

When Linguistics Meets Web Technologies. Recent advances in Modelling Linguistic Linked Open Data

Anas Fahad Khan^a, Christian Chiarcos^b, Thierry Declerck^c, Daniela Gifu^d,
Elena González-Blanco García^e, Jorge Gracia^f, Maxim Ionov^b, Penny Labropoulou^h,
Francesco Mambriniⁱ, John P. McCrae^j, Émilie Pagé-Perron^k, Marco Passarottiⁱ,
Salvador Ros Muñoz^l, Ciprian-Octavian Truică^m

^a *Istituto di Linguistica Computazionale «A. Zampolli», Consiglio Nazionale delle Ricerche, Italy*
E-mail: fahad.khan@ilc.cnr.it

^b *Applied Computational Linguistics Lab, Goethe-Universität Frankfurt am Main, Germany*
E-mail: chiarcos@informatik.uni-frankfurt.de,
E-mail: ionov@informatik.uni-frankfurt.de

^c *DFKI GmbH, Multilinguality and Language Technology, Saarbrücken, Germany*
E-mail: declerck@dfki.de

^d *Faculty of Computer Science, Alexandru Ioan Cuza University of Iasi, Romania*
E-mail: daniela.gifu@info.uaic.ro

^e *Laboratory of Innovation on Digital Humanities, IE University, Spain*
E-mail: egonzalezblanco@faculty.ie.edu

^f *Aragon Institute of Engineering Research, University of Zaragoza, Spain*
E-mail: jogracia@unizar.es

^g *Applied Computational Linguistics Lab, Goethe-Universität Frankfurt am Main, Germany*
E-mail: ionov@informatik.uni-frankfurt.de

^h *Institute for Language and Speech Processing, Athena Research Center, Greece*
E-mail: penny@athenarc.gr

ⁱ *CIRCSE Research Centre, Università Cattolica del Sacro Cuore, Milan, Italy*
E-mail: francesco.mambrini@unicatt.it,
E-mail: marco.passarotti@unicatt.it

^j *Insight SFI Research Centre for Data Analytics, Data Science Institute, National University of Ireland Galway, Ireland*
E-mail: john.mccrae@insight-centre.org

^k *Wolfson College, University of Oxford, United Kingdom*
E-mail: emilie.page-perron@wolfson.ox.ac.uk

^l *Laboratory of Innovation on Digital Humanities, National Distance Education University UNED, Spain*
E-mail: sros@scc.uned.es

^m *Computer Science and Engineering Department, Faculty of Automatic Control and Computers, University Politehnica of Bucharest, Romania*
E-mail: ciprian.truica@upb.ro

Abstract. This article provides an up-to-date and comprehensive survey of models (including vocabularies, taxonomies and ontologies) used for representing linguistic linked data (LLD). It focuses on the latest developments in the area and both builds

upon and complements previous works covering similar territory. The article begins with an overview of recent trends which have had an impact on linked data models and vocabularies, such as the growing influence of the FAIR guidelines, the funding of several major projects in which LLD is a key component, and the increasing importance of the relationship of the digital humanities with LLD. Next, we give an overview of some of the most well known vocabularies and models in LLD. After this we look at some of the latest developments in community standards and initiatives such as OntoLex-Lemon as well as recent work which has been carried out in corpora and annotation and LLD including a discussion of the LLD metadata vocabularies META-SHARE and *lime* and language identifiers. In the following part of the paper we look at work which has been realised in a number of recent projects and which has a significant impact on LLD vocabularies and models.

Keywords: linguistic linked data, FAIR, corpora, annotation, language resources, OntoLex-Lemon, Digital Humanities, metadata, models

1. Introduction

The growing popularity of linked data, and especially of linked *open* data (that is, linked data with an open license), as a means of publishing language resources (lexica, corpora, data categories, etc.) necessitates a greater emphasis on models for linguistic linked data (LLD) since these are key to what makes linked data resources so reusable and so interoperable (at a semantic level). The purpose of this article is to provide a comprehensive and up-to-date survey of models used for representing linguistic linked data. It will focus on the latest developments and will both build upon as well as trying to complement previous works covering similar territory by avoiding too much repetition and overlap with the latter.

In the following section, Section 2, we give an overview of a number of trends from the last few years which have had, or which are likely to have, a significant impact on the definition and/or use of LLD models. We relate these trends to the rest of the article by highlighting relevant sections of the article (in bold). This overview of trends will help to locate the present work within a wider research context, something that is extremely useful in an area as active as linguistic linked data, as well as assisting readers in navigating the rest of the article. Next, in Section 2.4, we compare the present article with other related work, including an earlier survey of LLD models, in order to help clarify the topics and approach of the present work. Section 3 gives an overview of the most widely used models in LLD. Then in Section 4, we look at recent developments in community standards and initiatives. These include the latest extensions of the OntoLex-Lemon model in Section 4.1, a discussion of relevant work in corpora and annotations in Section 4.2, and a section on metadata Section 4.3. Finally there is a sec-

tion discussing projects, Section 5, and the conclusion, Section 6.

2. Setting the Scene: An Overview of Relevant Trends for LLD

The trends we have decided to focus on in this overview are the FAIRification of data in **Section 2.1**, the importance of projects to LLD models in **Section 2.2**, and finally the increasing importance of Digital Humanities use cases in **Section 2.3**.

2.1. FAIR New World

With the growing importance of Open Science initiatives, and especially those promoting the FAIR guidelines (where FAIR stands for Findable, Accessible, Interoperable and Reusable) [1] – and the consequent emphasis on the modelling, creation and publication of language resources as FAIR digital resources – shared models and vocabularies have begun to take on an increasingly prominent role. Although the linguistic linked data community has been active in promoting shared RDF vocabularies and models for years, this new emphasis on FAIR is likely to have a considerable impact in several ways, not least in terms of the necessity for these models to demonstrate a greater coverage, and to be more interoperable one with another. We will look at one series of FAIR related recommendations for models in Section 3 and see how they might be applied to the case of LLD. However in the rest of the subsection we will take a closer look at the FAIR principles themselves and show why their widespread adoption is likely to lead to a greater role for LLD models and vocabularies in the future.

In *The FAIR Guiding Principles for scientific data management and stewardship* [1], the article which

1 first articulated the well known FAIR principles, the
 2 authors clearly state that the criteria proposed by these
 3 principles are intended both "for machines and peo-
 4 ple" and that they provide "'steps along a path' to ma-
 5 chine actionability", where the latter is understood to
 6 describe structured data that would allow a "computa-
 7 tional data explorer" to determine:

- 8 – The type of a "digital research object"
- 9 – Its usefulness with respect to tasks to be carried
 10 out
- 11 – Its usability especially with respect to licensing
 12 issues, represented in a way that would allow the
 13 agent to take "appropriate action".
 14

15 The current popularity of the FAIR principles and,
 16 in particular, their promotion by governments and re-
 17 search funding bodies, such as the European Commis-
 18 sion,¹ through several national and international ini-
 19 tiatives reflects a wider recognition of the potential of
 20 structured and machine actionable data in changing
 21 how research is carried out, and especially in helping
 22 to support open science practices. The FAIR ideal, in
 23 short, is to allow machines as much autonomy as pos-
 24 sible in working with data, by the expedient of render-
 25 ing as much of the semantics of that data explicit (and
 26 machine actionable) as possible.

27 Publishing data using a standardised data model like
 28 the Resource Description Framework² (RDF) which
 29 was specifically intended to facilitate interoperability
 30 and interlinking between datasets – along with the
 31 other standards proposed in the Semantic Web stack
 32 and the technical infrastructure which has been devel-
 33 oped in order to support it – obviously goes a long way
 34 towards facilitating the publication of datasets as FAIR
 35 data. In addition, however, it is also vital that there ex-
 36 ist specialised vocabularies/terminologies/models and
 37 data category registries in order to ensure a viable,
 38 domain-wide level of interoperability and re-usability
 39 of data. These former resources serve to describe the
 40 shared theoretical assumptions held by a community of
 41 experts with regard to the semantics of the terms used
 42 by that community, and do so in a form that is (to vary-
 43 ing extents) machine readable to computational agents.
 44 The following FAIR principles are especially salient
 45 here:
 46

- 47 – F2. data are described with rich metadata.

48
 49 ¹https://ec.europa.eu/info/sites/info/files/turning_fair_into_reality_0.pdf

50 ²<https://www.w3.org/TR/rdf-primer/>

- 1 – I1. (meta)data use a formal, accessible, shared,
 2 and broadly applicable language for knowledge
 3 representation.
- 4 – I2. (meta)data use vocabularies that follow FAIR
 5 principles.
 6

7 It is important to note that the emphasis placed
 8 on machine actionability in FAIR resources (that
 9 is, on enabling computational agents to find rele-
 10 vant datasets and resources and to take "appropriate
 11 action" when they find them) gives Semantic Web
 12 vocabularies/registries a substantial advantage over
 13 other (non-Semantic Web native) standards in the
 14 field of linguistics like the Text Encoding Initiative
 15 (TEI) guidelines³ [2], the Lexical Markup Framework
 16 (LMF) [3] or the Morpho-syntactic Annotation Frame-
 17 work (MAF) [4].

18 For a start, none of these other standards possess
 19 a 'native', widely-used, widely technically supported
 20 knowledge representation language for describing the
 21 semantics of vocabulary terms in a machine readable
 22 way, or at least nothing as powerful as the Web Ontol-
 23 ogy Language (OWL)⁴ or the Semantic Web Rule Lan-
 24 guage (SWRL)⁵. For instance., there is no standard-
 25 ised way of describing the meanings of morphemes,
 26 lexemes, lemmas, etc. in TEI in a machine actionable
 27 way.

28 The ability to give precise, axiomatic definitions of
 29 terms in a formal knowledge representation (KR) lan-
 30 guage (allied with already established conceptual mod-
 31 elling techniques and ontology engineering best prac-
 32 tises) is especially helpful in humanistic disciplines
 33 such as linguistics or literary scholarship, where there
 34 can often be quite different definitions of the same
 35 or similar core concepts, e.g., with respect to differ-
 36 ent scholarly traditions or schools of thought. Using
 37 a machine readable description in OWL, once again,
 38 in conjunction with an ontology modelling methodol-
 39 ogy such as OntoClean [5], and together a more hu-
 40 man readable description given as documentation, can
 41 help to clarify (according to the expressive limitations
 42 of OWL) what we mean when we use a concept like
 43 'Sense' or 'Morpheme' in a dataset. It also facilitates
 44 the machine readable description of the relationship
 45 between different definitions of concepts across lan-
 46 guages or traditions.
 47

48
 49 ³<https://tei-c.org/guidelines/>

50 ⁴<https://www.w3.org/TR/2012/REC-owl2-overview-20121211/>

51 ⁵<https://www.w3.org/Submission/SWRL/>

Secondly, thanks to the use of a shared data model and a powerful native linking mechanism, linguistic linked data datasets can be easily (and in a standard way) integrated with/enriched by (linked data) datasets belonging to other disciplines, for instance geographical and historical datasets or gazettiers and authority lists. OWL, and vocabularies, such as PROV-O,⁶ also allow us to add information pertaining to when something happened, or whether we are describing a hypothesis⁷ or not (in which case, also who made it and when). Once again all of these things can be described in a way that makes the semantics of the information (relatively) explicit and machine actionable through the use of pre-existing standards and technologies including the Semantic Web query language *SPARQL Protocol and RDF Query Language* (SPARQL) as well as freely available Semantic Web reasoning engines.

Moreover the pursuit of the FAIR ideal has opened the way to new means of publishing datasets which offer enhanced opportunities for the re-use of such data in an automatic or semi-automatic way. These include for instance *nanopublications*, *cardinal assertions* and *knowlets*⁸. The potential of these new publishing approaches for discovering new facts as well as for comparing concepts and tracking how single concepts change are well described in [6].

The field of language resources offers us a rich array of highly structured kinds of datasets, structured according to a series of widely shared conventions (this is what makes the definition of models and vocabularies for lexica, corpora, etc, so viable in the first place) – something that would seem to lend itself well to making such resources FAIR in the machine-oriented spirit of the original description of those principles as well as to the new data publication approaches previously mentioned. However, the better and more expressive the underlying models are the more effective they will be.

⁶<https://www.w3.org/TR/prov-o/>

⁷For which one could use the Semantic Web ontology CRMInf<http://www.cidoc-crm.org/crminf/>.

