extends it with additional info (script, region of use, dialect). It can be consumed more easily from a Google sheet generated in Feb 2018.²⁵ Language tags are relevant in our context as some information (e.g., company names, street addresses) may be available in different languages.
- Person Core Vocabulary²⁶ aims at describing natural persons with a minimum set of classes and properties and is developed under the ISA Programme of the European Union. It is essential for representing people for example playing different roles in an organization.
- The Simple Event Model ontology (SEM) [20] is created for modelling events in a variety of domains and it is relevant for capturing different events in the lifetime of a company.

3. euBusinessGraph Ontology Development

In order to design the euBusinessGraph ontology, we applied common techniques recommended by well established ontology development methods [21, 22]. We used a bottom-up approach by identifying the scope and user group of the ontology, requirements, and ontological and non-ontological resources (some of which are referred to in Section 2).

One of the main resources used during the ontology development was company data that was provided by four company data providers and that needed to be harmonized before further processing. The data providers were OpenCorporates²⁷, SpazioDati²⁸, Brønnøysund Register Centre²⁹, and Ontotext³⁰. The

²¹<http://xmlns.com/foaf/spec>

²²<https://dublincore.org>

²³<https://wiki.dbpedia.org>

²⁴<https://www.iana.org/assignments/language-tags/language-tags.xml>

²⁵<https://docs.google.com/open?id=1M1yv9aBUmc-NyCJX69vOLUmH2uIglSwmDwgRgByI1AI>

²⁶<https://www.w3.org/ns/person>

²⁷<https://opencorporates.com>

²⁸<http://spaziodati.eu>

²⁹<http://www.brreg.no>

³⁰<https://www.ontotext.com>

data made available by the data providers originally came from both official sources (e.g., national and regional company registers) and unofficial sources (e.g., the corporate web, business-centric news aggregators and social networks). In the following we provide a brief description of the data provisioned by the four data providers:

- OpenCorporates provides core company data on over 180 million entities, obtained from more than 130 company registers around the world. The data is sourced only from official public sources and full provenance is provided. The depth of data varies from jurisdiction to jurisdiction, sometimes including directors and officers, industry codes, even occasionally shareholders and ultimate beneficial owners.
- SpazioDati integrates detailed up-to-date company and contact information on legal entities in Italy and the United Kingdom. Their dataset contains basic firmographics about more than 11 million business entities in both jurisdictions and information about 13 million directors and managers. Data comes from both authoritative sources (e.g., Registro imprese, the Italian Register of Companies and all the regional chambers of commerce) and non-authoritative sources (e.g., company websites, social media accounts, and business-centric news websites).
- Brønnøysund Register Centre (Brønnøysundregistrene) maintains the Norwegian Central Coordinating Register for Legal Entities (Enhetsregisteret)³¹—a database that contains information on all legal entities in Norway such as commercial enterprises and governmental agencies. It also includes business sole proprietorships, associations and other economic entities without registration duty that have chosen to join the register on a voluntary basis.
- Ontotext extracted data from the Bulgarian Trade Register. This register provides a centralized database whose purpose is to facilitate the start-up of businesses in Bulgaria, as well as to curb corruption practices.

These data sources were analyzed to determine the scope and requirements of the ontology. They cover official company information in Bulgaria, Norway, Italy and the United Kingdom, with additional unofficial information for the later two jurisdictions.

3.1. Scope and Requirements

After an analysis of the data provided by the different providers and the information available therein, we identified the major concerns that the ontology should address. Figure 1 provides an overview of the different types of information found during the data analysis, organized according to the type of entity being described (*Registered Organization* and *Officer*). In addition, the ontology needed to cover the description of dataset offerings by individual data providers (*Dataset*) and the description of identifier systems used to uniquely identify companies (*Identifier System*).

We identified target domains for our ontology, which primarily map to the business information sector, the marketing and sales sector, and the business publishing industry interested in new innovative data-driven products and services. Users working with data in these domains will benefit from a common representation that covers the types of information contributed by the different data providers. This common representation will also ease the task of data providers and aggregators who need to validate, transform and clean the data by providing a single ontology to target. The fact that there is a single ontology that provides a common representation will also benefit service developers who need to reference

³¹<https://data.brreg.no/enhetsregisteret/oppslag/enheter>

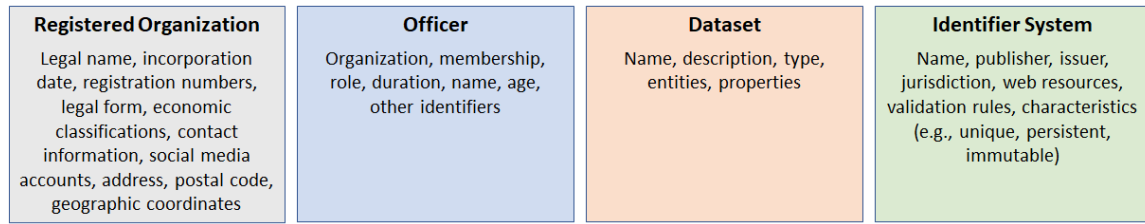


Fig. 1. Overview of the scope of the euBusinessGraph ontology.

company information to implement their services. To this end, the ontology has to capture the properties of the different identifiers that can be used to link the different entities being represented, providing machine readable descriptions for the identifier systems in use, including support for describing rules for validation and normalization of company and company-related identifiers.

Taking into account the needs of the intended users of the ontology and after the analysis of the data provided, we elicited the following requirements:

- (1) To capture the concept of a company, representing the different types or legal forms that companies can take, their jurisdictions and registration information, legal and alternative names official and secondary locations, prevalent economic activity, web keywords and social media accounts, among others;
- (2) To capture the concept of company officers, their roles and officerships, including temporal information to be able to represent these officerships through time;
- (3) To promote the use of the integrated data by reusing existing vocabularies as often as possible;
- (4) To provide machine-readable descriptions of the properties of the different systems of identifiers available to external applications and services, so algorithms can be developed to select and prioritise the most suitable identifiers for the task;
- (5) To provide validation and cleaning rules for identifiers to help their usage in unstructured data; and
- (6) To allow for extensibility, including vocabularies that describe additional properties of company and company-related entities that are not covered by the model but are available from the company data providers as unique or differentiating features.

Given the key requirements and the particular characteristics of the underlying datasets described at the beginning of this section, the ontology must be able to cover competency questions such as:

- (1) What companies are relevant to the search keywords “Opel” and “Car company”?
- (2) What kind of company identifier is the name “Opel”? What kind of identifier is “Opel Group GmbH”?
- (3) What are alternative names for the company registered as “Adam Opel GmbH”?
- (4) What is the company type of the company “Adam Opel GmbH”?
- (5) What jurisdiction does the company “Adam Opel GmbH” belong to?
- (6) Is “Bahnhofsplatz, 65423 Rüsselsheim am Mein” the address of the company “Adam Opel GmbH”?
- (7) Does the company “Adam Opel GmbH” have other locations?
- (8) Who are key managers of the company “Adam Opel GmbH”?

- (9) What is the Wikipedia page of the company “Adam Opel GmbH”?
- (10) What are the economic activities registered for the company “Adam Opel GmbH”?
- (11) Is the company “Adam Opel GmbH” publicly traded?
- (12) What additional information is available for the company “Adam Opel GmbH” from the different providers?

3.2. Ontology Development

The ontology development process was guided by the need to harmonize and integrate datasets with different sets of attributes, different representations for the same entity and in some cases close but not entirely similar semantics. Figure 2 depicts the four phases of the ontology development process in which we (a) gathered data from all company data providers that include natural language descriptions and example instances of each data attribute they provided, (b) analyzed attribute descriptions, refining them with additional notes describing their scope and using this information to group similar attributes, (c) analyzed identifiers and their identifier systems to produce machine readable descriptions of their properties, and (d) carried out manual reconciliation with the aim to reuse existing vocabularies.

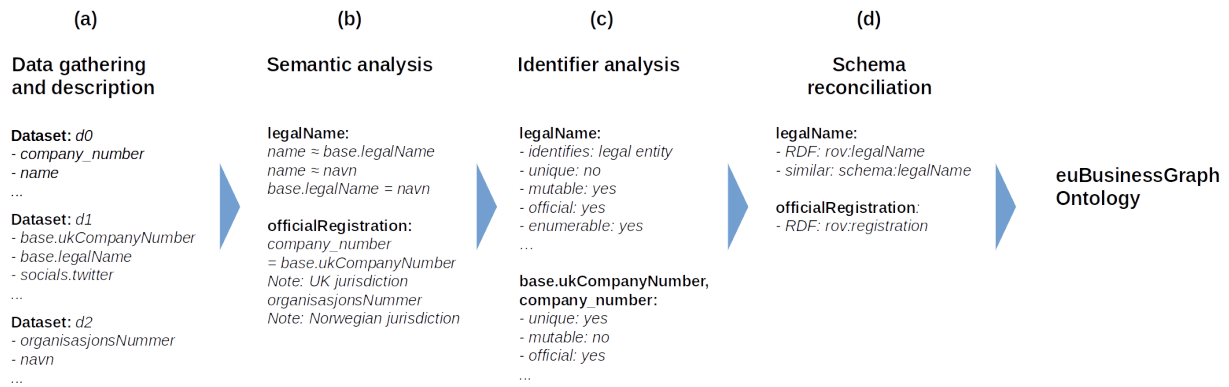


Fig. 2. Phases of the euBusinessGraph ontology development process.

There are differences in the types of information available from source to source (e.g., one dataset contains only official information from the national registers, while another integrates contact information parsed from company websites), differences in the way the same bit of information is represented by each provider (e.g., addresses as strings or as complex objects with separate attributes for street number, name and municipality) and differences in semantics for closely related concepts that may appear to be the same (e.g., information about officerships and their durations that contain references to possibly ambiguous officer names versus log entries that link person identification numbers to roles in different companies through time).

In the first phase of the ontology development process, as shown in Figure 2(a), each data provider provided a description of the dataset they shared. This data analysis focused on identifying the different attributes present and the way in which they were represented. Each attribute was described, adding notes and example uses that clarified the semantics as deemed appropriate. In this phase we already identified similar or even *same-as* candidates (e.g., *company_number*, *base.ukCompanyNumber*, *organisasjonsNummer* in Figure 2(a)). Moreover, each provider specified to which extent a particular attribute

was shared, in one of three modalities: (i) fully available, (ii) fully available to perform entity matching, but not available in any other case, and (iii) fully available for matching but available in reduced form for other purposes (e.g., address information without street numbers). Analyzing the descriptions provided in the previous phase, we identified a common subset shared by all contributed datasets. This common subset contained attributes that represented the same or very similar concepts in all datasets, which allowed us to group attributes from different providers accordingly (see similar attributes grouped under the *legalName* label across different providers in Figure 2(b)).

In the next phase, exemplified in Figure 2(c), we performed a different analysis to assess the suitability of each attribute to work as an identifier of the instance it described. The analysis contained a heterogeneous group of attributes with identifying characteristics: identifiers for geographical entities, legal entities, company headquarters and secondary sites, company websites, among others. Within the provided data, we found several ways to identify an instance in a group of similar instances (e.g., registration numbers and legal names are two different and useful ways to identify a company). Some identifiers are ambiguous in nature, such as company names, while others can be used to uniquely refer to a company, as is often the case with company registration numbers. The expectation is that the former will often be found in unstructured texts while the latter will be useful to annotate those unstructured texts to link to the corresponding instance being referred to. Some identifiers belong to official registers while others are self-issued and not centralized (e.g., websites). Some identifiers are subject to particular geographic jurisdictions (e.g., company registrations in local trade registers), or belong to special registers that attest that companies belong to a certain class (e.g., register of startup companies). In other cases, identifiers simply indicate the database in which the company information can be found (e.g., identification codes issued by data providers such as OpenCorporates, codes issued by other companies that aggregate company data such as Dun & Bradstreet), the website of a company or the various associated social network identifiers (e.g., a company's Facebook page or Twitter handle).

In light of the varied nature of the identifiers available, it was determined that the semantic model should also represent key aspects of the different identifier systems in use. These key aspects should encode expectations of the identifiers issued under each system and provide readily available rules to aid in validation and transformation of these identifiers. The expectations should help to determine the suitability of a particular indicator for common use cases that included publishing, reconciliation and matching within unstructured text. Additionally, the semantic model should provide links to information about issuing authorities and maintainers, revisions, databases and other resources.

In the last phase of the development process, as exemplified in Figure 2, we searched within existing vocabularies for all the concepts identified in the common subset aiming to reuse whenever possible. Examples of reuse from appropriate ontologies include W3C Org, RegOrg, Location, Person (not W3C), schema.org and ADMS datasets and identifiers.

Differences in the ways each provider decided to share the various attributes present in their datasets made it necessary to understand the scope of the ontology as early in the process as possible. In this way, it was possible to determine what to cover while having a clear path for extensibility.

4. Ontology Overview

The euBusinessGraph ontology is composed of 20 classes, 33 object properties, and 56 data properties that make it possible to represent basic company-related data. Figure 3 gives an overview of the ontology, depicting the main classes and their relationships (i.e., object properties). The ontology covers the following areas:

- (1) **Registered Organization:** The focal point of the ontology is companies that are registered as legal entities. Companies gain legal entity status by the act of registration. The class `RegisteredOrganization` is used to represent such a company. A company can have several `Sites`, for which the official registered site where legal papers can be served is captured by the object property `hasRegisteredSite`. A site can have an `Address`. Moreover, a company can have several different `Resources` associated in order to capture, e.g., `url` and `email` information.
- (2) **Identifier System:** A company can have several `Identifiers`, for which the official registration is captured by the object property `registration`. An identifier is part of an `IdentifierSystem`. Both the `Identifier` and the `IdentifierSystem` can have a creator of either a type `Person` or a type `Organization`. The `IdentifierSystem` also has additional `IdentifierWebResources` and `WebResources` information associated.
- (3) **Officer:** A company has associated officers, e.g., directors. The class `Membership` is used to associate officer data. It connects a `RegisteredOrganization` with a `Person` through a `Role`.
- (4) **Dataset:** Finally, in order to capture information about datasets that are offered by company data providers, we include the class `Dataset` that can have relevant `WebResources` information associated.