⁸*Nanopublications* are defined as the "smallest possible machine readable graph-like structure that represents a meaningful assertion" [6] and consist of publishing a single subject-predicate-object triple with full provenance information; a generalisation of this idea is that of the *cardinal assertion* where a single assertion is associated with more than one provenance graph. A *knowlet* consists of a collection of multiple cardinal assertions, with the same subject concept [6] and can be viewed as locating that concept in a rich 'conceptual space'. For instance, this could be a cloud of predicates centered around a word or a sense.

In order to ensure the continued effectiveness of linked data and the Semantic Web in facilitating the creation of FAIR resources, it is vital that pre-existing vocabularies/models/data registries be re-used whenever possible in the modelling of user data; this of course also means ensuring that these models have sufficient coverage and defining extensions when this is not the case, as well as creating training materials suitable for different groups of users. Part of the intention of this article, together with the foundational work carried out in [7], is to provide an overview of what exists out there in terms of LLD-focused models, to look at the areas which are receiving most attention in order to highlight those which are so far underrepresented. In addition in **Section 3** we look at the most well known LLD models in the light of a recent series of recommendations on the publication of models as FAIR resources.

2.2. The Importance of Projects and Community Initiatives in LLD

One significant indicator of the success which LLD has had in the last few years is the variety of new funded projects which have included the publication of linguistic datasets as linked data as a core theme. These include projects at a continental or transnational level, notably European H2020 projects, ERCs and COST actions, as well as projects at the national and regional levels. Arguably, this recent success in obtaining project funding reflects a much wider recognition of the importance of linked data as a means of ensuring the interoperability and accessibility of language resources both to the research community and to a wider public. In addition, it also demonstrates the continuing maturation of the field as LLD continues to be applied to new domains and use cases, in many cases, within the context of the projects alluded to above. In addition, these projects also offer us clear examples of the use of the LLD vocabularies and models we will look at 'in the wild' so to speak and demonstrating their application to a wide number of medium to large scale datasets.

We have therefore decided to dedicate a section of the current article, **Section 5** to a detailed discussion of the current situation as regards research projects and LLD models and vocabularies. This includes a detailed overview of the area, **Section 5.1** along with an extended descriptions of a number of projects which we regard as the most significant from the point of view of LLD models and vocabularies. These are (in order of

1 appearance); the **Linked Open Dictionaries (LiODi)**
2 project (**Section 5.2.1**); the **Poetry Standardization**
3 **and Linked Open Data (POSTDATA)** project; the
4 **LiLa: Linking Latin** ERC project (**Section 5.2.4**); the
5 **Prêt-à-LLOD** project (**Section 5.2.5**); the **European**
6 **network for Web-centred linguistic data science**
7 (**NexusLinguarum**) COST action (**Section 5.2.6**). A
8 list of all the projects described in Section 5 can be
9 found in Table 3.

10 Note, however, that although the projects which we
11 will discuss in Section 5 have, in many cases, set the
12 agenda for the development of LLD models and vocabularies,
13 much of the actual work on the definition of these resources
14 was carried out – and is being carried out – within
15 community groups, such as the W3C OntoLex group. We
16 therefore include an update on community standards and
17 initiatives in **Section 4**. These include a subsection on
18 the latest activities in the OntoLex group (**Section 4.1**);
19 a discussion of recent work on LLD models for corpora
20 and annotation (**Section 4.2**); and similarly for what
21 concerns models and vocabularies for LLD resource
22 metadata (**Section 4.3**).

23 Section 5.1.2 features a discussion of the relationship
24 between community initiatives and projects.

25 2.3. The Relationship of LLD to the Digital 26 Humanities

27
28
29 Several of the projects which we will discuss in this
30 article are related to the area of Digital Humanities
31 (DH). This is the third major trend which we want
32 to highlight here, since it represents a move away (or
33 rather a branching off) from LLD's beginnings in
34 computational linguistics and natural language processing
35 (although these latter two still perhaps represent the
36 majority of applications of LLD), something that calls
37 for a shift in emphases in the definition and coverage
38 of LLD models. This overlap between LLD and DH is
39 particularly apparent in the modelling of corpora
40 annotation (**Section 4.2**) and in support for lexicographic
41 use cases (see **Section 4.1.1** and **Section 5.2.3**). Indeed
42 one obvious example of these shared concerns is the
43 publication of retro-digitised dictionaries as LLD
44 lexica (a major theme of the ELEXIS project, see
45 **Section 5.2.3**). The latter use case confronts us with
46 the challenge of formally modelling both the *content*
47 of a lexicographic work, that is the linguistic
48 descriptions which it contains, as well as those aspects
49 which pertain to it as a *physical text* to be represented
50 in digital form. In the latter case, this includes the
51 representation of (elements of) the *form* of the text, i.e., its structural

1 layout and overall visual appearance.⁹ In fact, as we
2 discuss in our description of the OntoLex Lexicography
3 module in **Section 4.1.1** even the structural division
4 of lexicographic works into textual units such as
5 entries and senses is not always isomorphic to the
6 representation of the lexical content of those units
7 using OntoLex-Lemon classes such as `LexicalEntry` and
8 `LexicalSense`.

9 We may also wish to model different aspects of the
10 history of the lexicographic work as physical text.¹⁰
11 All of this calls for a much richer provision of
12 metadata categories than had previously been considered
13 for LLD lexica, both at the level of the whole work
14 as well as at the level of the entry. It also requires
15 the capacity to model salient aspects of the same
16 artefact or resource at different levels of description
17 (something which is indeed offered by the OntoLex
18 Lexicography module, see **Section 4.1.1**). We discuss
19 metadata challenges in humanities use cases in **Section 4.3**.
20 A related topic is the relationship between notions
21 such as *word* from the lexical/linguistic and the
22 philological points of view and, more broadly speaking,
23 the relationship between linguistic and philological
24 annotations of text is a topic which is just starting
25 to gain attention within the context of LLD. It is
26 being studied both at the level of community initiatives
27 (see **Section 4.2**) as well as in projects such as
28 LiLa (see **Section 5.2.4**) as well as POSTDATA
29 (**Section 5.2.2**).

30 An additional series of challenges arises in the
31 consideration of resources for classical and historical
32 languages, or indeed, historical stages of modern
33 languages. For instance in the case of lexical
34 resources for historical languages we often come up
35 against the necessity of having to model attestations
36 (something that is discussed in **Section 4.1.3**) which
37 sometimes cite reconstructed texts, as well as the
38 desirability of being able to represent different
39 scholarly and philological hypotheses for instance
40 when it comes to modelling etymologies. The LiLa
41 project [9] (**Section 5.2.4** for a more detailed
42 description) provides a good exam-

43 ⁹Encompassing what the TEI dictionary chapter guidelines call
44 the typographical and editorial views. See [https://www.tei-c.org/
45 release/doc/tei-p5-doc/en/html/DI.html#DIMV](https://www.tei-c.org/release/doc/tei-p5-doc/en/html/DI.html#DIMV)

46 ¹⁰For example, in the case of older resources, annotating instances
47 where the content has been superseded by subsequent scholarly
48 work. Or we might want to track the evolution of a historically
49 significant lexicographic work over the course of a number of
50 editions, in order to see, for example, how changes in entries
51 reflected both linguistic and wider, non-linguistic trends. This
52 was in fact one of the motivations behind the Nénufar project [8],
53 described in Section 5.1.1.

ple of the challenges and opportunities of adopting the LLD model to represent linguistic (meta)data for both lexical and textual resources for a classical language (Latin).

One extremely important (non RDF-based) standard for encoding documents in the Digital Humanities is **TEI/XML**. In the current article we discuss the relationship between TEI and RDF-based annotation approaches in Section 4.2.1, and introduce the new lexicographic TEI-based standard **TEI Lex-0** and describe current work on a crosswalk between OntoLex-Lemon and the latter in **Section 5.2.3**.

Finally, see **Section 5.1.1** for an overview of a number of projects combining DH and LLD.

2.4. Related Work

The current work is intended, among other things, to both complement as well as to update a previous general survey on models for representing LLD, published by Bosque-Gil et al. in 2018 [7]. Although we are now only two years on from the publication of that article, we feel that enough has happened in the intervening time period to justify a new survey article. In addition we believe that we cover a much wider range of topics than the previous article and that our focus is also quite different. Broadly speaking, that previous work offered a classification of various different LLD vocabularies according to the different levels of linguistic description that they covered. The current paper however concentrates more on the use of LLD vocabularies in practise and on their availability (this is very much how we have approached the survey in **Section 3**). Moreover, the present article includes a detailed discussion of recent work in the use of LLD models and vocabularies in corpora and annotation, **Section 4.2**, as well as an extensive section on metadata, **Section 4.3**, neither of which were given the same detailed level of coverage in [7]. Additionally we also cover the following initiatives which were not discussed in [7] because they had not yet gotten underway:

- The development of new OntoLex-Lemon modules for morphology Section 4.1.2 and frequency, attestations, and corpus Information, described in **Section 4.1.3**
- An important new initiative in aligning LLD vocabularies for corpora and annotation, described in **Section 4.2.5**.

In what follows we will assume that the reader already has some grounding in linked data in general – includ-

ing a basic familiarity with the Resource Description Framework (RDF), RDF Schema (RDFS) and the Web Ontology Language (OWL) – and linguistic linked data in particular. The recently published *Linguistic linked data: representation, generation and applications* [10] should however give the interested reader who is missing this minimal background a comprehensive introduction to and overview of the latter field, focusing on more established models and vocabularies and their application rather than on recent developments. Another important new book on the topic of LLD and which has relevance to the current work is the collected volume *Development of linguistic linked open data resources for collaborative data-intensive research in the language sciences* [11] which aims to describe major developments since 2015. It consists mostly of position papers by researchers from the linguistics and the language resource communities.

3. LLD Models: An Overview

Summary The current section will give an overview of some of the most well known and/or widely used models and vocabularies in LLD. A summary of the models discussed in the current section (and in the whole article) can be found in **Tables 1 and 2** (with Table 1 dealing with published LLD models/vocabularies and 2 with models/vocabularies that are currently unavailable or no longer updated). An account of some the latest developments with regards to these models, on the other hand, can be found in Section 4.

We will classify each of the models described in this section according to the scheme given in the linguistic LOD cloud diagram¹¹ (the cloud itself is described in [12]), namely:

- Corpora (and Linguistic Annotations)(Section 3.1)
- Lexicons and Dictionaries (Section 3.2)
- Terminologies, Thesauri and Knowledge Bases (Section 3.3)
- Linguistic Resource Metadata (Section 3.4)
- Linguistic Data Categories (Section 3.5)
- Typological Databases (Section 3.6)

For each category we list the most prominent and/or widely used LLD models/vocabularies that belong to that category (the relevant section is given in paren-

¹¹<http://linguistic-lod.org/llod-cloud>

theses after each category in the list above). These models were either originally designed to help encode that kind of dataset or have been widely appropriated for that end; in the case of the category *Linguistic Data Categories* we list linked data linguistic data categories. For instance, the OntoLex-Lemon model falls under *Lexicons and Dictionaries* since it was initially conceived as a means of enriching ontologies with lexical information, that is, of lexicalising ontological concepts, but subsequently gained popularity as a means of encoding linked data lexica although it can also be used for modelling and publishing other kinds of datasets. Tables 1 and 2 give a summary of the LLD vocabularies and models covered in this paper (with the relevant sections of the article listed).

Below we describe our methodology for the rest of the section. In Section 3.7 we discuss tools and platforms for the publication of LLD.

Our Approach to Classification

This section is intended as an overview so we do not give a detailed description of single models. Several of these models are described in more detail in the rest of the article, or in the Appendix in the case of OntoLex-Lemon. Others can be found in the previous survey paper, [7]. Instead we will describe here them on the basis of a number of criteria many of which are related to their status as FAIR resources, and in particular to their status as FAIR models and vocabulary. In particular, we will refer to a recent draft survey on FAIR Semantics [13], the result of a dedicated brainstorming workshop of the FAIRsFAIR project.¹² This report outlined a number of recommendations and best practices for FAIR semantic artefacts where these are defined as "machine -actionable and -readable formalisation[s] of a conceptualisation enabling sharing and reuse by humans and machines" (the term includes: taxonomies, thesauri, ontologies).