Further details about the `Registered Organization`, `Identifier System`, `Officer` and `Dataset` ontology areas, covering the full set of classes, object properties and data properties, are given in Sections 4.1, 4.2, 4.3 and 4.4 respectively. Moreover, Section 4.5 presents validation rules for the ontology.

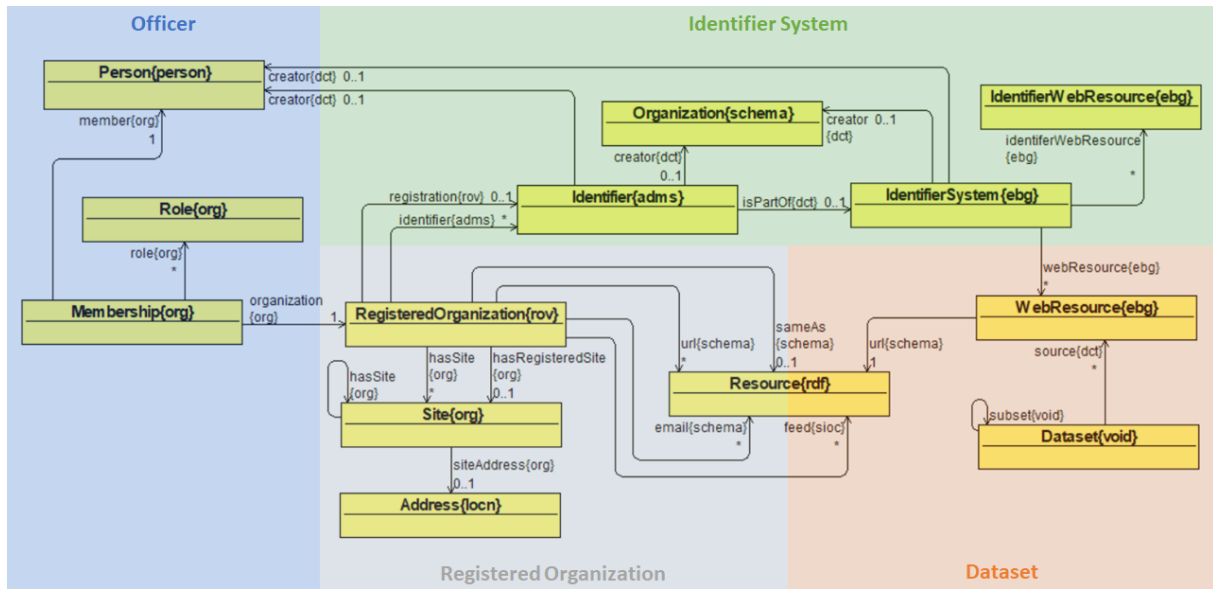


Fig. 3. euBusinessGraph ontology overview: Main classes and their relationships.

The class diagrams (depicting the ontology classes, object properties and data properties) and the object diagrams (depicting instances of the ontology classes and properties) in this section were cre-

ated using the Graphical Ontology Editor (OWLGrEd)³². An overview of the graphical elements in OWLGrEd for visualizing ontologies can be found in [23]. OWLGrEd expresses classes, namespaces, object properties, data properties and their data types, as well as cardinality in a visual manner. The notation `RegisteredOrganization{rov}` on a class refers to the term `RegisteredOrganization` defined in the namespace `rov`. The notation `legalName{rov}:string{xsd}[1..*]` on a data property refers to the term `legalName` defined in the namespace `rov`, that has the datatype `string` defined in the namespace `xsd`, and a cardinality of `1..*` (i.e., one or more). For simplicity, in the ontology descriptions in this section we omit namespaces if the context is given.

The ontology was defined as a Resource Description Framework (RDF) data model. We used the Terse RDF Triple Language (Turtle) syntax as the file format for the ontology. We reused classes and properties from existing ontologies and nomenclatures where appropriate in order to build our own ontology. Table 1 lists the prefixes and namespaces used in the euBusinessGraph ontology.

Table 1
Prefixes and namespaces used in the euBusinessGraph ontology

prefix	schema	namespace
adms	Asset Description Metadata Schema	http://www.w3.org/ns/adms#
dbo	DBpedia	http://dbpedia.org/ontology/
dct	DCMI Metadata Terms	http://purl.org/dc/terms/
ebg	The euBusinessGraph Ontology	http://data.businessgraph.io/ontology#
foaf	Friend of a Friend	http://xmlns.com/foaf/0.1/
locn	ISA Programme Location Core Vocabulary	http://www.w3.org/ns/locn#
ngeo	NeoGeo Geometry Ontology	http://geovocab.org/geometry#
nuts	EU NUTS classification as Linked Data	http://nuts.geovocab.org/id/
org	The Organization Ontology	http://www.w3.org/ns/org#
person	Core Person Vocabulary	http://www.w3.org/ns/person#
ramon	Reference And Management Of Nomenclatures	http://rdfdata.eionet.europa.eu/ramon/ontology/
rov	Registered Organization Vocabulary	http://www.w3.org/ns/regorg#
schema	Schema.org	http://schema.org/
sem	The Simple Event Model Ontology	http://semanticweb.cs.vu.nl/2009/11/sem/
skos	Simple Knowledge Organization System RDF Schema	http://www.w3.org/2004/02/skos/core#
time	Time Ontology in OWL	http://www.w3.org/2006/time#
void	Vocabulary of Interlinked Datasets	http://rdfs.org/ns/void#
xsd	XML Schema	http://www.w3.org/2001/XMLSchema#

The ontology uses `domainIncludes{schema}` and `rangeIncludes{schema}`, which are polymorphic and describe which properties are applicable to a class, rather than `domain{rdfs}` and `range{rdfs}`, which are monomorphic and prescribe what classes must be applied to each node using a property. We find that this enables more flexible reuse and combination of different ontologies.

Availability of the ontology and related materials. The ontology, datasets and examples described in this article are released as open source on the euBusinessGraph GitHub repository³³. The repository contains the ontology source file³⁴, the ontology reference documentation³⁵ generated with pyLODE³⁶,

³²<http://owlgred.lumii.lv>

³³<https://github.com/euBusinessGraph/eubg-data>

³⁴<https://raw.githubusercontent.com/euBusinessGraph/eubg-data/master/model/ebg-ontology.ttl>

³⁵<https://rawcdn.githack.com/euBusinessGraph/eubg-data/master/ontology/doc.html>

³⁶<https://github.com/RDFLib/pyLODE>

and the sources for the full example³⁷ used throughout this article. Additional materials related to the ontology include a semantic model with informative descriptions [24], a poster [25], and the ontology home page³⁸.

4.1. Registered Organization

Registered organizations are the main entities for which information is captured in the euBusinessGraph ontology. The ontology is not concerned with unregistered informal groups. Registered organizations gain legal entity status by the act of registration and are distinct from the broader concept of organizations, groups or, in some jurisdictions, sole traders. Figure 4 shows the classes and properties for representing core data about a registered organization. The class `RegisteredOrganization` contains names and other basic information about an organization such as `legalName` and `jurisdiction` (see Section 4.1.1), supports different types of classifications such as `orgActivity`, `orgType` and `orgStatus` (see Section 4.1.2). An organization can have several online resources associated such as email (see Section 4.1.3). A registered organization has a public site/address where legal papers can be served, and possible other sites/addresses. The sites/addresses are represented using the classes `Site` and `Address` (see Section 4.1.4). The object property `registration` denotes the identifier of a company. The identifier system is described in further details in Section 4.2.

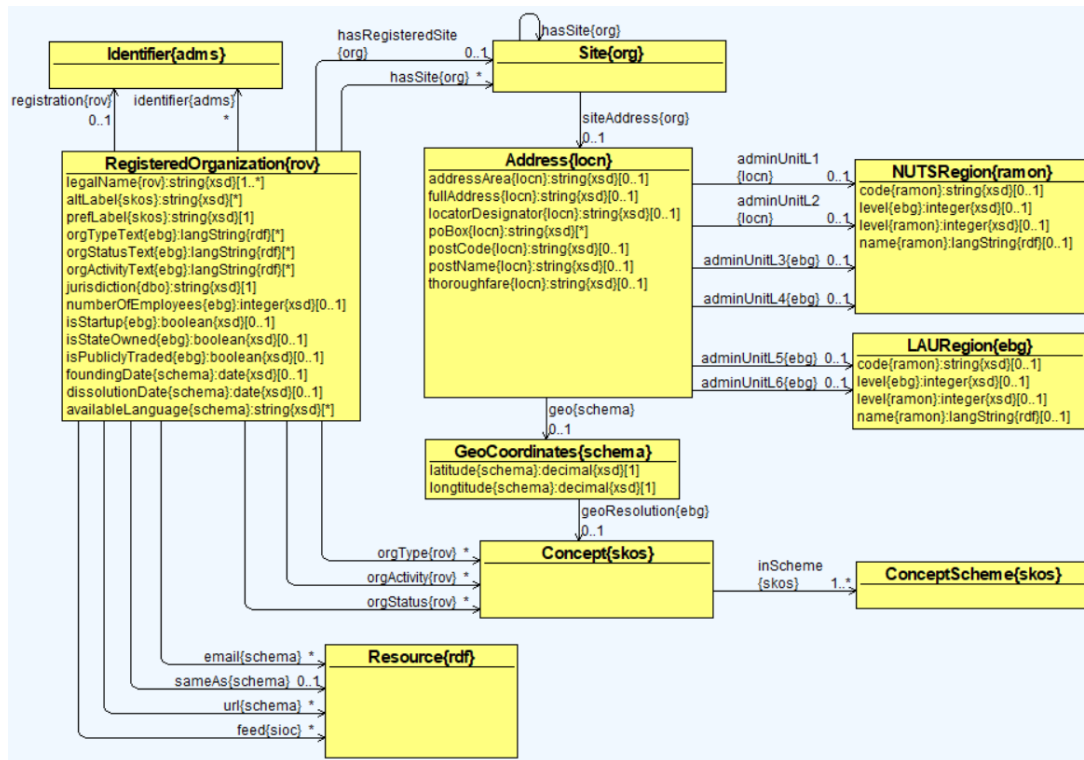


Fig. 4. Registered organization: Main classes and properties.

³⁷<https://github.com/euBusinessGraph/eubg-data/tree/master/example>

³⁸<https://www.eubusinessgraph.eu/eubusinessgraph-ontology-for-company-data>

4.1.1. Names and Other Basic Information

The ontology adopts two different name types for a registered organization, namely formal legal names and informal alternative names, e.g., a trading name. In addition we code a single name as the preferred name of the organization. The `RegisteredOrganization` class has the following data properties to record names:

- `legalName`: The legal name of the company, i.e., the official name of a company. A company may have more than one legal name, particularly in jurisdictions with more than one official language (e.g., Belgium). Some registries also treat a transliterated name as official, i.e., conversion of a legal name in one alphabet to another, e.g., from Russian to Latin.
- `altLabel`: Alternative names, e.g., an informal or popular name of the company. We also use this for former names.
- `prefLabel`: A single preferred name of a company.

The ontology defines the following data properties for capturing additional basic information about an organization:

- `jurisdiction`: Jurisdiction in which the company is registered.
- `numberOfEmployees`: The number of employees in the company.
- `isStartup`: Whether the company is a startup.
- `isStateOwned`: Whether this company is owned by the government, a government agency, municipality, city or other public entity. In many cases it is not possible to compute this attribute without access to a shareholder register, so it may be missing.
- `isPubliclyTraded`: Whether the company is publicly traded (listed at a stock exchange).
- `foundingDate`: Date when the company was created.
- `dissolutionDate`: Date the company was dissolved or removed from register.
- `availableLanguage`: Languages used by the company.

4.1.2. Classifications

Three types of classifications are defined in the ontology for representing the company type, company status and company activity. These are modelled as SKOS concept schemes. Alternatively, a free text field can be used. The `RegisteredOrganization` class has the following object properties and data properties to support the three classification types:

- `orgType`: Company type (legal form of the entity). There is no set of company types that is standardized across jurisdictions. Each jurisdiction will thus have a limited set of recognized company types. These should be expressed in a consistent manner in a SKOS concept scheme. Values are taken from the euBusinessGraph company type concept scheme³⁹ that covers jurisdictions NO, UK, IT and BG defined in collaboration with the data providers.
- `orgTypeText`: Company type (legal form of the entity) given in the form of free text.

³⁹<https://raw.githubusercontent.com/euBusinessGraph/eubg-data/master/data/lookups/EBG-company-type.ttl>

- `orgStatus`: The operational and/or legal registration status of the entity, e.g., whether a company is active or not. There is no globally accepted list of company states. For inactive, some providers look at hard evidence (i.e., that the company was deregistered), others at dissolution date in the past, or an extended period of inactivity (dormant). Because of this, a user cannot assume that active and inactive are opposites. A best practice for recording status levels is to use the relevant jurisdiction's terms and to encode these in a SKOS concept scheme. Values are taken from the euBusinessGraph company status concept scheme⁴⁰ that covers jurisdictions NO, GB, BG and statuses from data providers OpenCorporate, and SpazioDati and also from LEI. This concept scheme was defined in collaboration with the data providers.
- `orgStatusText`: Company status as it comes from a data provider (free text).
- `orgActivity`: Economic activity is recorded using a controlled vocabulary based on EC NACE 2. Values are taken from the euBusinessGraph NACE concept scheme⁴¹ which implements the NACE 2 vocabulary.
- `orgActivityText`: Economic activity of the organization (free text).

4.1.3. Online Resources

We represent commonly used electronic resources and channels (website, Wikipedia, email, news feed) as specific object properties of a company pointing to a `Resource` class:

- `email`: Email that is officially registered and with the same validity as certified mail.
- `sameAs`: Wikipedia page pertaining to the company.
- `url`: Website pertaining to the company or URL of a web resource.
- `feed`: URL of RSS/Atom feed pertaining to the company.

4.1.4. Sites and Addresses

Physical presence of companies is defined via addresses. We model `Address` in a structured way using a set of attributes such as country, macroregion, province, etc. Addresses may have geographic coordinates specified with a different resolution level. Least precise geographic coordinates are resolved at the level of a country, while most precise are geographical points that specify location up to a street and house number. We also enable data providers to provide full addresses in the form of a free text, which is essentially a string that combines all attributes together into a human-readable format. To provide RDF binding for the attributes, we considered two ontologies: Schema.org and the ISA Programme Location Core Vocabulary. We chose the latter as it has structured attributes, among which `fullAddress{locn}` that specifies the full address in a free-text form. However, to represent geographic coordinates, Schema.org was used as it provides a simpler way to model geographic coordinates via two properties (`latitude{schema}` and `longitude{schema}`).

We distinguish between registered, and other kinds of addresses. Many jurisdictions have the concept of registered address, i.e., the legal address where summons, subpoenas and other legal documents can be sent. An address is modelled using the `Site` and `Address` classes. A `Site` of a company is connected using the object property `hasSite`. A registered site is additionally connected using the object property `hasRegisteredSite`. A `Site` connects to an `Address` through the object property `siteAddress`.

The class `Address` represents a mailing or physical address of the company and has the following properties:

⁴⁰<https://github.com/euBusinessGraph/eubg-data/blob/master/data/lookups/EBG-company-status.ttl>

⁴¹<https://raw.githubusercontent.com/euBusinessGraph/eubg-data/master/data/NACE/nace.ttl>

- `fullAddress`: Full address, free text.
- `adminUnitL1`: Country of the address.
- `adminUnitL2`: NUTS1 region of the address.
- `adminUnitL3`: NUTS2 region of the address.
- `adminUnitL4`: NUTS3 region of the address.
- `adminUnitL5`: LAU1 region of the address. Some countries (e.g., Bulgaria) use both LAU1 and LAU2 levels. Others (e.g., Italy) use only LAU2.
- `adminUnitL6`: LAU2 region of the address.
- `postName`: Locality/city/settlement of the address, free text.
- `addressArea`: Part of a city, village or neighbourhood.
- `thoroughfare`: Street name (and optionally number).
- `locatorDesignator`: Street number and/or building name.
- `postcode`: Postal code of the address.
- `poBox`: Some addresses are associated with a PO box instead of a street address.

NUTS values are assigned using the EU NUTS classification as Linked Data (NUTS-RDF) datasets⁴². The NUTS-RDF datasets cover 34 European countries and use the `NUTSRegion` class to represent the NUTS regions. In order to represent the lower-level LAU regions we introduced the `LAURegion` class and created our own set of LAU-RDF datasets⁴³ covering 32 jurisdictions (including all of the EU and EEA), 26 languages, and both LAU territorial levels (`lau4`, `lau5`). LAU-RDF datasets were created from the official Eurostat Excel spreadsheet for 2016⁴⁴ for EU, and our own research on some other jurisdictions.

4.1.5. Example

Figure 5 is an object diagram depicting how the ontology is used to represent company data about the legal entity `OpenCorporates`. Each object (depicted as a green rectangle) is an instance of a class defined in the ontology. The objects have data properties according to the class definitions. The data properties are assigned values depicted using the notation `data property = "value"`. Some properties are mandatory (multiplicity of 1..*) whereas others are optional (cardinality of 0..* or *). Not all information about a company is available from a data provider. Thus an object will only contain the data properties that we are able to retrieve from the data provider. This may vary greatly from data provider to data provider, and from jurisdiction to jurisdiction.

Another example showing company data about the legal entity `SpazioDati` can be found in Section 5.1 (see Figure 15), where information about mapping of data from a data provider to the ontology is also discussed.

4.2. Identifier System

Mechanisms to identify companies in various data sources are essential in integration of data about companies across data sources. A proper understanding of what kind of systems of identifiers can be used for companies is thus necessary in this context. We analyzed various types of identifiers commonly used for companies and collected various properties of the systems they are part of. We modelled identifiers and identifier systems explicitly in the ontology as shown in Figure 6.

⁴²<http://nuts.geovocab.org>

⁴³<https://github.com/euBusinessGraph/eubg-data/tree/master/data/LAU/rdf>

⁴⁴https://ec.europa.eu/eurostat/documents/345175/501971/EU-28_LAU_2016

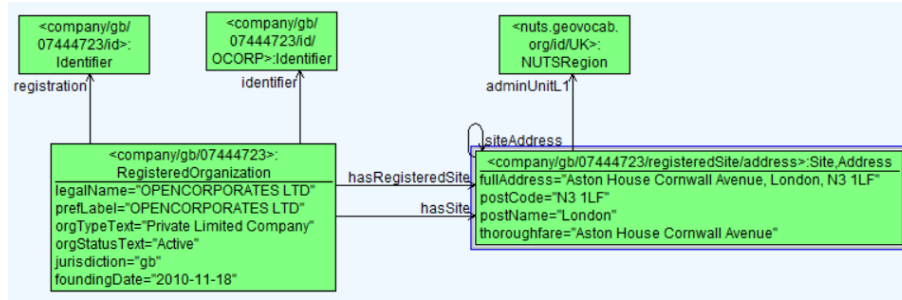


Fig. 5. Example of company representation for OpenCorporates.

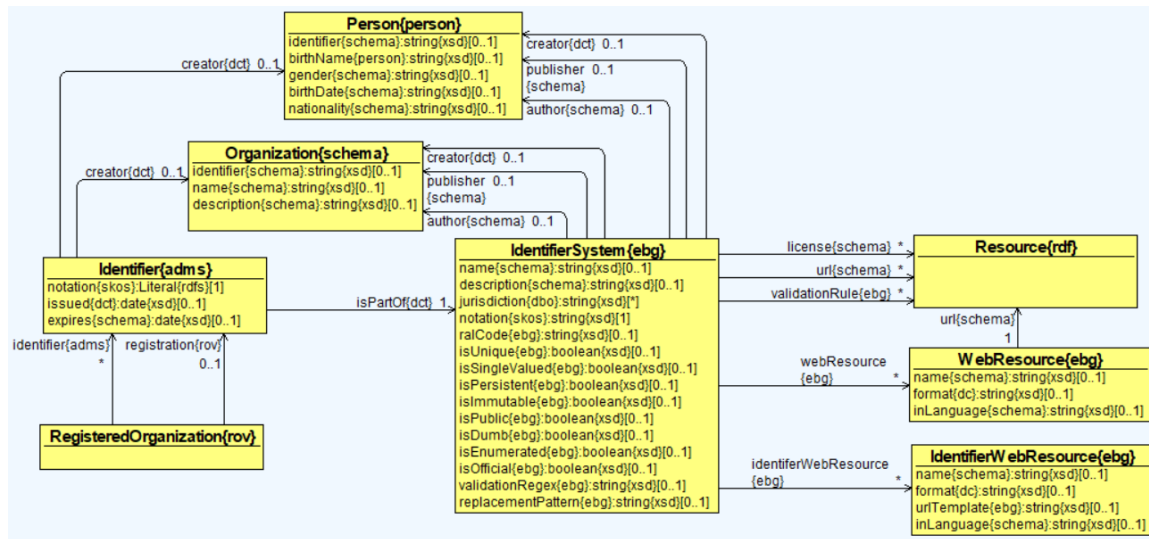


Fig. 6. Classes, object properties and data properties for representing identifier systems and identifiers.

A `RegisteredOrganization` can have several `Identifiers` issued by different issuers for different purposes. This is modelled by having each company identifier belong to an `IdentifierSystem` (see Section 4.2.1). In this way, we can differentiate between an “official registration” in official business registers and “alternative registrations” in other kinds of registers. While they have the same nature, only the former can be used to uniquely identify a company in each jurisdiction, and to confirm existence of the company as a legal entity in this jurisdiction. Other registrations may not be unique or persistent. The ontology models the different cases through properties that describe the lifecycle of each identifier issued and by encoding a series of characteristics of the identifier system to which the identifier belongs (see Section 4.2.2). Additionally, we model Web resources (see Section 4.2.3) that are frequently found for identifier systems such as search endpoints, templates for building identifier URLs (through which company information can be reached) and other resources that describe the system’s rules. Finally, the model captures the representation of different agents (see Section 4.2.4) that are in charge of setting and maintaining rules, issuing identifiers and publishing identifier databases.

4.2.1. Identifier and Identifier System

The `Identifier` class represents a company identifier. It has the following object and data properties:

- `isPartOf`: System the identifier is a part of.
- `creator`: The issuer of the identifier. In many countries there is a single registry although in others, such as Spain and Germany, multiple registries exist. If the system has an issuer, in most cases the identifier issuer will coincide with that issuer.
- `notation`: Literal value of the identifier.
- `issued`: Date when the identifier was issued.
- `expires`: Date when the identifier expires.

The `IdentifierSystem` class represents a system managed by a publisher (e.g., a register or agency) that is used to issue identifiers to companies. Many registers keep several identifier systems. There can be three different types of agents related to a system. This is modelled using three different object properties:

- `author`: The author who is in charge of specifying the rules and organization of the system.
- `creator`: The issuer who issues identifiers and then keeps them in a database (register).
- `publisher`: The publisher who publishes the identifier database (register) in some form.

4.2.2. Identifier System Properties and Characteristics

Identifier systems have some basic properties:

- `name`: Name of the identifier system.
- `description`: Description of the identifier system.
- `jurisdiction`: Jurisdiction to which the identifier system applies.
- `notation`: Short mnemonic code for the identifier system, used in its URL. Also used in identifier URLs that are part of the system. Issued locally by euBusinessGraph. For identifier systems published by the sole or preferred official register in a jurisdiction, we use the jurisdiction code (e.g., “BG”, “GB”). For others, if the identifier system has no explicit name, we use a short mnemonic code of the publisher: upper-case for company registers (e.g., “OCORP” for Open-Corporates, “SDATI” for SpazioDati, “BRC” for Brønnøysund Register Centre, “RAL”, “EU”, “BRIS”), mixed-case for social network registers (e.g., “Twitter”, “Facebook”).
- `ralCode`: GLEI RAL code for the identifier system.
- `url`: Various websites of the identifier system and/or its associated issuer and register, e.g., home page, search, download.
- `license`: License that applies to the system.
- `webResource`: Web resource(s) associated with an identifier system.
- `identifierWebResource`: Identifier Web resource(s) associated with an identifier system.

Identifier systems have some boolean characteristics (flags) that represent expectations about their identifiers. Some systems have exceptions, i.e., identifiers that don’t satisfy the expectations. Each flag is set to “true” in the desirable (positive) case. We strive to provide all flags for each system, but in some cases the flag could be omitted (e.g., if there is not enough information):

- `isUnique`: Whether each identifier in the system relates to only one entity.
- `isSingleValued`: Whether each entity has only one identifier in the system.
- `isPersistent`: Whether identifiers can be removed from the register (e.g., when a company is dissolved).

- `isImmutable`: Whether identifiers can change.
- `isPublic`: Whether identifiers from the system are available for public use: consulting, search or download.
- `isPublic`: Whether identifiers from the system are available for public use: consulting, search or download.
- `isDumb`: “Intelligent” or “smart” identifiers contain built-in “intelligence” (semantic information) embedded in the identifier. This is increasingly considered bad practice, since when the attributes change the identifier must also change, making it unreliable, particularly as a foreign key. “Dumb” identifiers on the other hand contain no intelligence and will not change.
- `isEnumerated`: Whether the system has an issuer, and issued identifiers are kept in a database (register).
- `isOfficial`: Whether the system is considered the official one in all jurisdictions in which it applies.

Identifier systems are associated with some properties that can be useful for identifier validation:

- `validationRule`: URL providing human or machine-readable rule(s) for validating identifiers in the system.
- `validationRegex`: Regular expression for validating identifier values of that system.
- `replacementPattern`: Pattern to use together with the `validationRegex` to normalize identifier values by removing optional decorations.

4.2.3. Web Resources

A Web resource is a URL complemented with a MIME type to specify what the URL is about. These web resources are used for identifier systems (e.g., to provide the search or download URL) and per-company, as a URL template in which to substitute the identifier value. There can be several MIME types because some URLs return various resource types using content negotiation. The class `WebResource` has the following object and data properties:

- `url`: URL of the Web resource.
- `name`: Name or short (generic) description of the resource.
- `format`: MIME type(s) of the resource. If several are provided, the server must provide all these resource types using content negotiation.
- `inLanguage`: Language of the Web resource.

The class `IdentifierWebResource` has the mandatory data property `urlTemplate` in addition to the three data properties defined for `WebResource` (i.e., excluding `url`). The property `urlTemplate` specifies a template that can be used uniformly to build URLs for all identifiers in the system. The template value can have placeholders that should be interpreted as follows:

- If it has a placeholder `{ }`, substitute the identifier value there.
- If it has placeholders like `$1`, `$2`, ..., substitute the groups extracted by the `validationRegex` of the `IdentifierSystem`.

4.2.4. Agents

We represent an agent using either a `Person` or `Organization` class, depending on the type of agent. For both types, we define the `identifier` data property which can be assigned a textual identifier or a URL value. For `Organization`, we additionally assign values to the data properties `name` and `description`. For `Person`, we introduce a set of data properties (see Section 4.3 for further details).

4.2.5. Example

An example of an identifier system is shown in Figure 7, illustrating the ATOKA identifier system that was created by SpazioDati. Full representation of all the Italian identifier systems (i.e., ATOKA, REA, Tax and VAT) referenced by the company SpazioDati in Figure 15 are available in RDF-format on GitHub⁴⁵.