From all of the recommendations listed in [13] we have selected the following subset on the basis of their salience to the set of models and vocabularies under discussion (with justifications for recommendations based on those given in [13]):

- (P-Rec 2) Ensure there is a separate URI for the metadata and that they are published separately; this helps in making the resource more findable and supports the extraction of this metadata.

- (P-Rec 4) Publish semantic artefacts and their contents in a semantic repository: in order to be able to exploit repository technologies for findability and re-use of semantic artefacts ;
- (P-Rec 6) Retrieval through search engines ;
- (P-Rec 10) Use a foundational ontology to align semantic artefacts (this enhances re-usability);
- (P-Rec 13) Create documented crosswalks and bridges
- (P-Rec 16) Ensure clear licensing of semantic artefacts.

To start with the recommendations (P-Rec 2), (P-Rec 4), and (P-Rec 10) have been followed by *none* of the models/vocabularies which we look at below. Following these three recommendations would, however, greatly help to make these resources (and the datasets they help to encode) more FAIR and we regard their adoption as desirable future objectives for the models and vocabularies listed below, bringing them into line with the latest thinking on making such kinds of resource FAIR.¹³ In terms of the recommendation (P-Rec 13) at the time of writing we can only mention ongoing efforts at developing a TEI Lex-0/OntoLex-Lemon crosswalk described in Section 5.2.3.

We will use (P-Rec 6) and (P-Rec 16) to help us to analyse the models and vocabularies to follow. For instance several of the models mentioned do exist on the Linked Open Vocabulary (LOV)¹⁴ search engine¹⁵ [15] and the DBpedia archive ontology archive.¹⁷ In cases where licensing information is available as machine actionable metadata, using properties like DCT:license and URI's such as <https://creativecommons.org/publicdomain/zero/1.0/> we will point this out as it enhances the re-usability of those resources.

In addition to the written descriptions of different LLD models given below, we also give a tabular summary of the most significant/stable/widely available¹⁸ of these models in Table 1. This also points, in relevant cases, to other parts of the sections of the pa-

¹³The adoption of foundational ontologies, for instance, might help to alleviate some of the problems raised by the proliferation of independently developments as described in [7].

¹⁴<https://lov.linkeddata.es/dataset/lov>

¹⁵Note that the LOV site provides a list of criteria for inclusion on their search engine [14]¹⁶

¹⁷<http://archivo.dbpedia.org/>

¹⁸Several of the models which are described in the rest of the section aren't available, at least anymore, but may be interesting for historical reasons.

¹²<https://www.fairsfair.eu/>

per where a more in-depth description of said model is given. Every one of the models listed in the table at is an OWL ontology.. We will also list the other vocabularies which they make use of (aside from the vocabularies OWL, RDF, and RDFS which are common to all of the vocabularies on the list). These include the well known ontologies/vocabularies: XML Schema Definition¹⁹ (XSD); the Friend of a Friend Ontology²⁰ (FOAF); the Simple Knowledge Organisation System²¹ (SKOS); Dublin Core²² (DC); Dublin Core Metadata Initiative (DCMI) Metadata Terms;²³ the Data Catalog Vocabulary²⁴ (DCAT), described also in Section 4.3; the PROV Ontology²⁵ (PROV-O).

In addition the table also mentions the following vocabularies.

- Activity Streams(AS): a vocabulary for activity streams.²⁶
- GOLD: an ontology for describing linguistic data, which is described in Section 3.5.
- MARL: a vocabulary for describing and annotating subjective opinions.²⁷
- ITSRDF: an ontology used within the Internationalization Tag Set.²⁸
- The Creative Commons vocabulary²⁹ (CC).
- VANN: a vocabulary for annotating vocabulary descriptions.³⁰
- SKOS-XL: an extension of SKOS with extra support for “describing and linking lexical entities”.³¹ SKOS and SKOS-XL are, along with *lemon* and its successor *OntoLex-Lemon*, amongst the most well known ways of enriching linked data taxonomies and conceptual hierarchies with linguistic information. We will look at the use of a SKOS-XL vocabulary in the context of a project on the classification of folk tales in Section 5.

¹⁹<https://www.w3.org/TR/xmlschema-0/>

²⁰<http://xmlns.com/foaf/spec/>

²¹<https://www.w3.org/2004/02/skos/>

²²<https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>

²³<https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>

²⁴<https://www.w3.org/TR/vocab-dcat-2/>

²⁵<https://www.w3.org/TR/prov-o/>

²⁶<https://www.w3.org/TR/activitystreams-vocabulary/>

²⁷<http://www.gsi.dit.upm.es/ontologies/marl/>

²⁸<https://www.w3.org/TR/its20/>

²⁹<https://creativecommons.org/ns/>

³⁰<https://vocab.org/vann/>

³¹<https://www.w3.org/TR/skos-reference/skos-xl.html>

3.1. Vocabularies and Models for Corpora and Linguistic Annotations

Linguistic annotation, e.g. for digital editions, corpora, and linking texts with external resources has long been a topic of interest in the context of RDF and linked data. Coexisting with relational databases, XML-based formats (most notably, TEI, see 4.2) or simply text-based formats, RDF-based annotation models have been steadily undergoing development and are increasingly being used in research and industry. Currently, there are two primary RDF vocabularies widely used for text annotations: **NLP Interchange Format** (NIF),³² used mostly in the language technology sector and **Web Annotation**,³³ formerly known as *Open Annotation* (abbreviated here as OA), used in digital humanities, life sciences and bioinformatics. Both models have their advantages and shortcomings, and a number of proposals to extend these have been proposed. Most importantly, there is a need for synchronization between the two. Both are available in LOV³⁴ and *archivo*³⁵ (the NIF core in the case of NIF³⁶). The Web Annotation model, although it is covered by a W3C software and document notice and license, does not express this information in the form of triples in the resource metadata; NIF on the other hand does express licensing information as machine actionable metadata.

More details about both models and their recent developments are described in Section 4.2. Other vocabularies described in that section include POWLA, CoNLL-RDF and Ligt. The first of these, POWLA,³⁷ is available on *archivo*,³⁸ the only one of the three to be so available. CoNLL-RDF³⁹ has version info as a string using the owl:versionInfo property and is covered by a CC-BY 4.0 license as specified in the LICENSE.data page.⁴⁰

³²<https://nif.readthedocs.io/en/latest/>

³³<https://www.w3.org/TR/annotation-model/>

³⁴<https://lov.linkeddata.es/dataset/lov/vocabs/nif> and <https://lov.linkeddata.es/dataset/lov/vocabs/oa>

³⁵<http://archivo.dbpedia.org/info?o=http://www.w3.org/ns/oa>

³⁶<http://archivo.dbpedia.org/info?o=http://persistence.uni-leipzig.org/nlp2rdf/ontologies/nif-core>

³⁷<http://purl.org/powla/powla.owl>

³⁸<https://archivo.dbpedia.org/info?o=http://purl.org/powla/powla.owl>

³⁹<http://purl.org/acoli/conll#>

⁴⁰<https://github.com/acoli-repo/conll-rdf/blob/master/LICENSE.data.txt>

| Summary | | | | | |
|-----------------------------------|---|--|--|--|---|
| Name | Other Vocabularies/Models Used | LLO Category | Licenses | Versions (at time of writing 26/07/21) | Extended Coverage in Current Article |
| OntoLex-Lemon | CC, DC, FOAF, SKOS, XSD | Lexicons and Dictionaries | CC0 1.0 | Version 1.0, 2016 (but this is closely based on the prior <i>lemon</i> model [16]) | Section 3.2, Section 4.3.3 and Appendix A |
| Lexicog (OntoLex-Lemon) | DC, LexInfo, SKOS, VOID, XSD | Lexicons and Dictionaries | CC0 | Version 1.0, (2019-03-08) | Section 3.2 and Section 4.1.1 |
| MMoOn | DC, FOAF, GOLD, LexVo, OntoLex-Lemon, SKOS, XSD | Terminologies, Thesauri and KBs (Morphology) | CC-BY 4.0 | Version 1.0, 2016 | Section 3.3 |
| Web Annotation Data Model (OA) | AS, FOAF, PROV, SKOS, XSD | Corpora and Linguistic Annotations | W3C Software and Document Notice and License | Version "2016-11-12T21:28:11Z" | Section 3.1 and Section 4.2.3 |
| NLP Interchange Format (NIF Core) | DC, DCTERMS, ITS RDF, levont, MARL, OA, PROV, SKOS, VANN, XSD | Corpora and Linguistic Annotations | Apache 2.0 and CC-BY 3.0 | Version 2.1.0 | Section 3.1 and Section 4.2.2 |
| POWLA | FOAF, DC, DCT, | Corpora and Linguistic Annotations | NA | Last Updated 2018-04-03 | Section 4.2 |
| CoNLL-RDF | DC, NIF Core, XSD | Corpora and Linguistic Annotations | Apache 2.0 and CC-BY 4.0 | Last Updated 2020-05-26 | Section 4.2.4 |
| Ligt | DC, NIF Core, OA | Corpora and Linguistic Annotations | NA | Version 0.2 (2020-05-26) | Section 4.2.4 |
| META-SHARE | CC, DC, DCAT, FOAF, SKOS, XSD | Linguistic Resource Metadata | CC-BY 4.0 | Version 2.0 (pre-release) | Section 3.4 and Section 4.3.2 |
| OLiA | DCT, FOAF, SKOS | Linguistic Data Categories | CC-BY-SA 3.0 | Version last updated 27/02/20 | Section 3.5 |
| LexInfo | CC, Ontolex, TERMS, VANN | Linguistic Data Categories | CC-BY 4.0 | Version 3.0, 14/06/2014 | Section 3.5 |
| LexVo | FOAF, SKOS, SKOSXL, XSD | Typological Databases | CC-BY-SA3.0 | Version 2013-02-09 | Section 3.6 |

Table 1
Summary of published LLD vocabularies

3.2. Lexicons and Dictionaries

The most well known model for the creation and publication of lexica and dictionaries as linked data is the **OntoLex-Lemon model**⁴¹ [17], an output of

⁴¹The URI for OntoLex-Lemon is: <http://www.w3.org/ns/lemon/ontolex> and the OntoLex-Lemon guidelines can be found at <https://www.w3.org/2016/05/ontolex/>.

the W3C ontolex working group which manages its ongoing development and further extension (see Appendix A for an introduction to the model with examples and Section 4.1 for extensions and further developments). It is based on a previous model, the **Lexicon Model for ONtologies (*lemon*)** [16]. Like its predecessor, OntoLex-Lemon was designed with the intention of enriching ontologies with linguistic information and not of modelling dictionaries and lexicons

| Summary | | | |
|---------------------------|------------------------------------|--------------------------------------|--------------------------------------|
| Name | LLO Category | Status (at time of writing 26/07/21) | Extended Coverage in Current Article |
| OntoLex-Lemon: FrAC | Lexicons and Dictionaries | Under Development | Section 4.1.3 |
| OntoLex-Lemon: Morphology | Lexicons and Dictionaries | Under Development | Section 4.1.2 |
| PHOIBLE | Terminologies, Thesauri and KBs | Unavailable | Section 3.3 |
| FRED | Corpora and Linguistic Annotations | Project Specific Vocabulary | Section 4.2 |
| NAF | Corpora and Linguistic Annotations | Project Specific Vocabulary | Section 4.2 |
| GOLD | Linguistic Data Categories | No Longer Updated | Section 3.5 |

Table 2

Other LLD Vocabularies Discussed in this Paper

per se. Thanks to its popularity however, it has come to take on the status of a de facto standard for the modelling and codification of lexical resources in RDF (including, for instance, retrodigitized dictionaries and wordnets) in general. Resources which have been modelled using OntoLex-Lemon include: the LLD version of the Princeton Wordnet,⁴² DBnary (the linked data version of Wiktionary) [18], and the massive multilingual knowledge graph Babelnet [19]. The OntoLex-Lemon model is modular and consists of a core module along with modules for *Syntax and Semantics*,⁴³ *Decomposition*,⁴⁴ and *Variation and Translation*,⁴⁵ as well as a dedicated metadata module, *lime*⁴⁶ (all of these modules are described in Appendix A, except for *lime* which is described in Section 4.3.3).