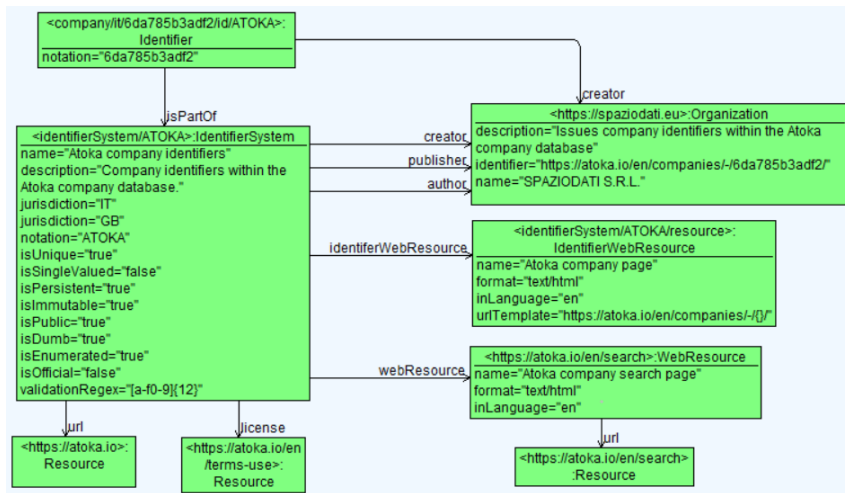


Fig. 7. Example of representing the ATOKA identifier system created by SpazioDati.

Another example of identifier systems is shown in Figure 8, illustrating the OpenCorporates identifier system for which OpenCorporates is the publisher and the official UK identifier system for which Companies House is the publisher.

4.3. Officer

We use the membership model⁴⁶ of the W3C Organization Ontology in a straightforward way to represent officer data. An officer is represented using a `Person` class for which the properties `identifier` and `birthName` are mandatory. The identifier may come from official registries or be derived from these. Additionally, other properties may be present such as `gender`, `birthDate` and `nationality`.

An officer is a natural person (as opposed to a legal person) that has a high-level management role in a company (e.g., the CEO, treasurer and chief financial officer). Despite their high status, they typically serve at the will of the company directors, who can fire or replace them. Officers can also be shareholders and directors but don't necessarily have to be. They have the authority to act on behalf of the corporation, including contract authority.

⁴⁵<https://github.com/euBusinessGraph/eubg-data/tree/master/example>

⁴⁶<https://www.w3.org/TR/vocab-org/#membership-roles-posts-and-reporting>

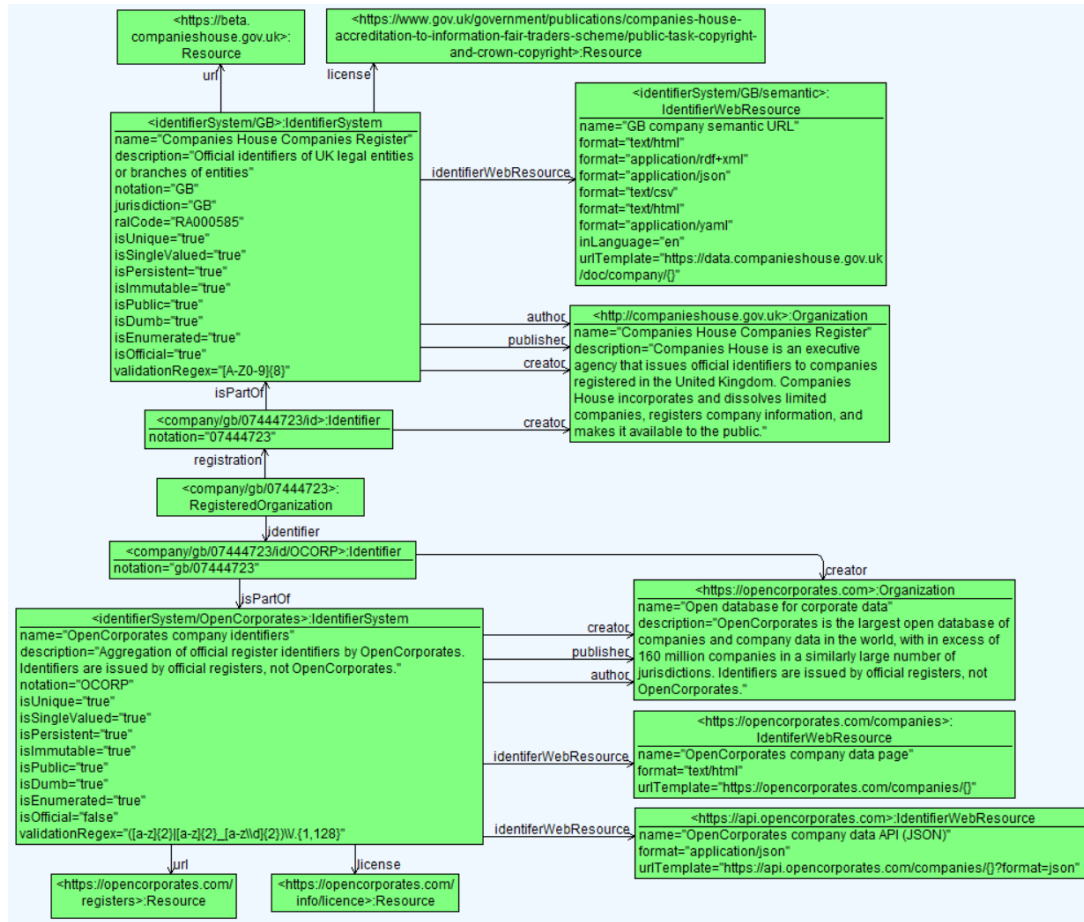


Fig. 8. Example of representing the OpenCorporates identifier system published by OpenCorporates.

A *Membership* describes the relation between an officer and the company in which they occupy a position. The *Role* defines the position the officer fulfills according to the membership. Ideally, the roles should be defined according to a SKOS concept scheme. We have not defined a global set of officer roles as this may vary per jurisdiction and/or provider. Thus, we also introduced the data property *rolePositionText* in the *Membership* class in order to capture the role as free text.

The membership interval is defined by the *memberDuring* object property that points to an *Interval*. The interval has a beginning and an end date. For open intervals only the beginning is mandatory. These dates are defined by the class *Instant* which has the data property *inXSDDate*.

4.3.1. Example

An example of the CEO role using SKOS concepts defined by the Atoka *IdentifierSystem* for the company SpazioDati is shown in Figure 10.

An example of officer roles using the free text data property *rolePositionText* for the company OpenCorporates is shown in Figure 11.

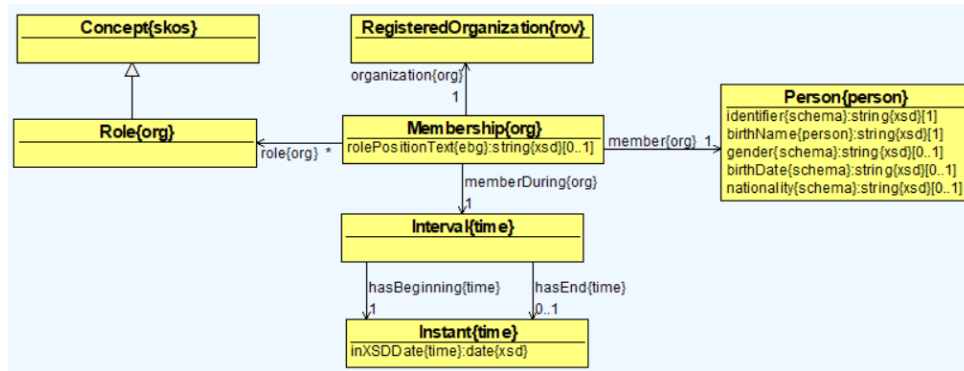


Fig. 9. Classes, object properties and data properties for representing officers.

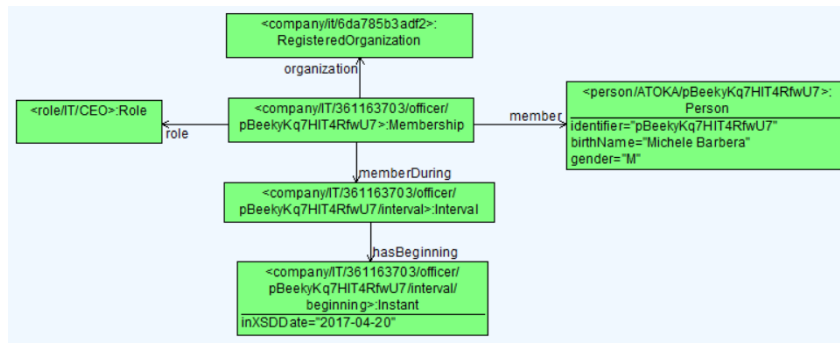


Fig. 10. Example of officer representation for the company SpazioDati.

4.4. Dataset

Data consumers need to know how many companies are included in a data provider dataset, from which jurisdictions, and what depth of data is included (e.g., which properties, addresses with what geo resolution, etc.). We thus need to express both metadata about the dataset itself, and fine-grained statistics about the content of a dataset, e.g.,:

- Publisher, source, last modified, license, home page, download distribution, etc.
- Subsets of data by kind of entity (e.g., companies vs. addresses), field coverage (which fields are included in which subsets), and entity characteristics (e.g., Italian companies, startups, startups in Italy).
- Count of entities in a dataset or subset.

After an analysis of various dataset description ontologies, we decided on using VOID with some extensions (see Figure 12). VOID describes RDF datasets in terms of entities (i.e., number of triples), property (i.e., used to list the properties available in the dataset), etc. The Dataset has a void:subset relation that is used to describe a dataset polyhierarchy. For each data provider we can capture their full dataset and the respective subsets. For each dataset the dct:publisher, dct:type and dct:license have to be captured.

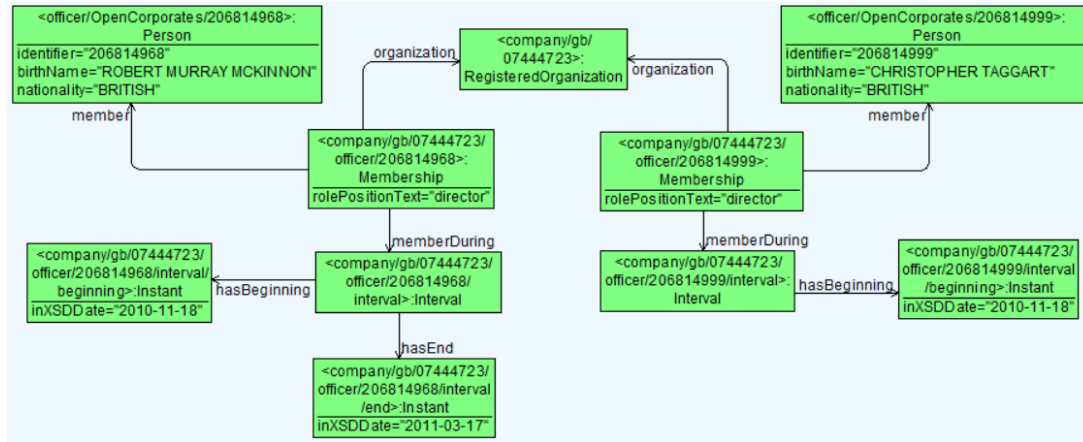


Fig. 11. Example of officer representation for the company OpenCorporates.

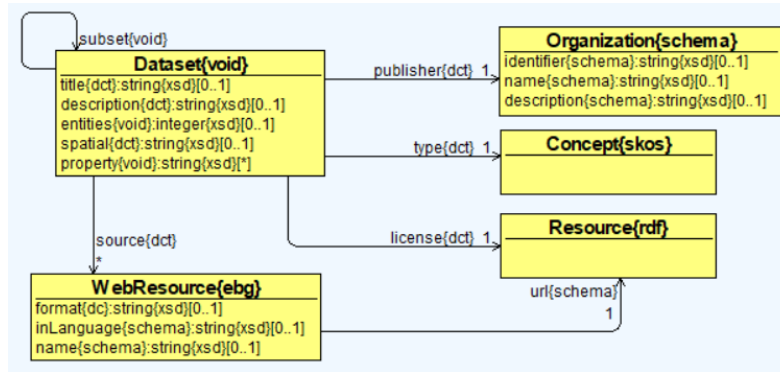


Fig. 12. Classes, object properties and data properties for representing datasets.

4.4.1. Example

Figure 13 shows an example of the datasets provided by SpazioDati. The main dataset `<dataset/SDATI>` consists of two subsets, namely `<dataset/SDATI/IT>` and `<dataset/SDATI/GB>`. For each subset we specify the number of entities and the properties that are available.

4.5. Validation Rules

In order to ensure that data can be correctly published according to the ontology, we devised a set of data validation rules that are associated with the ontology. The types of validations rules considered are as follows:

- **Data completeness:** Specifies that a given set of business attributes must be present (e.g., attribute `legalName` must be available).
- **Accuracy** Describes that data values must be correct (e.g., values of attribute `jurisdiction` must be included in the list of recognized nations available on Wikipedia⁴⁷).