OntoLex-Lemon is available on LOV as its predecessor *lemon*.⁴⁷ All of its separate modules are listed separately however:⁴⁸ the *core*,⁴⁹ *lime*,⁵⁰ *vartrans*,⁵¹

⁴²<http://wordnet-rdf.princeton.edu/about>

⁴³<http://www.w3.org/ns/lemon/synsem>

⁴⁴<http://www.w3.org/ns/lemon/decomp>

⁴⁵<http://www.w3.org/ns/lemon/vartrans>

⁴⁶<http://www.w3.org/ns/lemon/lime>

⁴⁷<https://lov.linkeddata.es/dataset/lov/vocabs/lemon>

⁴⁸See the Appendix for a description of each, aside from *lime* described in Section 4.3.3

⁴⁹<https://lov.linkeddata.es/dataset/lov/vocabs/ontolex>

⁵⁰<https://lov.linkeddata.es/dataset/lov/vocabs/lime>

⁵¹<https://lov.linkeddata.es/dataset/lov/vocabs/vartrans>

synsem,⁵² the *decomp* module.⁵³ Three of its modules are available on *archivo*, the *core*:⁵⁴ the *lime* metadata module⁵⁵ and the *Variation and Translation* module.⁵⁶ All of the OntoLex modules have their licenses (CC 1.0) described with RDF triples using the CC vocabulary⁵⁷ with a URI as an object. Version information is described using *owl:versionInfo*.

The OntoLex-Lemon Lexicography module⁵⁸ (described in more detail in Section 4.1.1) was published separately from OntoLex-Lemon. It is not available on LOV yet, however it is available on *archivo*.⁵⁹ The license (CC-Zero) is described with RDF triples using the CC vocabulary⁶⁰ and DC⁶¹ with a URI as an object. Version information is described using *owl:versionInfo*.

⁵²<https://lov.linkeddata.es/dataset/lov/vocabs/synsem>

⁵³<https://lov.linkeddata.es/dataset/lov/vocabs/lexdcp>

⁵⁴<https://archivo.dbpedia.org/info?o=http://www.w3.org/ns/lemon/ontolex>

⁵⁵<http://archivo.dbpedia.org/info?o=http://www.w3.org/ns/lemon/lime>

⁵⁶<http://archivo.dbpedia.org/info?o=http://www.w3.org/ns/lemon/vartrans>

⁵⁷Using the *cc:license* property

⁵⁸The guidelines for the module can be found at <https://www.w3.org/2019/09/lexicog/>, the URL for the module is at <http://www.w3.org/ns/lemon/lexicog#>

⁵⁹<https://archivo.dbpedia.org/info?o=http://www.w3.org/ns/lemon/lexicog>

⁶⁰Using the *cc:license* property

⁶¹using *dc:rights*

3.3. Vocabularies for Terminologies, Thesauri and Knowledge Bases

The **Simple Knowledge Organisation System (SKOS)** is a W3C recommendation for the creation of terminologies and thesauri, or more broadly speaking, knowledge organisation systems.⁶² We will not go into any depth into it here since it is a general purpose vocabulary which is applied well beyond the domain of language resources.

In terms of specialised vocabularies or models for the modelling of linguistic knowledge bases – and aside from linguistic data category registries which will be discussed in Section 3.5 – we can list two. The first is **MMoOn ontology**⁶³ which was designed for the creation of detailed morphological inventories [20]. It does not currently seem to be available on any semantic repositories/archives/search engines but it does have its own dedicated website⁶⁴ which offers a SPARQL endpoint (although this was down at the time of writing). Its license information (it has a CC-BY 4.0 license) is available as triples using `dct:license` with a URI as an object.

PHOIBLE is an RDF model for creating phonological inventories [7]. As of the time of writing, PHOIBLE data was no longer available as a complete RDF graph, but only in its native (XML) format from which RDF fragments are dynamically generated. The original data remains publicly available,⁶⁵ but on the PHOIBLE website, it is only possible to browse and export selected content into RDF/XML.⁶⁶ Since it no longer provides resolvable URIs for its components, PHOIBLE data does not fit within the narrower scope of LLD vocabularies anymore. It does, however, maintain a non-standard way of linking, as it has been absorbed into the Cross-Linguistic Linked Data infrastructure [21, CLLD] (along with other resources from the typology domain). CLLD datasets and their RDF exports continue to be available as open data under <https://clld.org/>, see below for additional details, Section 3.6.

⁶²<https://www.w3.org/2004/02/skos/>

⁶³<https://github.com/MMoOn-Project/MMoOn/blob/master/core.ttl>

⁶⁴<https://mmoon.org/>

⁶⁵<https://github.com/clld/phoible/tree/master/phoible/static/data>

⁶⁶See, for example, <https://phoible.org/inventories/view/161>.

3.4. Linguistic Resource Metadata

Due to the importance of the topic we give a much fuller overview in Section 4.3; here we will only look at accessibility issues for the two models for language resource metadata which we mention there. These are the METASHARE ontology⁶⁷ and *lime*. The latter has been previously introduced and is described in more detail in Section 4.3.3. The former is currently in its pre-release version 2.0 (the last update being 2020-03-20). Its license information (it has a CC-BY 4.0 license) is available as triples using `dct:license` with a URI as an object.

3.5. Linguistic Data Categories

History

As of 2010, two major repositories were in widespread use by different communities for addressing the harmonization and linking of linguistic resources via their data categories. In computational lexicography and language technology, the most widely applied terminology repository was **ISOcat** [22] which provided human-readable and XML-encoded information about linguistic data categories that were relevant for linguistic annotation, the encoding of electronic dictionaries and language resource metadata via persistent URIs.

In the field of language documentation and typology, the **General Ontology of Linguistic Description (GOLD)** emerged in the early 2000s [23], having been originally developed in the context of the project Endangered Metadata for Endangered Languages Data (E-MELD, 2002-2007).⁶⁸ GOLD stood out in particular because of its excellent coverage of low resource languages. In the RELISH project, a curated mirror of GOLD-2010 was incorporated into ISOcat [24]. Unfortunately, since then, GOLD development has stalled and, while the resource is still being maintained by the LinguistList (along with the data from related projects) and still remains accessible,⁶⁹ it has not been updated since [25] (and for this reason we have not included it in our summary table). In parts, its function seems to have been taken over by ISOcat, but it is worth pointing out here that the ISOcat registry exists only as a static, archived resource, but no longer as an operational system.

⁶⁷<http://www.meta-share.org/ontologies/meta-share/meta-share-ontology.owl/documentation/index-en.html>

⁶⁸<http://emeld.org/>

⁶⁹<https://linguistlist.org/projects/gold.cfm>

1 The Current Situation

2 The ‘official’ successor of ISOcat, the CLARIN
3 Concept Registry is briefly discussed in Section 4.3
4 below (it is not strictly speaking a linked data vocab-
5 ulary). Another one of its successors is the **Lex-**
6 **Info ontology**,⁷⁰ the data category register used in
7 OntoLex-Lemon and which re-appropriates many of
8 the concepts which were contained in ISOcat for use
9 within the lexical domain (dictionaries, terminologies,
10 lexica). Currently in its third version, LexInfo can
11 be found both on the LOV search engine⁷¹ and on
12 *archivo*,⁷² it appears both times however in its second
13 version. Version 3.0 is under development since late
14 2019 in a community-guided process via GitHub,⁷³
15 and is not registered with either service, yet. LexInfo
16 has a (CC-BY 4.0) license. This is described with RDF
17 triples using the CC vocabulary and DCT with a URI
18 as an object in both cases. Version information is de-
19 scribed using owl:versionInfo.

20 For linguistic data categories in linguistic annota-
21 tion (of corpora and by NLP tools), a separate termi-
22 nology repository exists with the **Ontologies of Lin-**
23 **guistic Annotation** [26, OLiA].⁷⁴ OLiA has been de-
24 veloped since 2005 in an effort to link community-
25 maintained terminology repositories such as GOLD,
26 ISOcat or the CLARIN Concept Registry with an-
27 notation schemes and domain- or community-specific
28 models such as LexInfo or the Universal Depend-
29 encies specifications by means of an intermediate “Ref-
30 erence Model”. OLiA consists of a set of modular, in-
31 terlinked ontologies and is designed as a native linked
32 data resource. Its primary contributions are to provide
33 machine-readable documentation of annotation guide-
34 lines and a linking with and among other terminology
35 repositories. It has been suggested that such a collec-
36 tion of linking models, developed in an open source
37 process via GitHub, may be capable of circumvent-
38 ing some of the pitfalls of earlier, monolithic solutions
39 of the ISOcat era [27]. At the moment, OLiA cov-
40 ers annotation schemes for more than 100 languages,
41 for morphosyntax, syntax, discourse and aspects of se-
42 mantics and morphology. OLiA has a (CC-BY 4.0) li-
43 cense; this is described using the Dublin Core property
44 license with a URI as an object.

46 ⁷⁰<https://lexinfo.net/>

47 ⁷¹<https://lov.linkeddata.es/dataset/lov/vocabs/lexinfo>

48 ⁷²[http://archivo.dbpedia.org/info?o=http://www.lexinfo.net/](http://archivo.dbpedia.org/info?o=http://www.lexinfo.net/ontology/2.0/lexinfo)
49 [ontology/2.0/lexinfo](http://www.lexinfo.net/ontology/2.0/lexinfo)

50 ⁷³It will be the first version that is compliant with OntoLex-
51 Lemon.

⁷⁴<http://purl.org/olia>

3.6. Vocabularies for Typological Datasets

Relevant Resources and Initiatives

4 Linguistic typology is commonly defined as the
5 field of linguistics that studies and classifies languages
6 based on their structural features [28]. The field of lin-
7 guistic typology has natural ties with language docu-
8 mentation, and accordingly, considerable work on lin-
9 guistic typology and linked data has been conducted in
10 the context of the GOLD ontology (see above, Section
11 3.5). We can identify the following relevant datasets.

12 One of the main contributors and advisors to the sci-
13 entific study of typology is the **Association for Lin-**
14 **guistic Typology (ALT)**.⁷⁵ They facilitate the descrip-
15 tion of the typological patterns underlying datasets.
16 One of the most well-known resources that ALT makes
17 available is the **World Atlas of Language Structures**
18 **(WALS)**⁷⁶ [29, 30] which is a large database of phono-
19 logical, grammatical, and lexical properties of lan-
20 guages gathered from descriptive materials. This re-
21 source can both be used interactively online and can be
22 downloaded. The **CLLD**⁷⁷ (**Cross-Linguistic Linked**
23 **Data**) project integrates WALS, thus, offering a frame-
24 work that structures this typological dataset using the
25 Linked Data principles.

26 Another collection that provides web-based access
27 to a large collection of typological datasets is the **Ty-**
28 **poological Database System (TDS)** [31, 32]. The main
29 goals of TDS are to offer users a linguistic knowledge
30 base and content metadata. The knowledge base in-
31 cludes a general ontology and dictionary of linguis-
32 tic terminology, while the metadata describes the con-
33 tent of the term ontology databases. TDS supports a
34 unified querying across all the typological resources
35 hosted with the help of an integrated ontology. The
36 **Clarín Virtual Language Observatory (VLO)**⁷⁸ in-
37 corporates TDS among its repositories.

38 Finally, another group of datasets relevant for typo-
39 logical research include large-scale collections of lex-
40 ical data, as provided, for example by **PanLex**⁷⁹ and
41 **Starling**.⁸⁰ An early RDF edition of PanLex has been
42 described by [33] and was incorporated in the initial
43 version of the Linguistic Linked Open Data cloud dia-
44 gram. At the time of writing, however, this early RDF

46 ⁷⁵<https://linguistic-typology.org/>

47 ⁷⁶<https://wals.info/>

48 ⁷⁷<https://clld.org/>

49 ⁷⁸<https://vlo.clarin.eu/>

50 ⁷⁹<http://panlex.org>

51 ⁸⁰<https://starling.rinet.ru/>

version does not seem to be accessible anymore. Instead, CSV and JSON dumps are being provided from the PanLex website. On this basis [34] describe a fresh OntoLex-Lemon edition of PanLex (and other) data as part of the **ACoLi Dictionary Graph**.⁸¹ However, they currently do not provide resolvable URIs, but rather redirect to the original PanLex page. The authors mention that linking would be a future direction, and in preparation for this, they provide a **TIAD-TSV** edition of the data along with the OntoLex edition, with the goal to adapt techniques for lexical linking developed in the context of, for example, the on-going series of Shared Tasks on Translation Inference Across Dictionaries.⁸² As for modelling requirements of lexical datasets in linguistic typology, these are not fundamentally different from other forms of lexical data, but they adopt OntoLex, resp. its predecessor, *lemon*, see above. They do, however, require greater depth with respect to identifying and distinguishing language varieties. This was one of the driving forces behind the development of Glottolog.

Vocabularies for Typological Datasets

In terms of linked data vocabularies and models which are relevant for the creation of typological databases we can identify **LexVo**⁸³ [35]. LexVo bridges the gap between linguistic typology and the LOD community and brings together language resources and the entity relationships provided through the Linked Data Web and the Semantic Web. The project manages to link a large variety of resources on the Web, besides providing global IDs (URIs) for language-related objects. LexVo is available on *archivo*⁸⁴ but is not yet available on LOV. Further discussion of LexVo can be found in Section 4.3.4

3.7. Excursus: Tools and Platforms for the Publishing of LLD

The availability of tools and platforms for the editing, conversion and publication of LLD resources (on the basis of the models which we discuss in this article) is important for the adoption of those models amongst a wider community of end users. It is especially important for users who are unfamiliar with the technical details of linked data and the Semantic Web and who

are yet highly motivated to create and/or make use of linked data resources. Such software is especially helpful when it comes to the validation and post-editing of language resources which have been generated automatically or semi-automatically. They also help to make the information in these resources accessible for those who may not yet have learned SPARQL in order to browse or carry out (simple) queries on them.

In terms of existing tools or software which offer dedicated provision for the models which we look at in this article, we can mention **VocBench** and **LexO** for OntoLex-Lemon. Both of these are web-based and allow for the collaborative development of computational resources. In the case of the **VocBench** platform, currently in its third release [36], these resources can be OWL ontologies and SKOS thesauri as well OntoLex-Lemon lexicons. LexO focuses on lexical and ontological resources and was originally developed in the context of DitMaO a project on the medico-botanical terminology of the Old Occitan language [37]. A first, generally available version of LexO, LexO-lite, is available at the <https://github.com/andreabellandi/LexO-lite>.

Finally, we should mention **LLODifier**⁸⁵ a suite of tools for creating and working with LLD which is currently being developed by the Applied Computational Linguistics Lab of the Goethe University Frankfurt. These include **vis** for working with NIF and **unimorph** which works with CoNLL-RDF.

4. An Overview of Developments in LLD Community Initiatives and Standards

Summary and Overview The current section comprises an extensive overview of developments in various different LLD community initiatives and standards relating to LLD models and vocabularies. In particular, it focuses on the three areas that we believe have been the most active in the last few years (the first two of the following) or that are starting to gain greater prominence (the third): lexical resources (Section 4.1), annotation and corpora (Section 4.2), and finally metadata (Section 4.3). We have referred to these as community standards/initiatives because they have been pursued or developed as community efforts rather than within a single research group or project.

Membership in these community initiatives is (often) open to all, rather than being limited to members

⁸¹Data available under <https://github.com/acoli-repo/acoli-dicts>.

⁸²<https://tiad2021.unizar.es/>

⁸³<http://lexvo.org/>

⁸⁴<http://archivo.dbpedia.org/info?o=http://lexvo.org/ontology>

⁸⁵<https://github.com/acoli-repo/LLODifier>

1 of a project or experts nominated by a standards body.
 2 The intention being to allow for the participation of a
 3 wider range of stakeholders as well as encouraging the
 4 collection of a wider variety of use-cases than might
 5 otherwise be possible.

6 One of the most notable community efforts in the
 7 context of LLOD (that is Linguistic Linked Open
 8 Data) is the Open Linguistics Working Group (OWLG)
 9 of Open Knowledge International⁸⁶ that introduced
 10 the vision of a Linguistic Linked Open Data (LLOD)
 11 cloud in 2011 [38], and whose activities, most no-
 12 tably the organization of the long-standing series of
 13 international Workshops on Linked Data in Linguis-
 14 tics (LDL, since 2012), as well as the publication
 15 of the first collected volume on the topic of Linked
 16 Data in Linguistics [39], ultimately led to the imple-
 17 mentation of LLOD cloud in 2012, celebrated with
 18 a special issue of the Semantic Web Journal pub-
 19 lished in 2015 [40]. The LLOD cloud, now hosted under
 20 <http://linguistic-lod.org/>, was enthusiastically em-
 21 braced, the linguistics category became a top-level cat-
 22 egory in the 2014 LOD cloud diagram, and since 2018,
 23 it represented the first LOD domain sub-cloud.

24 At the same time, a number of more specialized ini-
 25 tiatives emerged, as mentioned below, for which the
 26 Open Linguistics Working Group acted and continues
 27 to act as an umbrella that facilitates information ex-
 28 change among them and between them and the broader
 29 circles of linguists interested in linked data technol-
 30 ogies and knowledge engineers interested in language.
 31 Currently, main activities of the OWLG are the orga-
 32 nization of workshops on Linked Data in Linguistics
 33 (LDL), the coordination of datathons such as Multi-
 34 lingual Linked Open Data for Enterprises (MLODE
 35 2012, 2013) and the Summer Datathon in Linguis-
 36 tic Linked Open Data (SD-LLOD, 2015, 2017, 2019),
 37 maintaining the Linguistic Linked Open Data (LLOD)
 38 cloud diagram⁸⁷ and continued information exchange
 39 via mailing list⁸⁸

40 Over the years, however, the focus of discussion
 41 moved from the OWLG to more specialized mailing
 42 lists and communities. At the time of writing, partic-
 43 ularly active community groups concerned with data
 44 modelling include

- the W3C Community Group Ontology-Lexica,⁸⁹ 1
 originally working on ontology lexicalization, the 2
 group extended their activities after the publica- 3
 tion of the OntoLex vocabulary (May 2016) and 4
 now represents the main locus to discuss the mod- 5
 elling of lexical resources with web standards and 6
 in LL(O)D. 7
- the W3C Community Group Linked Data for 8
 Language Technology,⁹⁰ with a focus on lan- 9
 guage resource metadata and linguistic annota- 10
 tion with W3C standards 11

12 Most recently, these activities have converged in
 13 funded networks, especially, the Cost Action NexusLin-
 14 guarum. 14

15 Also, Linked Data plays a certain role in the context
 16 of older standardization initiatives, e.g., the TEI Con-
 17 sortium,⁹¹ or the ISO TC37/SC4 committee.⁹² 17

18 We take the standards and initiatives proposed by
 19 these communities as our basis of the topics in this sec-
 20 tion, but we will also look at significant developments
 21 respecting these standards and initiatives outside and
 22 independent of these groups (see Section 4.1.4) in the
 23 interests of completeness and to understand current
 24 trends. 24

25 Note that our intention has been to go for complete-
 26 ness in our description of these developments. At the
 27 same time, however, as has been mentioned we have
 28 tried not to include too much material that was already
 29 available in the sources cited in Section 2.4. 29

30 In addition, a discussion of the relationship between
 31 community initiatives and projects can be found in Sec-
 32 tion 5.1.2 below. 30

31 4.1. Lexical Resources: OntoLex-Lemon and its 32 Extensions 31

33 *Summary* In this section we will describe some of
 34 the most recent work that has been carried out on the
 35 OntoLex-Lemon model,⁹³ both within the ambit of the
 36 W3C Ontolex group as well as outside of it. With re-
 37 gards to the former we will discuss three of the lat-
 38 est extensions to the model (one of which has been
 39 published and two which are still currently under de-
 40 velopment) in Sections 4.1.1, 4.1.2, and 4.1.3. In Sec-
 41 tion 4.1.4 we look at a number of new extensions to
 42 41

43 ⁸⁶<https://linguistics.okfn.org/>

44 ⁸⁷<http://linguistic-lod.org/>

45 ⁸⁸Since early 2020, the mailing list operates via <https://groups.google.com/g/open-linguistics>. Earlier messages are archived under
 46 <https://lists-archive.okfn.org/pipermail/open-linguistics/>.

47 ⁸⁹<https://www.w3.org/community/ontolex/>

48 ⁹⁰<https://www.w3.org/community/ld4lt>

49 ⁹¹<https://tei-c.org/>

50 ⁹²<https://www.iso.org/committee/297592.html>

51 ⁹³An introduction to the model is given in Appendix A

1 OntoLex-Lemon which have emerged independently
2 of the W3C Ontolex group over the last two years and
3 which moreover have not been discussed in [7] (for
4 an in-depth discussion of such developments prior to
5 2018 please refer to the latter paper).

6 Note that the use of OntoLex-Lemon in a number of
7 different projects is described in Section 5.

8 4.1.1. The OntoLex Lexicography Module (lexicog)

9 As mentioned previously, *lemon* and its successor
10 OntoLex-Lemon have been widely adopted for the
11 modelling and publishing of lexica and dictionaries as
12 linked data. The core module has proven to be reason-
13 ably effective in capturing some of the most typi-
14 cal kinds of *lexical* information contained in dictio-
15 naries and lexical resources in general (e.g., [41–45]).
16 However, there are certain fairly common situations in
17 which the model falls short, most notably in the repre-
18 sentation of certain elements of dictionaries and other
19 lexicographic datasets [46]. This, however, is not sur-
20 prising given that *lemon* was initially conceived as a
21 model for grounding ontologies with linguistic infor-
22 mation and not for modelling lexical resources per se.

23 In order to adapt OntoLex-Lemon to the modelling
24 necessities and particularities of dictionaries and other
25 lexicographic resources, the W3C Ontolex commu-
26 nity group developed a new **OntoLex Lexicography
27 Module (lexicog)**.⁹⁴ This module was the result of
28 collaborative work with contributions from lexicogra-
29 phers, computer scientists, dictionary industry practi-
30 tioners, and other stakeholders and was first released
31 in September 2019. As stated in the specification, the
32 *lexicog* module “overcome[s] the limitations of *lemon*
33 when modelling lexicographic information as linked
34 data in a way that is agnostic to the underlying lexico-
35 graphic view and minimises information loss”.

36 The idea is to keep purely lexical content separate
37 from lexicographic (textual) content. For that purpose,
38 new ontology elements have been added that reflect
39 the dictionary structure (e.g., sense ordering, entry hi-
40 erarchies, etc.) and complement the OntoLex-Lemon
41 model. The *lexicog* module have been validated with
42 real enterprise-level dictionary data [47] and can be
43 considered in a stable status right now.

44 In *lexicog* the structural organisation of a lexico-
45 graphic resource is now associated with the class Lexi-
46 cographic Resource (a subclass of the VoID⁹⁵ class
47 Dataset) whereas the lexical content is (as previously)

1 associated with the *lime* class Lexicon (see Section
2 4.3.3). The former is described as representing “a col-
3 lection of lexicographic entries[...]in accord with the
4 lexicographic criteria followed in the development of
5 that resource”.⁹⁶ These lexicographic entries are rep-
6 resented in their turn by another new *lexicog* class,
7 namely, the class Entry, which is defined as being a
8 “structural element that represents a lexicographic arti-
9 cle or record as it arranged in a source lexicographic
10 resource”⁹⁷ (emphasis ours). Furthermore an Entry is
11 related to its source Lexicographic Resource via the
12 object property entry. The class Entry is a subclass
13 of the more general class Lexicographic Component,
14 defined as “a structural element that represents the
15 (sub-)structures of lexicographic articles providing in-
16 formation about entries, senses or sub-entries”, mem-
17 bers of this class “can be arranged in a specific order
18 and/or hierarchy”.⁹⁸ That is, Lexicographic Component
19 allows for the representation of the ordering of senses
20 in an entry (and even potentially entries if this is re-
21 quired), the arrangement of senses and sub-senses in
22 a hierarchy, etc. in a published lexicographic resource
23 (by making use of the classes and properties we have
24 looked at above, along with the *lexicog* object property
25 subComponent), separately from the representation of
26 the same resource as lexical content. Finally,⁹⁹ we need
27 some way of linking together these two levels of repre-
28 sentation. This is provided by the *lexicog* object prop-
29 erty describes which relates individuals of class Lexi-
30 cographic Component, which belong to a specific lex-
31 icographic resource, “to an element that represents the
32 actual information provided by that component in the
33 lexicographic resource”.¹⁰⁰ See Figure 1.