⁴⁷https://en.wikipedia.org/wiki/List_of_sovereign_states

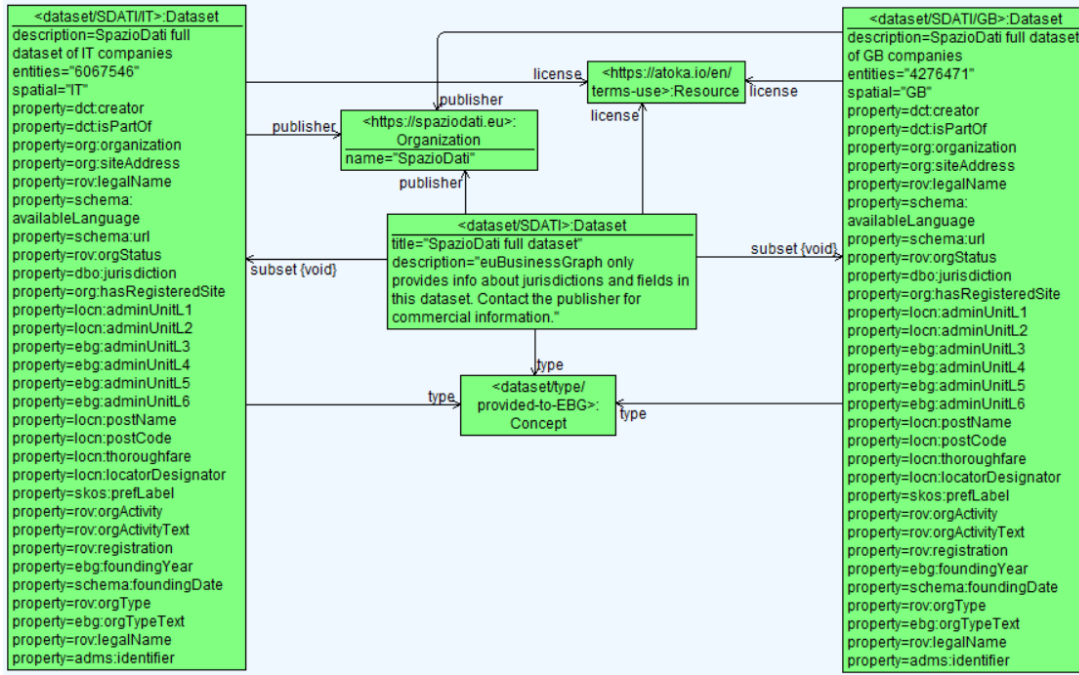


Fig. 13. Example of datasets provided by SpazioDati.

- **Precision:** Specifies that all data values for a business attribute must be as precise as required by the attribute's business requirements, intended meaning, intended usage, and precision in the real world.
- **Consistency:** Specifies that certain business attributes must follow a given pattern (e.g., age and dateOfBirth attributes are connected by the following rule $age = year(today) - year(dateOfBirth)$).
- **Temporal dimension:** Refers to the temporal dimension of data, such as volatility (the average time between update of data), timeliness (the average age of values), or currency (when data is entered in the system). An example of such a rule would be "the last modification date of attribute companyRevenue must be more recent than a year ago".

There are several possible ways to describe data validation rules, ranging from an algorithmic style such as:

```
legalName EXISTS AND len(trim(legalName)) <> 0
```

to a semantic based definition by using the SHACL [26] (Shapes Constraint Language) notation. SHACL is a language for validating RDF data graphs against a set of conditions that are provided as shapes and other constructs expressed in the form of an RDF graph (i.e., a shapes graph). ShEx [27] (Shape Expression) is a similar high-level language that can be used to validate RDF graph data. Both SHACL and ShEx use RDF syntax, and share the mechanisms of shape constraints, node constraints, property constraints, cardinalities, and logical operators. Examples of SHACL and ShEx shapes for the eu-BusinessGraph ontology are available in the Github repository⁴⁸. Figure 14 shows an example of how

⁴⁸<https://github.com/euBusinessGraph/eubg-data/tree/master/model>

SHACL validation shapes can be defined for a company URI node and two corresponding attributes (i.e., `legalName` and `orgActivity`). The `legalName` pattern requires the legal name to be canonicalized, i.e., not have leading, trailing or consecutive spaces (denoted as underscores below).

```
ebgsh:Company a sh:NodeShape;
  sh:targetClass rov:RegisteredOrganization;
  sh:closed true;
  sh:nodeKind sh:IRI;
  sh:pattern "^http://data.businessgraph.io/company/[A-Z]{2}/.+/" ;
  sh:property [sh:path rov:legalName;
    sh:or ([sh:datatype xsd:string] [sh:datatype rdf:langString]);
    sh:not ([sh:pattern "^_|\_|\$|_|_{2}"]); sh:minCount 1];
  sh:property [sh:path rov:orgActivity;
    sh:nodeKind sh:IRI;
    sh:pattern "^http://data.businessgraph.io/nace/.+"];
```

Fig. 14. Example of SHACL shape used to validate RDF company data.

5. Examples of Use of the euBusinessGraph Ontology

We present examples of how the euBusinessGraph ontology was used. We will first describe the approach on how the ontology was used to harmonize and make available company data from various data providers, resulting in the development of a business knowledge graph (Section 5.1 and Section 5.2). We will then show how this knowledge graph was used in the euBusinessGraph marketplace for basic company data—a place where data consumers can search, analyse, and compare data from various providers (Section 5.3). Finally, we provide an example how the ontology was used in the area of public procurement (Section 5.4), and how it was extended in the domain of financial transactions (Section 5.5).

5.1. Overview of Data Mapping Approach

In order to develop the euBusinessGraph knowledge graph harmonizing data from various data providers, we devised a data mapping approach that was used to convert company data from CSV and JSON sources into RDF conforming to the ontology. In the following, we describe the mapping notation and provide specific examples showing how the mapping rules were used. Actual mappings for data are publicly available via the DataGraft platform⁴⁹ [28, 29].

Figure 15 shows an instance diagram of the formal ontology that represents a specific company (i.e., SpazioDati) that is generated from raw JSON data, and provides an overview of typical attributes that we want to map from a JSON data format to the ontology. The first step of the mapping process is to select attributes (e.g., `base.legalName`) from the original data source (e.g., JSON file from data provider), and construct parameter names (e.g., `legalName`) so that we can reference the attribute values in the definition of the mapping functions, as exemplified in Table 2. When defining the mappings, we assume that the input data is a set of attribute-value pairs. Mapping parameters in Table 2 that are specified as lower-case italic letters refer to a string or number value (e.g., *legalName* refers to “SpazioDati S.R.L” in the data provider’s raw data source files), while parameters denoted in upper-case letters refer to SKOS concept schemes that were defined as part of the RDF generation process. As an example of the use of concept schemes, the mapping parameter `ORGACTIVITY` will refer to a URI that uses a classification vocabulary to represent the data attribute (e.g., the URI `<nace:62.01>` uses a controlled vocabulary⁵⁰

⁴⁹<https://datagraft.io>

⁵⁰<https://github.com/euBusinessGraph/eubg-data/blob/master/data/NACE/nace.ttl>

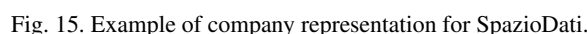


Table 2

Mapping parameter	Data provider's JSON data attribute
<i>id</i>	id
<i>legalName</i>	base.legalName
<i>jurisdiction</i>	base.country
ORGTYPE	base.legalForms[*].name
ORGACTIVITY	base.ateco[*].code
COUNTRY	base.registeredAddress.state
MACROREGION	base.registeredAddress.macroregion
REGION	base.registeredAddress.region
PROVINCE	base.registeredAddress.province
MUNICIPALITY	base.registeredAddress.municipality
<i>lat</i>	base.registeredAddress.lat
<i>lon</i>	base.registeredAddress.lon
LATLONPREC	base.registeredAddress.latlonPrecision

Next, Table 3 defines a set of helper functions for a subset of base URIs that will be used to map JSON data to RDF. The helper functions improve readability of mapping rules by reducing the text needed to refer to a specific URI. As an example, the helper function `curi` refers to the actual URI `http://data.businessgraph.io/company/IT/361163703`. To produce this URI, mapping parameters listed in *italic* (e.g., *jurisdiction* and *id*) will be replaced by the actual values (e.g., “IT” and “361163703”) from the source JSON data. Furthermore, the mapping definitions may contain input parameters denoted in **bold** that refer to another function that was defined as part of the mapping process (e.g., **ebg-comp** points to the URI `http://data.businessgraph.io/company`). After the set

of helper functions were defined, mapping rules were constructed for each of the data provider JSON attributes listed in Table 2. The resulting mapping rules are described in Table 4.

Table 3
Helper functions used to create base URIs

Helper function	Definition	Comments
ebg-comp	<code>http://data.businessgraph.io/company</code>	Base company URI
curi	<code>ebg-comp/jurisdiction/id</code>	Company URI
ciduri	<code>curi/id</code>	Company identifier URI
cadruri	<code>curi/address</code>	Company address URI
guri	<code>cadruri/geo</code>	Geographic coordinate URI

Using the mapping rules from Table 4 to transform JSON data to RDF for a specific company (e.g., SpazioDati) from data provider SpazioDati, will result in the subset of RDF triples listed below (e.g., applying the mapping function `<curi> rov:legalName "legalName"` to the source JSON data from the data provider):

```
<company/IT/361163703> rov:legalName "SPAZIODATI S.R.L." .
```

Table 4
Mapping functions for a subset of company data attributes

Scope of mapping function	Definition	Comments
Company URI node	<code><curi> rdf:type rov:RegisteredOrganization .</code>	Company class
	<code><curi> rov:registration <ciduri> .</code>	Company identifier triple
	<code><curi> org:hasRegisteredSite <cadruri> .</code>	Company address triple
	<code><curi> schema:geo <guri> .</code>	Company geo-coordinate triple
	<code><curi> rov:legalName "legalName" .</code>	Legal name
	<code><curi> dbo:jurisdiction "jurisdiction" .</code>	Jurisdiction
	<code><curi> rov:orgType ORGTYPE .</code>	Organization type
Identifier URI node	<code><curi> rov:orgActivity ORGACTIVITY .</code>	Economic activity
Address URI node	<code><ciduri> rdf:type adms:Identifier .</code>	Identifier class
	<code><ciduri> skos:notation "id" .</code>	Identifier value
	<code><cadruri> rdf:type locn:Address .</code>	Address class
	<code><cadruri> rdf:type org:Site .</code>	Address type
	<code><cadruri> org:siteAddress <cadruri> .</code>	Self reference
	<code><cadruri> locn:adminUnitL1 COUNTRY .</code>	Country
	<code><cadruri> locn:adminUnitL2 MACROREGION .</code>	Macro region
Geo-coordinate URI node	<code><cadruri> ebg:adminUnitL3 REGION .</code>	Region
	<code><cadruri> ebg:adminUnitL4 PROVINCE .</code>	Province
	<code><cadruri> ebg:adminUnitL5 MUNICIPALITY .</code>	Municipality
Geo-coordinate URI node	<code><guri> rdf:type schema:GeoCoordinates .</code>	Geolocation class
	<code><guri> schema:latitude lat .</code>	Latitude
	<code><guri> schema:longitude lon .</code>	Longitude
	<code><guri> ebg:geoResolution LATLONPREC .</code>	Geo-ordinate resolution

The following set of RDF triples were generated by using the mapping approach described in this section. The first three triples are produced by mapping source data to the ontology by use of SKOS concept schemes for the attributes `orgType`, `orgStatus` and `orgActivity`. The subsequent four triples refer

to different identifier systems that are associated with the company. Next, the proceeding four triples define actual values for SpazioDati using the identifier system “ATOKA”. Finally, the last five RDF triples show how geographical information for SpazioDati is mapped to the ontology with NUTS and LAU classification schemes.

```
<company/IT/361163703> rov:orgType <type/IT/SR> .
<company/IT/361163703> rov:orgStatus <status/SDATI/active> .
<company/IT/361163703> rov:orgActivity <nace/62.01> .

<company/IT/361163703> adms:identifier <company/IT/361163703/id/ATOKA> .
<company/IT/361163703> adms:identifier <company/IT/361163703/id/REA> .
<company/IT/361163703> adms:identifier <company/IT/361163703/id/Tax> .
<company/IT/361163703> adms:identifier <company/IT/361163703/id/Vat> .

<company/IT/361163703/id/ATOKA> dct:isPartOf <identifier/ATOKA> .
<company/IT/361163703/id/ATOKA> skos:notation "6da785b3adf2" .
<company/IT/361163703/id/ATOKA> rdf:type adms:Identifier .
<company/IT/361163703/id/ATOKA> dct:creator https://atoka.io .

<company/IT/361163703/registeredSite> locn:adminUnitL1 <http://nuts.geovocab.org/id/IT> .
<company/IT/361163703/registeredSite> locn:adminUnitL2 <http://nuts.geovocab.org/id/ITD> .
<company/IT/361163703/registeredSite> ebg:adminUnitL3 <http://nuts.geovocab.org/id/ITD2> .
<company/IT/361163703/registeredSite> ebg:adminUnitL4 <http://nuts.geovocab.org/id/ITD20> .
<company/IT/361163703/registeredSite> ebg:adminUnitL5 <lau/IT-022205> .
```

5.2. Infrastructure for the Knowledge Graph Generation

A data provisioning infrastructure was developed to onboard data from various data providers. Using this approach, data source files from data providers were processed and mapped to the euBusinessGraph ontology using the mapping process discussed in the previous section. After transforming each dataset from a tabular format (i.e., CSV or JSON) to RDF, the resulting data was published to one named graph for each data provider jurisdiction in an enterprise semantic graph database, GraphDB⁵¹, hosted by Ontotext.

GraphDB is a service component on the Ontotext Platform⁵² that implements GraphQL querying over RDF data. GraphQL is a simple query language in which the shape of the returned data (JSON) closely mirrors the shape of the query. It is a framework through which one can build simple, uniform and even federated *facades* over heterogeneous and complex data stores. Unlike traditional REST endpoints, one GraphQL query can access one or several data stores, and gets exactly the data that it has requested. Thus it is developer-friendly and has found a wide following with application developers. GraphQL Introspection is a standard way for the client to discover the schema of a GraphQL endpoint, enabling tools like *GraphiQL* to offer strong query completion features. The author of [30] describes an example of querying data about Star Wars and compares SPARQL to live GraphQL queries. The Ontotext platform uses a simple YAML-based language called Semantic Objects Modeling Language (SOML)⁵³ to describe a semantic model, generate a GraphQL schema and querying capabilities over it. The platform also has important features such as data mutations, user management (Fusion Auth), access control, deployment and monitoring.