34 As an example we can take the entry for the Ital-
35 ian word *chiaro* ‘clear’ in the popular Italian dictionary
36 *Treccani*.¹⁰¹ The latter lists the word as both an adject-
37 ive and a masculine noun, as well as having an adverb
38 as a subcomponent of the entry, in addition to listing
39 the related adverb *chiaramente* ‘clearly’ as a related
40 entry as well as the diminutive *chiarretto*.

41 The first two of the (four) subsenses of the entry are
42 classed as adjectives, the third as a noun, and the fourth

94 <https://www.w3.org/2019/09/lexicog/>

95 <https://www.w3.org/TR/void/>

96 <https://www.w3.org/2019/09/lexicog/#lexicographic-resource>

97 <https://www.w3.org/2019/09/lexicog/#Entry>

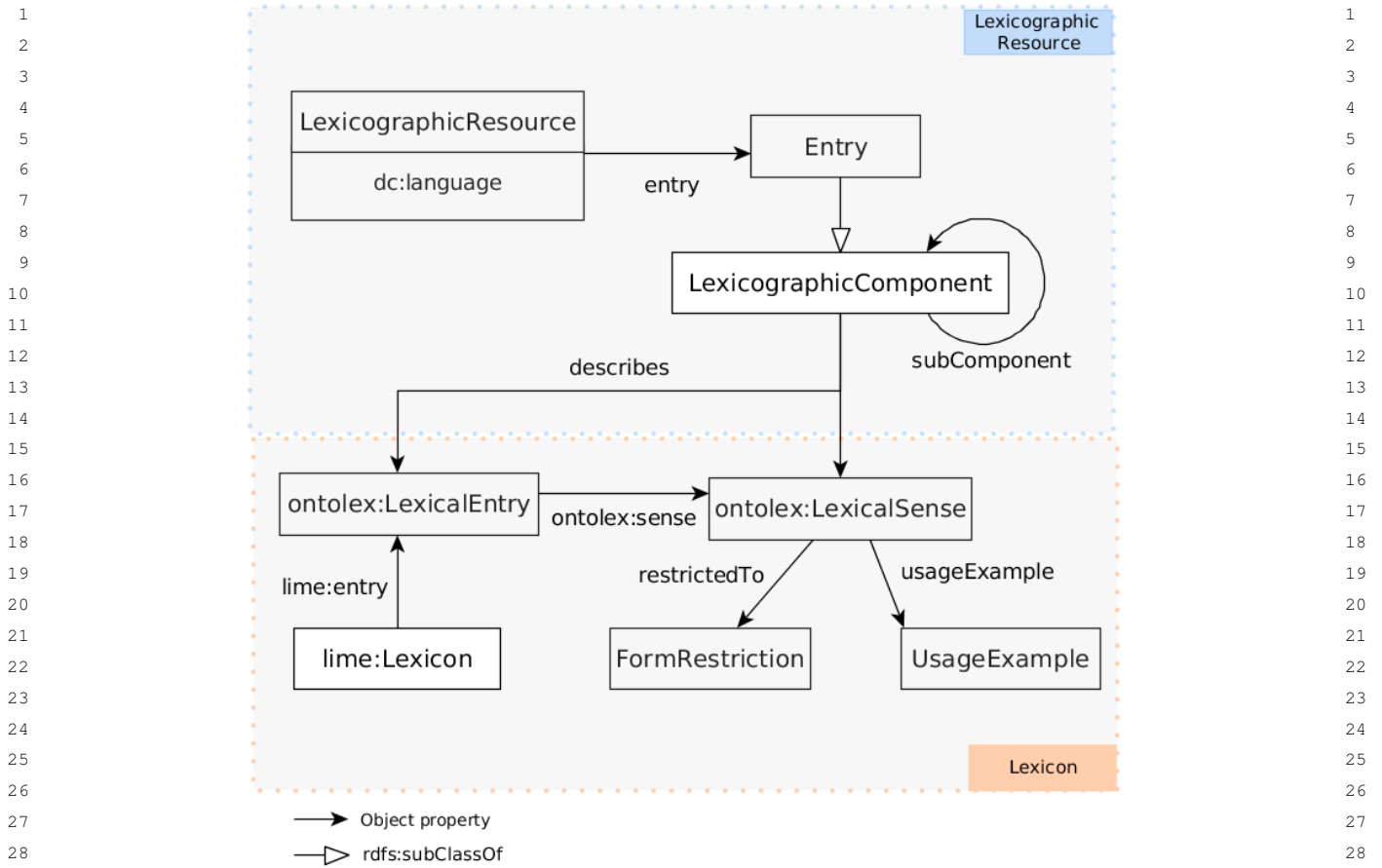
98 <https://www.w3.org/2019/09/lexicog/>

#lexicographic-component

99 Note that we do not cover all the classes and properties in the module in this brief description, please see the guidelines for a comprehensive description with examples.

100 <https://www.w3.org/2019/09/lexicog/#describes-0>

101 <https://www.treccani.it/vocabolario/chiaro/>

Fig. 1. The *lexicog* module (taken from the guidelines).

as an adverb. We will simplify this for the purposes of exposition by assuming that the first subsense is an adjective, the second a noun, and the third an adverb. This can be represented in the following way. First we represent the encoding of the Treccani dictionary structure itself, and the different sub-components of the entry for *chiaro*:

```

:treccaniRDF a lexicog:LexicographicResource;
  dc:language "it" ;
  lexicog:entry :chiaro_entry .

:chiaro_entry a lexicog:Entry ;
  rdfs:member :chiaro_1_comp,
    :chiaro_2_comp,
    :chiaro_3_comp.

:chiaro1_comp a lexicog:LexicographicComponent;
:chiaro2_comp a lexicog:LexicographicComponent;
:chiaro3_comp a lexicog:LexicographicComponent.

```

Next we encode a lexicon which represents the content of the resource in the last listing.

```

:myItLexicon a lime:Lexicon;

```

```

lime:language "it" ;
lime:entry :chiaro1_adj, :chiaro2_n,
:chiaro3_adv .

```

```

:chiaro1_adj a ontolex:LexicalEntry .
:chiaro2_n a ontolex:LexicalEntry .
:chiaro3_adv a ontolex:LexicalEntry .

```

Finally we bring the two resources together using the *describes* property.

```

:chiaro1_comp lexicog:describes :chiaro1_adj .
:chiaro2_comp lexicog:describes :chiaro2_n.
:chiaro3_comp lexicog:describes :chiaro3_adv

```

4.1.2. OntoLex Morphology Module

Since November 2018, the W3C OntoLex community group has been developing another extension of the core model that would allow for better representation of morphological data in lexical resources.

Morphology plays an important role in many languages, and its description has also played an impor-

tant role in the work of lexicographers. The extent of its presence in concrete resources can vary, ranging from the sporadic indication of certain specific forms in a dictionary (e.g. plural form for some nouns) to electronic resources which provide tables with entire inflectional paradigms for every word.¹⁰²

The core OntoLex-Lemon model, together with LexInfo (see Section 3.5), provides the means of encoding basic morphological information: for lexical entries, morphosyntactic categories such as part of speech can be provided and basic inflection information (i.e. morphological relationship between a lexical entry and forms) can be modelled by creating any additional inflected forms with corresponding morphosyntactic features (e.g. case, number, etc.). However this only covers a small portion of the potential morphological data present in many lexical resources. Neither derivation (i.e. morphological relationships between lexical entries) nor additional inflectional information (e.g. declension type for Latin nouns) can be properly modelled with the core model. The new **OntoLex Morphology** module has been proposed to address these limitations. The scope of the module is threefold:

- *Representing derivation*: for a more sophisticated description of the decomposition of lexical entries;
- *Representing inflection*: introducing new elements to represent paradigms and wordform-building patterns;
- Providing means to *create wordforms automatically* based on lexical entries, their paradigms and inflection patterns.

Figure 2 presents the diagram for the module.

The central class of the module, and which is used in the representation of both derivation and inflection, is *Morph* with subclasses for different types of morphemes.

For derivation, elements from the *decomp* module are reused. A derived lexical entry has *Components* for each of the morphemes of which it consists. A *stem* corresponds to a different lexical entry whereas morphemes which do not correspond to any headwords, correspond to an object of a *Morph* class. A derived lexical entry has constituent properties pointing to objects of the *Component* class:

¹⁰²For example, *Wiktionary*, <https://en.wiktionary.org/wiki/Buch#Declension>.

```

:lex_drive_v a ontolex:LexicalEntry .
:lex_driver_n a ontolex:LexicalEntry ;
                decomp:constituent :component_drive,
                                :component_er .

:component_drive a decomp:Component ;
                decomp:correspondsTo :lex_drive_v .
:component_er a decomp:Component ;
                decomp:correspondsTo :suffix_er .

:suffix_er a morph:AffixMorph .

```

Inflection is modelled as follows: every instance of *Form* has properties *morph:consistsOf* which point to instances of *morph:Morph*.¹⁰³ These instances can have morphosyntactic properties expressed by linking to an external vocabulary, e.g. *LexInfo*:

```

:lex_drive_v a ontolex:LexicalEntry ;
                ontolex:otherForm :form_drives .

:form_drives a ontolex:Form ;
                consistsOf :stem_drive_v, :suff_s .

:suff_s a morph:AffixMorph ;
                lexinfo:number lexinfo:Plural .

```

The module¹⁰⁴ has not yet been published and is still very much under development by the W3C group. At the time of writing, a consensus was reached on the first two parts of the module, and their overview has been published in [48]. The third part, which concerns the automatic generation of forms is currently being discussed, and the next step will be validating the model by creating resources using the module.

4.1.3. *OntoLex-FrAC: Frequency, Attestations, Corpus Information*

In parallel with the development of the Morphology Module, the OntoLex W3C group has also started developing a separate module that would allow for the enrichment of lexical resources with information drawn from corpora and other language resources. Most notably, this includes the representation of attestations (used for instance as illustrative examples in a dictionary). These were originally discussed within lexicog (See 4.1.1), but this discussion quickly outgrew the confines of computational lexicography alone. Furthermore, it was observed that OntoLex lacked any support for corpus-based statistics, a cornerstone not only of empirical lexicography, but also of computational philology, corpus linguistics and lan-

¹⁰³One of the problems with this approach is that the order of the affixes is undefined, there are several possible solutions for this, e.g. a property *next* between two morphs, but currently there is no consensus in the community on how to model the order.

¹⁰⁴<https://www.w3.org/community/ontolex/wiki/Morphology>

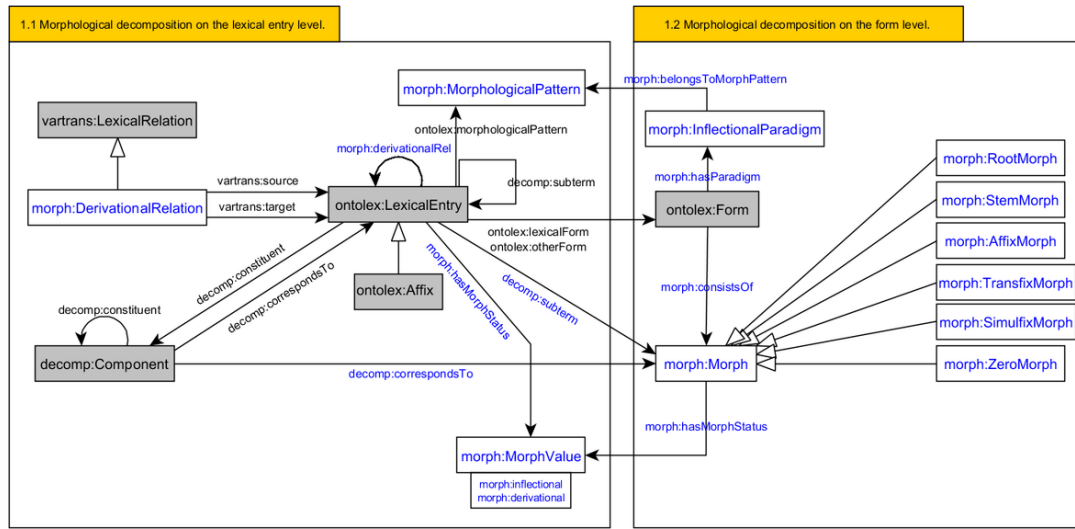


Fig. 2. Preliminary diagram for the Morphology Module.

guage technology, and thus, again, beyond the scope of the lexilog module. Finally, the community group felt the need to specifically address the requirements of modern language technology by extending its expressive power to corpus-based metrics and data structures like word embeddings, collocations, similarity scores and clusters, etc.

The development of the module has been use-case-based, which has dictated the order and development for various parts of the module.

So far, two papers have been published on parts of the module in the course of its development, and these parts can thus be considered to be relatively stable. They include the representation of (absolute) frequencies and attestations, and, by analogy, any use case that requires pointing from a lexical resource into an annotated corpus or other forms of external empirical evidence [49].

The central element which has been introduced in this module is `frac:Observable`. Since the type of elements for which corpus-based information can be provided is not limited to an entry, form, sense, or concept but can be any of these, `Observable` was introduced as a superclass for all these classes.

The module provides means to model only absolute frequency, because “relative frequencies can be derived if absolute frequencies and totals are known” [49, p. 2]. To represent frequency, a property `frequency` with an instance of `CorpusFrequency` as an object

should be defined. This instance must implement the properties `corpus` and `rdf:value`.¹⁰⁵

```

epsd:kalag_strong_v a ontolex:LexicalEntry;
  frac:frequency [
    a frac:CorpusFrequency;
    rdf:value "2398"^^xsd:int;
    frac:corpus
      <http://oracc.museum.upenn.edu/epsd2/pager>
  ] .

```

The usage recommendation is to define a subclass of `CorpusFrequency` for a specific corpus when representing frequency information for many elements in the same corpus.

For attestations, i.e. corpus evidence in FrAC, it is defined as “a special form of citation that provide evidence for the existence of a certain lexical phenomena; they can elucidate meaning or illustrate various linguistic features”. As with frequency, there is a class `Attestation`, an instance of which should be an object of a property attestation. It should have two properties: `attestationGloss` – text of the attestation and `locus` – location where the attestation can be found:

```

diamant:sense_1 a ontolex:LexicalSense;
  frac:attestation diamant:attestation_1 ;
  diamant:attestation_1 a frac:Attestation ;
  cito:hasCitedEntity diamant:cited_document_1 ;
  cito:hasCitingEntity diamant:sense_1;
  frac:locus diamant:locus_1 ;
  frac:quotation "... dat men licht yemant de cat
    aen het been kan werpen," .

```

¹⁰⁵Examples in this section are taken from [49].

The FrAC module does not provide an exhaustive vocabulary and instead promotes reuse of external vocabularies, such as CITO [50] for a citation object and NIF or WebAnnotation (see 4.2) to define a locus.

Another, more recent paper focused on representing embeddings in lexical resources is [51]. It should be noted that the term *embedding* is used here in a broader sense than is usual in the field of natural language processing, namely as a morphism $Y (f : X \rightarrow Y)$.¹⁰⁶ Therefore, the class *Embedding* has subclasses for modelling bags of words and time series.

The main motivation to model embeddings as a part of this module is to provide metadata as RDF for pre-computed embeddings, therefore a word vector itself is stored as a string with an embedding vector:

```
:embedding a
  frac:FixedSizeVector;
  dc:extent "300"^^xsd:int;
  rdf:value "0.145246 0.38873 ...";
```

As with modelling frequency, the recommendation is to define a subclass for the specific type of embedding concerned in order to make the RDF less verbose.

Figure 3 presents a diagram of the latest version of the module.

At the times of writing, module development is focused on collecting and modelling various use-cases. Among the many use-cases that were proposed during this phase, one stood out in particular and seemed to be more challenging than the others: this was related to the modelling of sign language data. Given the nature of the data (video clips with signs and/or time series of key coordinates for preprocessed data), it was decided that although the use-case was out of the scope of the FrAC module, it did indeed raise serious interest within the community, and therefore discussion on whether it will be developed as a separate module in the future, is now underway. The question of the scope of this new module and, more generally, its connection to OntoLex-Lemon, is currently subject to discussion.

4.1.4. Selected individual contributions

‘Unofficial’ OntoLex extensions pursued by individual research groups are manifold, and while these are not yet being pursued as candidates for future OntoLex modules in the W3C Community Group OntoLex, they may represent a nucleus and a cumulation point for future directions.

Selected recent extensions include *lemon-tree* [52], an OntoLex-Lemon and SKOS based model for pub-

lishing *topical thesauri*, where the latter are defined as lexical resources which are organised on the basis of meanings or topics.¹⁰⁷ The lemon-tree model has already been used to publish the Thesaurus of Old English [53] and reveals the flexibility of the OntoLex-Lemon/LLD approach in enabling the modelling of more specialised kind of linguistic information. As does *lemonEty* [54] another ‘unofficial’ extension of the OntoLex-Lemon model which has been proposed as a means of encoding etymological information. Namely, both the kinds of etymological information contained in lexica and dictionaries as well in other kinds of resources (such as articles or monographs). The *lemonEty* model does this by exploiting the graph-based structure of RDF data and by rendering explicit the status of etymologies as prospective hypotheses.

In both of these cases, the RDF data model along with the various different standards and technologies which make up the Semantic Web stack as a whole permit us to structure such information in salient and ‘meaningful’ ways that help to enhance the machine actionability of strongly heterogeneous linguistic data. This is something which the author of [55] attempts to demonstrate by way of a proposal to extend OntoLex-Lemon with temporal information in a ontologically well motivated way (while being careful to remain within the expressive limitations of RDF in order to exploit standard technologies for that framework including reasoning tools for OWL), and allowing the integration of lexical data with data relating to textual attestations and other historical information.

4.2. Annotation and Corpora

Summary In this section we give an overview of a number of LLD vocabularies for the annotation of texts. In Section 4.2.1 we give a detailed overview of this topic. Then we look at the two most popular such vocabularies *the NLP Interchange Format* (Section 4.2.2) and *Web Annotation* (Section 4.2.3). Next we look at two domain specific vocabularies, *Ligt* and *CoNLL-RDF* in Section 4.2.4. Finally in Section 4.2.5 we look at the prospects of a convergence between the vocabularies which we have discussed.

¹⁰⁶An injective structure-preserving map.

¹⁰⁷The lemon-tree specifications can be found here <https://ssstolk.github.io/onto/lemon-tree/>

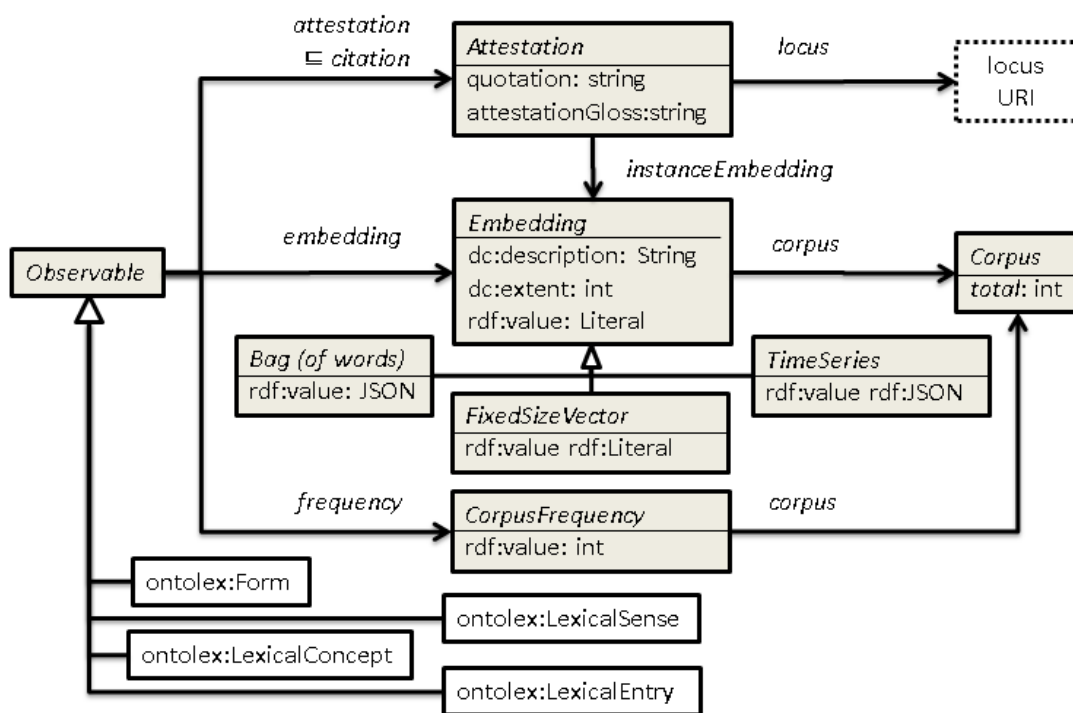


Fig. 3. Preliminary diagram for the FrAC Module.

4.2.1. Introduction and Overview

Linguistic annotation by NLP tools and within corpora has long been a topic of discussion within linked data and RDF circles, with different proposals grounded in traditions from natural language processing [56], web technologies [57], knowledge extraction [58], but also from linguistics [59], philology [60], and the development of corpus management systems [61, 62].

A practical introduction to the various different vocabularies used (by various different communities, for different purposes and according to different capabilities) for linguistic annotation in RDF today is given over the course of several chapters in [10]. In brief however the RDF vocabularies which are most widely used for this purpose are the NLP Interchange Format (NIF, in language technology) and Web Annotation (OA, in bioinformatics and digital humanities), as well as customizations of these. We describe NIF in Section 4.2.2 and Web Annotation in Section 4.2.3. In what follows we look at the relationship between RDF and two other pre-RDF vocabularies, namely LAF and TEI/XML. Next we look at some platform specific RDF vocabularies for annotations that have been de-

veloped over the years. Notable pre-RDF vocabularies include those developed by ISO TC37/SC4, in particular, the **Linguistic Annotation Framework (LAF)** that represents “universal” data structures shared by the various, domain- and application specific ISO standards [63]. Following the earlier insight that a labelled directed multigraph can represent any kind of linguistic annotation LAF produces concepts and definitions for four main aspects of linguistic annotation: **anchors and regions** elements in the primary data that annotations refer to, **markables (nodes)** elements that constitute and define the scope of the annotation by reference to anchors and regions; **values (labels)** elements that represent the content of a particular annotation. **relations (edges)** links (directed relations) that hold between two nodes and can be annotated in the same way as markables.

Note that in relation to Web Annotation *anchors* roughly correspond to Web Annotation selectors (or target URIs); *markables* roughly correspond to annotation elements; *values* to the body objects of Web Annotation. However in Web Annotation, relations as data

1 structures are not foreseen.¹⁰⁸ As for NIF, its relation
2 with LAF is more complex. Like Web Annotation,
3 NIF does not provide a counterpart of LAF relations,
4 but more importantly, the roles of regions and
5 markables is conflated in NIF: Every markable must
6 be a string (character span), and for every character
7 span, there exists exactly one potential markable (URI,
8 or, a number of URIs with different schemes that are
9 owl:sameAs).

10 At the moment, direct RDF serializations of the
11 LAF do not seem to be widely used in an LLOD context.
12 The reason is certainly that the dominant RDF
13 vocabularies for annotations, despite their deficiencies,
14 cover the large majority of use cases. One notable
15 RDF serialisation of LAF however is **POWLA** [64], an
16 OWL2/DL serialization of PAULA, a standoff-XML
17 format that implemented the LAF as originally described
18 by [65]. POWLA complements LAF core data structures
19 with formal axioms and a slightly more refined data
20 structures that support, for example, effective navigation
21 of tree annotations. On current applications of POWLA
22 see the CoNLL-RDF Tree Extension below.¹⁰⁹

23 It is also worth mentioning **TEI/XML** in the context
24 of this discussion. The standard, widely used in the
25 digital humanities and in computational philology,
26 only comes with partial support for RDF and does not
27 represent a publication format for Linked Data. Traditionally
28 there has been an acknowledgement on the part of the
29 TEI community of the value in being able to link from
30 a digital edition (or another TEI/XML document) to a
31 knowledge graph.¹¹⁰ Linking between (elements of)
32 electronic editions created with the TEI was addressed
33 by means of specialised XML attributes with narrowly
34 defined semantics. Accordingly, electronic editions
35 in TEI/XML do not normally qualify as

36
37
38
39
40
41 ¹⁰⁸Although Web Annotation lacks any formal counterpart of
42 edges or relations as defined by LAF there have been attempts to
43 define a vocabulary that extends Web Annotation with LAF data
44 categories [57], but this has apparently never been applied in practice.