⁵¹<http://graphdb.ontotext.com>

⁵²<http://platform.ontotext.com>

⁵³<http://platform.ontotext.com/soml>

In addition to GraphDB, the data provisioning infrastructure includes a set of data ingestion services and data preparation tools that can be used to simplify data cleaning and transformation from the various sources. The services include data interlinking tools for data transformation, enrichment, interlinking, and metadata generation processes in order to publish the business graph data as Linked Data.

Figure 16 illustrates the data provisioning process and the tools and services that are used to generate the business knowledge graph. Steps 1 and 2 of the illustration show that the core process of knowledge graph creation is executed by using the cloud-based data management platform DataGraft. Grafterizer⁵⁴ [31] is a framework (part of DataGraft) for interactive data cleaning and transformation, and RDF knowledge graph generation that is used together with the tabular annotation tool ASIA⁵⁵ [32] and ABSTAT⁵⁶ [33] to map company data to the euBusinessGraph ontology. Finally, in step 3, the RDF triples are published as a knowledge graph in GraphDB. Grafterizer, ASIA and ABSTAT were used to clean, transform, enrich and convert tabular data to RDF as part of the business knowledge graph construction. The eu-BusinessGraph ontology Github repository includes examples of a GraphQL query for some company data⁵⁷ (including auto-completion on Observation fields) and the corresponding result⁵⁸.

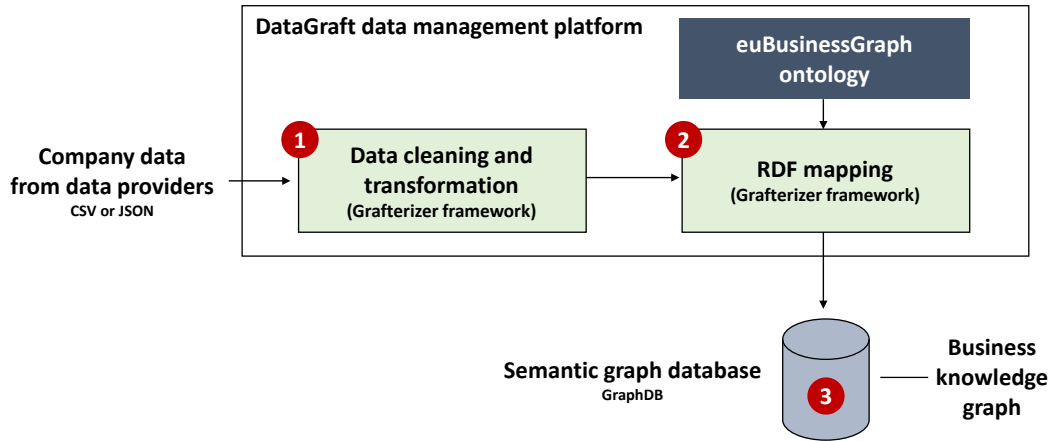


Fig. 16. The data provisioning process used to publish company data as part of the business knowledge graph.

Figures 17 and 18 show a specific example of how to map CSV data to RDF by using the tree mapping functionality in Grafterizer to build RDF triples. The following procedure exemplifies how the mapping rules defined in Section 5.1 can be used together with the infrastructure illustrated in Figure 16 to generate a company knowledge graph:

- (1) **Tabular transformation:** Figure 17 shows the first step of the process in which a raw CSV file is imported to the graphical user interface of Grafterizer. This step includes cleaning and transforming tabular data into a format that corresponds with the data validation rules described in Section 4.5.

⁵⁴<https://www.eubusinessgraph.eu/grafterizer-2-0>

⁵⁵<https://www.eubusinessgraph.eu/asia-2>

⁵⁶<https://www.eubusinessgraph.eu/abstat>

⁵⁷<https://github.com/euBusinessGraph/eubg-data/blob/master/example/GraphQL-Ontotext-query.png>

⁵⁸<https://github.com/euBusinessGraph/eubg-data/blob/master/example/GraphQL-Ontotext-result.png>

- (2) **RDF mapping:** Figure 18 illustrates the next step of the process where tabular data is ready to be mapped from the tabular format to the ontology by using the data mapping approach that was defined in Section 5.1 (e.g., the mapping function `<curi> rov:legalName "legalName"` is applied to the source input data by fetching the actual value from the tabular column "name"). This is a step-wise process in which each of the mapping rules are added in order to make the connection between the source data and the ontology to produce a full set of RDF triples.
- (3) **RDF storage:** Finally, the RDF data is uploaded and published to GraphDB to enable queries and create the foundation for the company data marketplace that will be described in the next section.

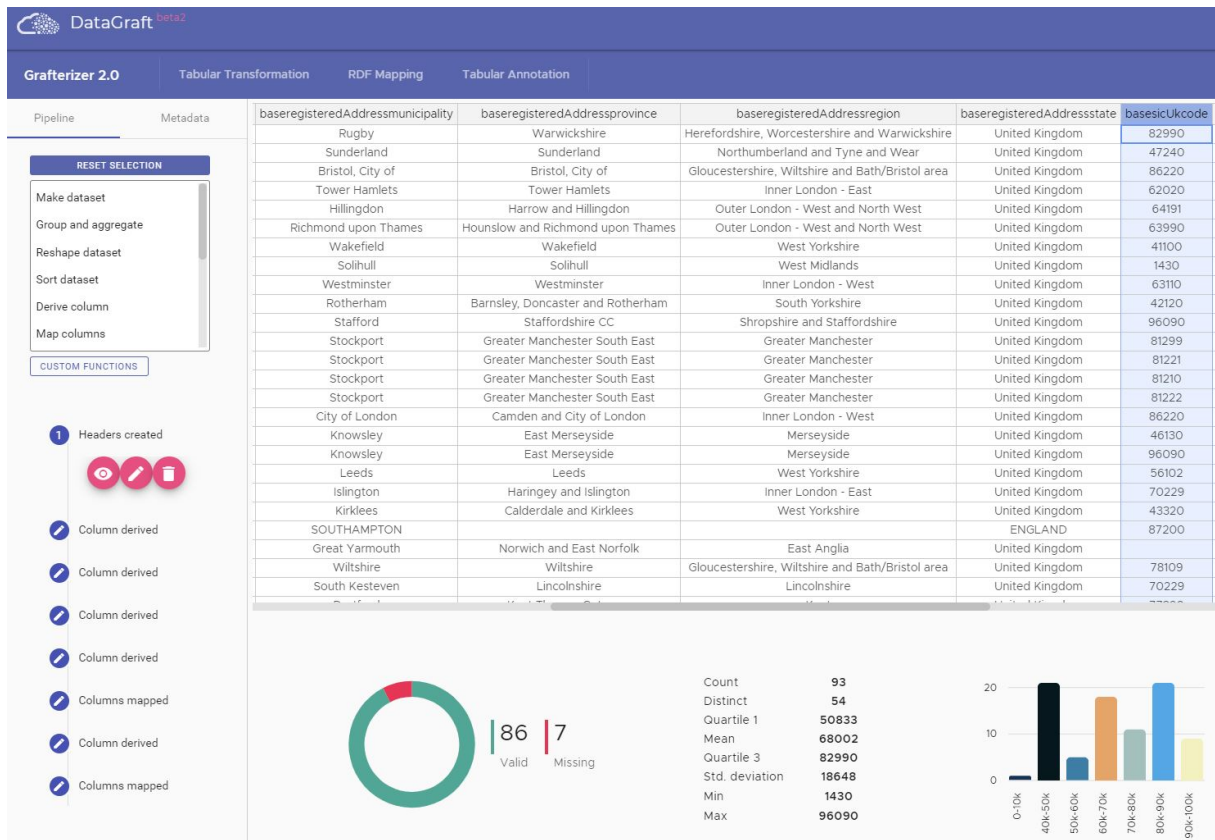


Fig. 17. Grafterizer user interface that shows the functionality for cleaning and transforming tabular data.

The repository hosted at GraphDB contains more than 1.4 Billion RDF triples of company data covering a subset of data from eight jurisdictions (i.e., countries). The RDF data was structured into named graphs for each data provider/jurisdiction to allow for duplicate triples of the same company from different providers. The named graphs `http://data.businessgraph.io/provider/sdati/uk` and `http://data.businessgraph.io/provider/ocorp/uk` for example can use the same company URI (e.g., `http://data.businessgraph.io/company/GB/02485441`) in the graph database without mingling the RDF statements from the two providers, and collapsing identical statements into

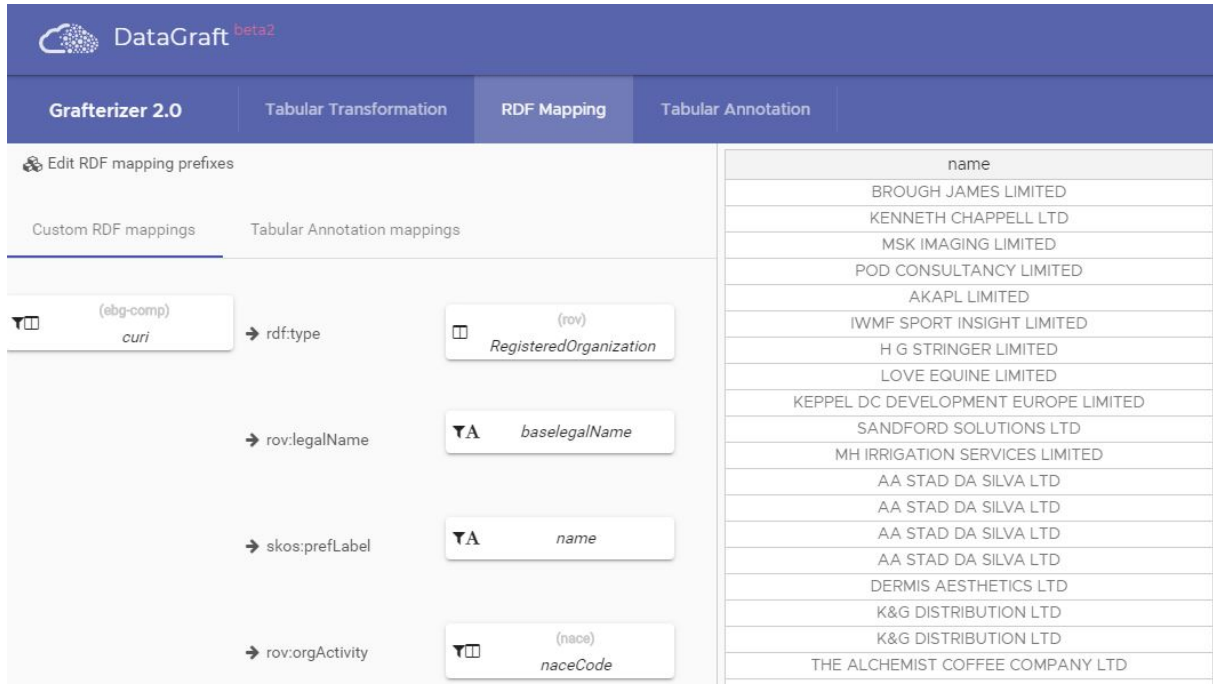


Fig. 18. Grafterizer user interface for the RDF mapping functionality.

one. As a result, several data providers can use the same identifier system for a specific company, and the repository currently contains named graphs for the following data providers and jurisdictions:

- Norway from provider BRC;
- Bulgaria from provider Ontotext;
- Italy from provider SpazioDati;
- UK from providers SpazioDati and OpenCorporates;
- Germany, France, Belgium and Luxembourg from provider OpenCorporates; and
- Norway from provider EVRY.

To demonstrate the data provisioning process and need for an ontology to structure company data, we chose to harmonize data at two levels of granularity. Data for jurisdictions Norway, Bulgaria, Italy, and UK were harmonized at a detailed level with regards to basic company attributes (e.g., name and founding date), identifier systems, and classification schemes (i.e., NACE, NUTS, LAU, organization types and organization status). Data for jurisdictions Germany, France, Belgium, and Luxembourg were harmonized with less detail (e.g., for jurisdiction Germany only highest level of NUTS classification is present for geographical location, and information about NACE economic classification is not available from data provider). The next section describes how the published knowledge graph was used to populate a marketplace for company data.

5.3. The euBusinessGraph Marketplace

A main motivation behind the development of a data marketplace for basic company data is the democratisation of the company information market, currently dominated by a few large international

players (e.g., Bisnode⁵⁹) that create a market barrier for smaller company data providers like OpenCorporates and SpazioDati. The intention of the marketplace is to enable such smaller players to join a common ecosystem to promote their data offerings, and for data consumers to have a central point where they could easily compare company data offerings. A public prototype of the data marketplace application⁶⁰ developed to showcase the use of the euBusinessGraph ontology is available online⁶¹.

The available data in the marketplace application includes the most central attributes that reflect how the ontology can be used to describe the semantic relations of company data. Each data provider URI in GraphDB is related to a dataset description that describes the data being offered in the marketplace by inserting `void:inDataset` for each `rov:RegisteredOrganization` in the graph database as illustrated in Figure 19.

```
base <http://data.businessgraph.io/>
prefix void: <http://rdfs.org/ns/void#>
prefix rov: <http://www.w3.org/ns/regorg#>
insert {
  graph ?g {?x void:inDataset ?d}
} where {
  values (?g ?d) {
    (<provider/ocorp/uk> <dataset/OCORP/EBG>)
    (<provider/ocorp/de> <dataset/OCORP/EBG>)
    (<provider/bgtr> <dataset/ONTO>)
    (<provider/brc> <dataset/BRC>)
    (<provider/sdati/it> <dataset/SDATI/EBG>)
    (<provider/sdati/uk> <dataset/SDATI/EBG>)
  }
  graph ?g {?x a rov:RegisteredOrganization}
}
```

Fig. 19. Linking data providers to dataset descriptions in the graph database.

As an example, the provider link `<provider/sdati/it>` points to subset `<dataset/SDATI/EBG>` which describes the subset of data from SpazioDati that is provided to the euBusinessGraph marketplace. Since SpazioDati can provide more detailed data about companies that is not available in the knowledge graph, the URI `<dataset/SDATI>` would include parts that are not provided to the marketplace, but only advertised in the marketplace application. On the other hand, all data from Brønnøysund Register Centre is open and fully provided to the business graph, and hence for `<dataset/BRC>` there is no need to describe subsets. Figure 21 shows how the ontology was used to differentiate between the data attributes that SpazioDati provides to the marketplace (e.g., the lower table) and all attributes available upon request (e.g., the upper table). Upon request, SpazioDati can provide detailed information about company officers, but this information is not fully provided to the knowledge graph.

Figure 20 shows how the ontology was used to represent company information in a consistent way for a subset of the company data attributes that are available from two data providers (i.e., OpenCorporates (OCORP) and SpazioDati (SDATI)) for jurisdiction GB (i.e., United Kingdom). Depending on the use case, data consumers have the opportunity to select the datasets that suit their needs. As an example, Figure 20 illustrates that OpenCorporates can provide information about dissolution date, while SpazioDati does not have this information. Other use cases open up for a combination of data from different data providers to achieve higher data coverage.

⁵⁹<http://www.bisnode.com>

⁶⁰<https://www.eubusinessgraph.eu/the-marketplace>

⁶¹<http://marketplace.businessgraph.io>

Select a jurisdiction to get an overview of which data individual data providers have for the selected jurisdiction.

Select Jurisdiction: **GB**

	legal name	type	type text	status	status text	economic activity	economic activity text	founding date	founding year	dissolution date	jurisdiction	registered site/address	identifier	full address
OCORP	●	●	●	●	●	●	●	●	●	●	●	●	●	●
SDATI	●	●	●	●		●		●	●		●	●	●	

Fig. 20. Availability of company data attributes from two different data providers for jurisdiction United Kingdom (GB).

SPAZIODATI
http://www.spaziodati.eu/en
SpazioDati is an innovative Italian company in the field of Marketing / B2B Leads Generation.

Data attributes available from provider:

	Is startup	Economic activity (NACE)	Founding date	Founding year	Dissolution date	Identifier	Identifier value	Full address	Geographical coordinates	Address NUTS level 1	Address NUTS level 2	Address NUTS level 3	Address NUTS level 4	Address LAU level 5	Address LAU level 6	Post name	Post code	Street address	Street number	Officer name	Officer role	We
GB	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●		●	●	●
IT	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●		●	●	●

*The data is provided under the license <https://atoka.io/en/terms-use/>.

Data attributes available from provider on euBusinessGraph marketplace:

	Is startup	Economic activity (NACE)	Founding date	Founding year	Dissolution date	Identifier	Identifier value	Full address	Geographical coordinates	Address NUTS level 1	Address NUTS level 2	Address NUTS level 3	Address NUTS level 4	Address LAU level 5	Address LAU level 6	Post name	Post code	Street address	Street number	Officer name	Officer role	We
GB	●	●	●	●	●	●	●	●		●	●	●	●	●	●	●	●	●				
IT	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●	●				

*The data is provided under the license <https://atoka.io/en/terms-use/>.

Fig. 21. Overview of company data attributes provided by SpazioDati for jurisdictions Italy and United Kingdom (GB).

The marketplace includes functionality for full-text advanced search and detailed faceted search for exploration of the company knowledge graph. Furthermore, the marketplace offers analytics services such as data aggregation and visualization (e.g., company activities per city), search for company news articles, and search for company events.

The ontology was used in the marketplace to realize use case scenarios such as:

- **Company search:** Find a specific company by displaying a page that describes available attributes of the company. The ontology enables search for detailed company information from different providers (e.g., SpazioDati and OpenCorporates), and facilitates data provenance, as the specific company data (i.e., for company APODACA LIMITED) from data provider OpenCorporates can be traced back to its sources (i.e., OpenCorporates and Companies House Register). In this specific

The screenshot displays the euBusinessGraph marketplace interface. The top navigation bar includes links for Home, Advanced Search (highlighted), Data Providers, Analytics Services, and About. The left sidebar contains four filter sections: PROVIDER, JURISDICTION, STATUS, and COMPANY TYPE. Each section lists various options with corresponding counts in small boxes. The main content area on the right shows a list of search results, each with a company name and its euBusinessGraph identifier.

Facet	Option	Count
PROVIDER	OpenCorporates dataset provided to EBG	28768868
	SpazioDati dataset provided to EBG	10167288
	BRC dataset provided to EBG	1016792
	EVERY dataset provided to EBG	1004155
	Ontotext dataset provided to EBG	950486
JURISDICTION	FR	11351031
	GB	11234043
	IT	6015977
	DE	5305727
	BE	1139250
	NO	1036862
	BG	950486
	LU	231214
	DK	1000
	STATUS	active
removed		2140579
inactive		597482
in liquidation		132223
bankrupt		4570
dormant		4201
in forced settlement or forced resolution		1174
COMPANY TYPE	Private Limited Company	9349648
	Other company type	868208
	Enterprise	509355
	organization	446851

Company Name	euBusinessGraph identifier
WANDSWORTH VISION	CE009473
THE GRACE MEMORIAL TRUST	CE013057
SKYE AND LOCHALSH MICRO ABATTOIR LIMITED	RS007702
10 DERBY ROAD MANAGEMENT LTD	11349461
STONEYBURN AND BENTS FUTURE VISION GROUP SCIO	CS002523
THE RELIGIOUS STUDIES PROJECTION ASSOCIATION (SCIO)	CS003013

Fig. 22. euBusinessGraph marketplace demonstrator that illustrates how the ontology was used to facilitate search and filtering on various facets such as company type and activity.

example, Companies House Register is the official source, while OpenCorporates is the unofficial data provider that uses data directly from the original Companies House Register sources.

- **Advanced company search:** Find how many companies are in a certain jurisdiction, active or inactive, registered in a certain year, with a certain type, in a certain location or are operating within a certain economic activity. This scenario is covered by allowing search for companies by certain criteria or facets and dynamic filtering of results. The search functionality of the marketplace demonstrates how the semantic model enables a uniform way of harmonizing and representing hierarchical facets for geographical location (i.e., NUTS and LAU) and economic classification (i.e., NACE). Hierarchical facets such as location and economic activity consist of several levels, allowing users to decide on the level of specificity of their search. The faceted search (Figure 22, left side) allows users to explore the knowledge graph and search for companies according to different criteria, such as provider, jurisdiction, company status and type. The full-text advanced

search (Figure 22, top page) will return a page where users can see all data that is available in the graph for a given company of interest, i.e., available data providers and identifiers, addresses, economic classifications, and company officers. In addition, companies are classified by NACE codes and linked to external systems, such as the national trade register of the company (e.g., Atoka⁶² and CompaniesHouse⁶³).

- ***Analytics related to company data:*** Find out how many companies are registered per year in a specific country and city, and are operating in a specific location. The marketplace application provides the ability to get basic statistics about the company data in the knowledge graph. A bar chart visualization filters information by country, city and activity and gives the user a visual representation of the data. By analysing the knowledge graph, we can get answers to questions such as a) which geographical areas in a country of interest have specific economic activities?, b) which geographical area has the lowest presence of companies in the accommodation sector?, c) which region has the highest number of companies?, and d) where do we find the highest number of new companies, registered the last two years?

5.4. Use of the euBusinessGraph Ontology in the Public Procurement Domain

Public procurement accounts for a substantial part of the public investment and global economy and therefore there is a need for better insight into, and management of government spending. In this respect, national, regional, local, and EU-wide public procurement portals were established to publish procurement notices regarding the purchase of work, goods or services from companies by public authorities in order to increase transparency, economic activity, and competitiveness [34]. However, the technical landscape is quite scattered and there are no common data formats and models used for exposing such data uniformly allowing advanced analytics and analysis, such as for fraud and trend detection. To this end, the euBusinessGraph ontology was used in the procurement domain, in the context of an project They-BuyForYou (TBFY)⁶⁴, for integrating public procurement and company data into the TBFY knowledge graph [35]. The resulting knowledge graph allows browsing, visualising, and analysing public EU-wide procurement data and enables a variety of business cases built on top of it by various stakeholders, such as buyers, suppliers, and policy makers.

The data integrated includes procurement data, provided by OpenOpps⁶⁵, and company data, provided by OpenCorporates. OpenOpps has gathered over 2M tender documents from more than 300 publishers through Web scraping and by using open APIs and provides the resulting data in Open Contracting Data Standard (OCDS)⁶⁶, while OpenCorporates uses its own ad-hoc schema. These two datasets are integrated through an ontology network. An ontology for procurement data was developed based on the OCDS standard [36] and the euBusinessGraph ontology was used for representing the company data. The two datasets are integrated through a reconciliation process [37]. Suppliers appearing in tender data are matched against company data provided by OpenCorporates. The matched company data is extracted and ingested to the TBFY knowledge graph. The current release of the TBFY knowledge graph includes 23M triples originating from tender data collected initially for the first quarter of 2019 and more data will be ingested.

⁶²<https://atoka.io/en>

⁶³<https://beta.companieshouse.gov.uk>

⁶⁴<http://theybuyforyou.eu>

⁶⁵<https://openopps.com>

⁶⁶<https://standard.open-contracting.org/latest/en>

5.5. Use of the euBusinessGraph Ontology for Financial Transactions

Company-related economic information is crucial to many business operations. It empowers customer relationship management, acquisition of new clients, marketing campaigns, supply chain management, market analysis, competitive intelligence, mergers and acquisitions, etc. In this respect, the euBusinessGraph ontology was used for matching and linking company-related economic information within the context of Ontotext's Intelligent Matching and Linking of Company Data (CIMA) project⁶⁷. CIMA aims to use AI/ML technologies for linking and harmonizing company-related business data from various sources. The project applies machine learning, semantic modeling and integration, entity matching, automatic classification, logical inference to make data richer, better harmonized, integrated, interlinked and easier to use. As part of the project, Ontotext is creating a Company Knowledge Graph (ONTO-CG) for demo purposes by integrating data from open and a few proprietary datasets. The emphasis of the project is on financial data, industrial classification, company size/importance observations (e.g., annual sales, number of employees, etc.).

ONTO-CG builds upon the euBusinessGraph ontology and adds the following:

- **IdentifierSystems**: The identifier idea is extended to record any kind of useful identification info in a generic way such as phone, email, and website; profile links and identifiers in various external systems such as Wikidata, DBpedia, Facebook, Thomson Reuters permid (TR), and ISO 10383 Market Identifier Code (MIC); and research-oriented identifiers such as CrossRef funder, and Global Research Identifier Database (GRID).
- **cg:StockExchange**: a stock exchange where companies can offer shares or other securities. We record MIC and TR exchange codes as identifiers.
- **cg:Event** and **cg:EventAppearance**: Conference, workshop, meetup, etc., where the work of a certain person or company may be highlighted.
- **gn:Feature**: While the euBusinessGraph geographic hierarchy is based on EuroStat NUTS and LAU, ONTO-CG uses Geonames locations to implement geographic matching, auto-completion and faceting.
- **cg:AcademicQualification**: Academic degree (completed or not) of a person at a school in an academic major.
- **qb:Observation**: Statistical or other observation about an object (typically company), such as annual sales, number of employees, etc. It may be for a particular year, point in time, or without date (current).
- **cg:Transaction**: Financial transaction that gives money to a company in return for shares or other consideration.
- **cg:OrganizationRelation**: Relation between two agents. For asymmetric relations two fields "agentMinor" (e.g., subsidiary, owned, supplier) and "agentMajor" (e.g., parent, owner, customer) are used, and for symmetric relations the field "agent" is used twice.
- **Sourcing (provenance) for each node**: This includes `void:Dataset`, `dataset` as source of entities, `void:Linkset`, `linkset` as source of identifiers (links), and `cg:SourceMatch`, cluster of matched lower-level entities as the source of a higher-level entity.

⁶⁷<https://www.ontotext.com/cima>

In addition to the above new classes, ONTO-CG adds a 2-level data model where data from individual datasets sits at a lower (KG-building) level, and after matching and data fusion is promoted at a higher (data consumption) level. It also provides various extra fields such as `cg:geoPrecision` (precision of geo coordinates in meters); various flags such as for organization (`cg:isResearch`), position (`cg:isCurrent`, `cg:isPrimary`), academic qualification (`cg:isCompleted`), and organization relation (`cg:isCurrent`); and business nomenclatures (`skos:ConceptScheme`) including such as organization type, legal form, investor type, position type, transaction type, and relation type.

6. Conclusion and Outlook

As part of the work in this article, the analysis of existing initiatives in the area of interoperability of company-related data revealed the fact that harmonization of company data was far from a solved problem. We argued for the importance of harmonised basic company data as a key enabler for different value chains in various sectors that depend on company information. In this article, we described the euBusinessGraph ontology for harmonizing basic company data as a lightweight mechanism for aggregating, linking, provisioning and analysing basic company data.

The euBusinessGraph ontology was developed following standard practices in ontology development, identifying the scope and competency questions with different stakeholders, identifying and reusing existing ontologies, and publishing the ontology according to existing best practices for Linked Data vocabulary publishing. We provided an overview of the ontology scope, the ontology development process, explanations of core concepts and relationships, and the implementation of the ontology. Furthermore, we provided examples where the ontology was used, among others, for publishing company data and for comparing company data from various data providers.

The euBusinessGraph ontology serves now as an asset not only for enabling various tasks related to basic company data but also on top of which more specific extensions can be built upon. As an example of such an extension, initial efforts have been made to capture events that happen during the lifetime of a company [38] and for representing the French register data in RDF [38, 39]. In additions to possible extensions of the ontology, other interesting directions for future work can be envisioned. For example, interlinking harmonized data from various data providers is an interesting topic for future work (preliminary work on interlinking company data harmonised using the euBusinessGraph ontology is reported in [40]). Extending the ontology with classification datasets for additional jurisdictions (e.g., Germany) will further increase the relevance of the business graph, and enable more precise queries to be executed on the harmonized data. This harmonization process includes describing supplementary identifier systems for company entities and officers for new data providers, as well as creating additional classification schemes for NACE, NUTS, LAU, organization types and organization status.

In the TheyBuyForYou project, the ontology will be used as a core component of the proposed procurement knowledge graph and the ontology network. Currently, on the one hand, more data is being reconciled and ingested into the TBFY knowledge graph and on the other hand more research and development work is being undertaken in order to improve the reconciliation process matching supplier data against company data. Essentially, it will demonstrate how one can integrate disparate but relevant data sources, pose interesting queries that were otherwise not possible to answer, and create new business scenarios. In CIMA (ONTO-CG), the euBusinessGraph semantic model is extended to cover financial transactions and innovation assessments, and prototypes and exploitable systems are built using the Ontotext Platform and GraphQL over RDF data integrated from numerous sources.

Acknowledgement

The work in this article was partly funded by the EC H2020 projects euBusinessGraph (grant 732003), EW-Shopp (grant 732590), TheyBuyForYou (grant 780247), and CIMA (Bulgarian grant BG16RFOP002-1.005-0168-C01). Special thanks to the members of the euBusinessGraph project consortium for stimulating discussions around various aspects of basic company information, especially to Tatiana Tarasova, Fredrik Seehusen, and David Norheim for their initial involvement in the development of the ontology.

References

- [1] M. Janssen, D. Konopnicki, J.L. Snowdon and A. Ojo, Driving public sector innovation using big and open linked data (BOLD), *Information Systems Frontiers* **19**(2) (2017), 189–195. doi:10.1007/s10796-017-9746-2.
- [2] T. Heath and C. Bizer, *Linked Data: Evolving the Web into a Global Data Space*, Morgan & Claypool, 2011.
- [3] S.K. Bansal and S. Kagemann, Integrating Big Data: A Semantic Extract-Transform-Load Framework, *IEEE Computer* **48**(3) (2015), 42–50. doi:10.1109/MC.2015.76.
- [4] M. Giese, A. Soylu, G. Vega-Gorgojo, A. Waaler, P. Haase, E. Jiménez-Ruiz, D. Lanti, M. Rezk, G. Xiao, Ö.L. Özçep and R. Rosati, Optique: Zooming in on Big Data, *IEEE Computer* **48**(3) (2015), 60–67. doi:10.1109/MC.2015.82.
- [5] D. Reynolds (ed.), "The Organization Ontology", World Wide Web Consortium (W3C), 2014. <https://www.w3.org/TR/vocab-org/>.
- [6] J.F. Muñoz-Soro, G. Esteban, O. Corcho and F. Seron, PPROC, an ontology for transparency in public procurement, *Semantic Web* **7**(3) (2016), 295–309. doi:10.3233/SW-150195.
- [7] Semantic Interoperability Community, "e-Government Core Vocabularies", European Commission - ISA Programme, 2019. <https://joinup.ec.europa.eu/solution/e-government-core-vocabularies>.
- [8] Working Group for Describing Public Services, "Core Public Service Vocabulary Application Profile (CPSV-AP)", European Commission - ISA² Programme, 2016. https://ec.europa.eu/isa2/solutions/core-public-service-vocabulary-application-profile-cpsv-ap_en.
- [9] R.V. Guha, D. Brickley and S. Macbeth, Schema.org: evolution of structured data on the web, *Communications of the ACM* **59**(2) (2016), 44–51. doi:10.1145/2844544.
- [10] M. Bennett, The financial industry business ontology: Best practice for big data, *Journal of Banking Regulation* **14**(3) (2013), 255–268. doi:10.1057/jbr.2013.13.
- [11] M. McDaniel and V.C. Storey, Evaluating Domain Ontologies: Clarification, Classification, and Challenges, *ACM Computing Survey* **52**(4) (2019), 70:1–70:44. doi:10.1145/3329124.
- [12] Department of Economic and Social Affairs, "International Standard Industrial Classification of All Economic Activities (ISIC)", United Nations, 2008. <https://unstats.un.org/unsd/classifications/Econ/isic>.
- [13] Eurostat, "Statistical classification of economic activities in the European Community (NACE)", European Commission, 2008. <https://ec.europa.eu/eurostat/en/web/products-manuals-and-guidelines/-/KS-RA-07-015>.
- [14] ISO/TC 68/SC 8 Technical Committee, "Financial services – Legal entity identifier (LEI)", International Organization for Standardization (ISO), 2019. <https://www.iso.org/standard/75998.html>.
- [15] Eurostat, "Methodological manual on territorial typologies", European Commission, 2019. doi:10.2785/930137. <https://ec.europa.eu/eurostat/web/products-manuals-and-guidelines/-/KS-GQ-18-008>.
- [16] EU ISA Programme Core Vocabularies Working Group, "ISA Programme Location Core Vocabulary", World Wide Web Consortium (W3C), 2015. <https://www.w3.org/ns/locn>.
- [17] M. Dekkers, "Asset Description Metadata Schema (ADMS)", World Wide Web Consortium (W3C), 2013. <https://www.w3.org/TR/vocab-adms/>.
- [18] K. Alexander, R. Cyganiak, M. Hausenblas and J. Zhao, "Describing Linked Datasets with the VoID Vocabulary", World Wide Web Consortium (W3C), 2011. <https://www.w3.org/TR/void/>.
- [19] T. Baker, S. Bechhofer, A. Isaac, A. Miles, G. Schreiber and E. Summers, Key choices in the design of Simple Knowledge Organization System (SKOS), *Journal of Web Semantics* **20** (2013), 35–49. doi:10.1016/j.websem.2013.05.001.
- [20] W.R. van Hage, V. Malaisé, R. Segers, L. Hollink and G. Schreiber, Design and use of the Simple Event Model (SEM), *Journal of Web Semantics* **9**(2) (2011), 128–136. doi:10.1016/j.websem.2011.03.003.
- [21] N.F. Noy and D.L. McGuinness, Ontology Development 101: A Guide to Creating Your First Ontology, Technical Report, Stanford Medical Informatics, 2001.

- [22] O. Corcho, M. Fernández-López and A. Gómez-Pérez, *Ontological Engineering: Principles, Methods, Tools and Languages*, in: *Ontologies for Software Engineering and Software Technology*, C. Calero, F. Ruiz and M. Piattini, eds, Springer Berlin Heidelberg, 2006, pp. 1–48. doi:10.1007/3-540-34518-3_1.
- [23] J. Barzdins, K. Cerans, R. Liepins and A. Sprogis, Advanced Ontology Visualization with OWLGrEd, in: *Proceedings of the 8th International Workshop on OWL: Experiences and Directions (OWLED 2011)*, CEUR Workshop Proceedings, Vol. 796, CEUR-WS.org, 2011. http://ceur-ws.org/Vol-796/owled2011_submission_7.pdf.
- [24] V. Alexiev, T. Tarasova, J. Paniagua, C. Taggart, B. Elvesæter, F. Seehusen, D. Roman and D. Norheim, *euBusinessGraph Semantic Data Model*, euBusinessGraph Consortium, 2018. https://docs.google.com/document/d/1dhMOTIIOC6dOK_jksJRX0CB-GIRoiYY6fWtCnZArUhU/edit.
- [25] V. Alexiev, A. Kiryakov and P. Tarkalanov, euBusinessGraph: Company and Economic Data for Innovative Products and Services, in: *Proceedings of the 13th International Conference on Semantic Systems (Semantics 2017)*, 2017. http://rawgit2.com/webdata/SEMANTICS2017-posters/master/papers_final/163_Alexiev/index.html.
- [26] H. Knublauch and D. Kontokostas (eds), "Shapes constraint language (SHACL)", World Wide Web Consortium (W3C), 2017. <https://www.w3.org/TR/shacl/>.
- [27] E. Prud'hommeaux, J.E. Labra Gayo and H. Solbrig, Shape expressions: an RDF validation and transformation language, in: *Proceedings of the 10th International Conference on Semantic Systems (SEM 2014)*, ACM, 2014, pp. 32–40.
- [28] D. Roman, N. Nikolov, A. Putlier, D. Sukhobok, B. Elvesæter, A. Berre, X. Ye, M. Dimitrov, A. Simov, M. Zarev, R. Moynihan, B. Roberts, I. Berlocher, S. Kim, T. Lee, A. Smith and T. Heath, DataGraft: One-stop-shop for open data management, *Semantic Web* 9(4) (2018), 393–411. doi:10.3233/SW-170263.
- [29] D. Roman, M. Dimitrov, N. Nikolov, A. Putlier, D. Sukhobok, B. Elvesæter, A. Berre, X. Ye, A. Simov and Y. Petkov, Datagraft: Simplifying open data publishing, in: *European Semantic Web Conference*, Springer, 2016, pp. 101–106.
- [30] J. Rayfield, A New Hope: The Rise of the Knowledge Graph. Navigating through the Star Wars universe with knowledge graphs, SPARQL and GraphQL, 2019. <https://www.ontotext.com/blog/the-rise-of-the-knowledge-graph/>.
- [31] D. Sukhobok, N. Nikolov, A. Pultier, X. Ye, A.J. Berre, R. Moynihan, B. Roberts, B. Elvesæter, M. Nivethika and D. Roman, Tabular Data Cleaning and Linked Data Generation with Grafterizer, in: *Proceedings of The Semantic Web - ESWC 2016 Satellite Events*, LNCS, Vol. 9989, Springer, 2016, pp. 134–139. doi:10.1007/978-3-319-47602-5_27.
- [32] V. Cutrona, M. Ciavotta, F.D. Paoli and M. Palmonari, ASIA: a Tool for Assisted Semantic Interpretation and Annotation of Tabular Data, in: *Proceedings of the ISWC 2019 Satellite Tracks (Posters & Demonstrations, Industry, and Outrageous Ideas) co-located with 18th International Semantic Web Conference (ISWC 2019)*, CEUR Workshop Proceedings, Vol. 2456, CEUR-WS.org, 2019, pp. 209–212. <http://ceur-ws.org/Vol-2456/paper54.pdf>.
- [33] R.A.A. Principe, B. Spahiu, M. Palmonari, A. Rula, F.D. Paoli and A. Maurino, ABSTAT 1.0: Compute, Manage and Share Semantic Profiles of RDF Knowledge Graphs, in: *Proceedings of The Semantic Web: ESWC 2018 Satellite Events - ESWC 2018 Satellite Events*, LNCS, Vol. 11155, Springer, 2018, pp. 170–175. doi:10.1007/978-3-319-98192-5_32.
- [34] E. Simperl, Ó. Corcho, M. Grobelnik, D. Roman, A. Soylu, M.J.F. Ruíz, S. Gatti, C. Taggart, U.S. Klima, A.F. Uliana, I. Makgill and T.C. Lech, Towards a Knowledge Graph Based Platform for Public Procurement, in: *Proceedings of the 12th International Conference on Metadata and Semantic Research (MTSR 2018)*, 2018, pp. 317–323. doi:10.1007/978-3-030-14401-2_29.
- [35] A. Soylu, Ó. Corcho, E. Simperl, D. Roman, F.Y. Martínez, C. Taggart, I. Makgill, B. Elvesæter, B. Symonds, H. McNally, G. Konstantinidis, Y. Zhao and T.C. Lech, Towards Integrating Public Procurement Data into a Semantic Knowledge Graph, in: *Proceedings of the Posters and Demonstrations Session of 21st International Conference on Knowledge Engineering and Knowledge Management (EKAW 2018)*, CEUR Workshop Proceedings, Vol. 2262, CEUR-WS.org, 2018. <http://ceur-ws.org/Vol-2262/ekaw-poster-01.pdf>.
- [36] A. Soylu, B. Elvesæter, P. Turk, D. Roman, Ó. Corcho, E. Simperl, G. Konstantinidis and T.C. Lech, Towards an Ontology for Public Procurement Based on the Open Contracting Data Standard, in: *Proceedings of the 18th IFIP WG 6.11 Conference on e-Business, e-Services, and e-Society (I3E 2019)*, Vol. 11701, 2019, pp. 230–237. doi:10.1007/978-3-030-29374-1_19.
- [37] A. Soylu, B. Elvesæter, P. Turk, D. Roman, Ó. Corcho, E. Simperl, I. Makgill, C. Taggart, M. Grobelnik and T.C. Lech, An Overview of the TBFY Knowledge Graph for Public Procurement, in: *Proceedings of the ISWC 2019 Satellite Tracks (Posters & Demonstrations, Industry, and Outrageous Ideas)*, CEUR Workshop Proceedings, Vol. 2456, CEUR-WS.org, 2019, pp. 53–56. <http://ceur-ws.org/Vol-2456/paper14.pdf>.
- [38] S.A.E. Kader, N. Nikolov, B.M. von Zernichow, V. Cutrona, B.E. M. Palmonari, A. Soylu, and D. Roman, Modeling and Publishing French Business Register (Sirene) Data as Linked Data Using the euBusinessGraph Ontology, in: *Proceedings of Semantic Statistics (SemStats 2019)*, 2019.
- [39] T. Ehrhart and R. Troncy, EURECOM at SemStats 2019, in: *Proceedings of Semantic Statistics (SemStats 2019)*, 2019.
- [40] A. Maurino, A. Rula, B.M. von Zernichow, M.S. Gomez, B. Elvesæter and D. Roman, Modelling and Linking Company Data in the euBusinessGraph Platform, in: *Proceedings of the 5th Workshop on Data Science for Macro-Modeling with Financial and Economic Datasets (DSMM 2019)*, ACM, 2019. doi:10.1145/3336499.3338012.