45 ¹⁰⁹Others include [61] utilised an RDF graph, with a RDF vocabulary
46 for nodes, labels and edges to express linguistic data structures
47 over a corpus backend natively based on an RDBMS; a prototypical
48 extension of Web Annotation with an RDF interpretation of the
49 LAF described by [57], which and the LAPPS Interchange Format,
50 conceptually and historically an instance of LAF, which has see the
51 discussion below on platform-specific vocabularies.

¹¹⁰This is useful for instance for managing prosopographical, bibliographical or geographical information

1 Linked Data, even if they used and provide resolvable
2 URIs (TEI pointers).¹¹¹

3 The annotation *of* rather than *within* TEI documents,
4 however, has been pursued by Pelagios/Pleiades, a community
5 interested in the annotation of historical documents and maps
6 with geographical identifiers and other forms of geoinformation
7 (though this does not yet run to linguistic annotations). One
8 result of these efforts is the development of a specialised
9 editor called Recogito, and its extension to TEI/XML.
10 In this case the annotation is not part of the TEI document,
11 but stored as standoff annotation in a JSON-LD format, and
12 thus, is in compliance with established web standards and
13 re-usable by external tools and addressable as Linked Data.
14 However, this approach is restricted to cases in which the
15 underlying TEI document is static and no longer changes.
16 ¹¹²Therefore, there is a need for encoding RDF triples
17 directly inline in a TEI document. Happily, it has been
18 demonstrated that this can be done in a W3C- and XML-
19 compliant way by incorporating RDFa attributes into
20 TEI [66, 67]. As a result and after more than a decade
21 of discussions, the TEI started in May 2020 to work on
22 a customization that allowed the use of RDFa in TEI
23 documents.¹¹³

24 Over the years, several platforms, projects and tools
25 have come up with their own approaches for modelling
26 annotations and corpora as linked data. Notable examples
27 include the RDF output of machine reading and NLP systems
28 such as FRED [68], NewsReader [69] or the LAPPS Grid [70].
29 **FRED** provides output based on NIF or EARMARK [71],
30 with annotations par-

31
32
33
34
35 ¹¹¹This may not considered to be drastic for electronic editions
36 of historical manuscripts which one could conceivably complement
37 with information drawn from the LLOD cloud. The situation is,
38 however, quite different for dictionaries whose content could easily
39 be made accessible and integrated with other lexical resources on
40 the LLOD cloud, e.g., for future linking. The situation has, however,
41 begun to change over the last few years and long-standing efforts
42 to develop technological bridges between both TEI and LOD are
43 beginning to yield concrete results. For instance, different tools
44 for the conversion of lexical resources in different TEI dialects to
45 OntoLex-Lemon have been presented in the last years. Among others,
46 this includes a converter for TEI Dict/FreeDict dialect, <https://github.com/acoli-repo/acoli-dicts/tree/master/stable/freedict> [34].
47 For ELEXIS related developments TEI/RDF related developments, see
48 Section 5.2.3.

¹¹²Otherwise, the efforts for synchronization will by far outweigh
49 any benefit that the use of W3C standards for encoding the
50 annotation brings

¹¹³For the current status of the discussion, cf. <https://github.com/TEIC/TEI/issues/311> and <https://github.com/TEIC/TEI/issues/1860>

tially grounded in DOLCE [72], but enriched with lexicalised ad hoc properties for aspects of annotation covered by these.¹¹⁴ The **NewsReader Annotation Format (or NLP Annotation Format) NAF**, is an XML-standoff format for which a NIF-inspired RDF export has been described [73], and LIF, the LAPPS Interchange Format [74], a JSON-LD format used for NLP workflows by the LAPPS Grid Galaxy Workflow Engine [75]¹¹⁵.

Both LIF and NAF-RDF are, however, not generic formats for linguistic annotations but rather, provide (relatively rich) inventories of vocabulary items for specific NLP tasks.¹¹⁶ Neither seem to have been used as a format for data publication, and we are not aware of their use independently from the software they have originally been created for or are being created by.

Aside from software- or platform-specific formats, a number of vocabularies has been developed that address specific problems or user communities. In Section 4.2.4 we describe two such examples, the *Ligt* vocabulary that addresses the gap that community standards such as NIF and Web Annotation have with respect to morphology and the annotation of morphologically rich languages, and CoNLL-RDF, a vocabulary and an associated library that aims to mirror the structure and contents of popular NLP formats in RDF and to provide a round-tripping between these formats and RDF graphs.

Finally, the prospects for convergency between the solutions discussed in the whole of Section 4.2 are described in Section 4.2.5. (Note that in this section, we only discuss vocabularies that define *data structures* for linguistic annotation by NLP tools and in annotated corpora. Linguistic categories and grammatical features, as well as other information that represents the content of an annotation are assumed to be provided by a(ny) repository of linguistic data categories (see above)).

¹¹⁴For the rendering of discourse relations, for example, it produces properties such as `fred:becauseOf` (apparently extrapolated from the surface string, so, not ontologically defined).

¹¹⁵A more recent development in this regard is that efforts have been undertaken to establish a clear relation between LIF and pre-RDF formats currently used by CLARIN [76].

¹¹⁶Historically, LIF is grounded in LAF concepts and has been developed by the same group of people, however, no attempt seems to have been made to maintain the level of genericity of the LAF. Instead, application-specific aspects seem to have driven LIF design.

4.2.2. NLP Interchange Format

The NLP Interchange Format (NIF),¹¹⁷ developed at AKSW Leipzig, was designed to facilitate the integration of NLP tools in knowledge extraction pipelines, as part of the building of a Semantic Web tool chain and a technology stack for language technology on the web [58]. NIF provides support for a broad range of frequently occurring NLP tasks such as part of speech tagging, lemmatization, entity linking, coreference resolution, sentiment analysis, and, to a limited extent, syntactic and semantic parsing. In addition to providing a technology for integrating NLP tools in semantic web annotations, NIF also provides specifications for web services.

A core feature of NIF is that it is grounded in a formal model of strings, and the obligatory use of String URIs as fragment identifiers for anything annotatable by NIF. Every element that can be annotated in NIF has to be a string.¹¹⁸ NIF does support different fragment identifier schemes, e.g., the offset-based scheme defined by RFC 5147. [77] As a consequence, any two annotations that cover the same string are bound to the same (or `owl:sameAs`) URI. While this has the advantage of being able to implicitly merge the output of different annotation tools, this limits the applicability of NIF to linguistically annotated corpora. As an example, NIF does not allow us to distinguish multiple syntactic phrases that cover the same token. Consider the sentence “Stay, they said.”¹¹⁹ The Stanford PCFG parser¹²⁰ analyzes *Stay* as a verb phrase contained in (and only constituent of) a sentence. In NIF, both would be conflated. Likewise, zero elements in syntactic and semantic annotation cannot be expressed. Another limitation of NIF is its insufficient support for annotating the internal structure of words. It is thus largely inapplicable to the annotation of morphologically rich languages. Overall, NIF fulfills its goals to provide RDF wrappers for off-the-shelf NLP tools, but it is not sufficient for richer annotations as are frequently found in linguistically annotated corpora. Nev-

¹¹⁷<https://persistence.uni-leipzig.org/nlp2rdf/ontologies/nif-core/nif-core.html>

¹¹⁸In particular, this includes the classes `nif:Phrase` and `nif:Word`. With the introduction of support for provenance annotations, NIF 2.0 also introduced `nif:Annotation` which can be attached as a property to a NIF string. However, it is to be noted that the linguistic data structures defined by NIF 2.0 are *not* subclasses of `nif:Annotation`, but of `nif:String`.

¹¹⁹From Stephen Dunn (2009), ‘Dont Do That’, poem published in the New Yorker, June 8, 2009.

¹²⁰<http://nlp.stanford.edu:8080/parser/index.jsp>

ertheless, NIF has been used as a publication format for corpora with entity annotations.¹²¹

NIF continues to be a popular component of the DBpedia technology stack. At the same time, active development of NIF seems to have slowed down since the mid-2010s, whereas limited progress on NIF standardization has been achieved. A notable exception in this regard is the development of the Internationalization Tag Set [78, ITS] that aims to facilitate the integration of automated processing of human language into core Web technologies. A major contribution of ITS 2.0 has been to add an RDF serialization into NIF as part of the standard.

More recent developments of NIF include extensions for provenance (NIF 2.1, 2016) and the development of novel NIF-based infrastructures around DBpedia and Wikidata [79]. In parallel to this, NIF has been the basis for the development of more specialised vocabularies, e.g., CoNLL-RDF for linguistic annotations originally provided in tabular formats, see Section 4.2.4.

4.2.3. Web Annotation

The Web Annotation Data Model is an RDF-based approach to standoff annotations (in which annotations and the material to be annotated are stored separately) proposed by the Open Annotation community.¹²² It is a flexible means of representing standoff annotation for any kind of document on the web. Although the most common use case of Web Annotation is the attaching of a piece of text to a single web resource, it is intended to be applicable across different media formats. So far Web Annotation has been primarily applied to linguistic annotations in the biomedical domain, although other notable applications include NLP [57] or Digital Humanities [82]. Web Annotation recommends the use of JSON-LD to add a layer of standoff annotations to documents and other resources accessible over the web, with primary data structures defined by the Web Annotation Data Model, formalised as an OWL ontology.

The core data structure of the Web Annotation Data Model is the annotation, i.e., instances of `oa:Annotation` that have an `oa:hasTarget` property that identifies the element that carries the annotation, and the `oa:has-`

¹²¹The most prominent example, the NIF edition of the Brown corpus published in 2015, formerly available from <http://brown.nlp2rdf.org/>, does not seem to be accessible anymore. Attempted to access on Jan 23, 2021.

¹²²The Web Annotation data model and vocabulary were published as W3C recommendations in 2017 [80, 81].

Source property that – optionally – provides a value for the annotation, e.g., as a literal. The target can be a URI (IRI) or a selector, i.e., a resource that identifies the annotated element in terms of its contextual properties, formalised in RDF, e.g., its offset or characteristics of the target format. By supporting user-defined selectors and a broad pool of pre-defined selectors for several media types, Web Annotation is applicable to any kind of media on the web. Targets can also be more compact String URIs, as introduced, for example, by NIF. NIF data structures can thus be used to complement Web Annotation [58].

Web Annotation can be used for any labelling or linking task, e.g., POS tagging, lemmatization, entity linking. It does, however, not support relational annotations such as syntax and semantics, nor (like NIF) the annotation of empty elements. The addition of such elements from LAF has been suggested [57], but does not seem to have been adopted, as labelling tasks dominate the current usage scenarios of Web Annotation.

Unlike NIF, Web Annotation is ideally suited for the annotation of multimedia content or entities that are manifested in different media simultaneously (e.g., in audio and transcript). As a result, it has become popular in the digital humanities, e.g., for the annotation of geographical entities with tools such as *Recogito* [83], especially since support for creating standoff annotations for static TEI/XML documents was added (around March 2018 [84, p.247]).

4.2.4. Domain-specific solutions: *Ligt* and *CoNLL-RDF*

Interlinear glossed text (IGT) is a notation where additional annotations are placed, as the name suggests, between the lines of a text to be annotated. With the purpose of helping readers to understand and interpret linguistic phenomena, the notation is frequently used in education and various language sciences such as language documentation, linguistic typology, and philological studies (for instance, it is commonly used to gloss linguistic examples). Moreover, IGT data can consist of different layers, including translation and transliteration layers and usually contains layers for ensuring morpheme-level alignment. This is not supported by any established vocabularies for representing annotations on linguistic corpora. And although there exist several specialised formats which are specifically designed for the storage and exchange of IGT, these formats are not re-used across different tools, limiting the reusability of annotated data. In order to help overcome this situation and improve data interoper-

ability, the RDF vocabulary **Ligt** [85] has been proposed for representing IGT as linked data. Ligt is a tool-agnostic representation model for IGT which in addition to structural interoperability also enables the use of LLD vocabularies and terminology repositories.

The Ligt vocabulary was developed as a generalisation over the data structures employed by established tools for creating IGT annotations, most notably Toolbox [86], FLEx [87] and Xigt [88].¹²³ Ligt is intended to facilitate a pivot format that faithfully captures linguistic information produced by these tools in a uniform way for subsequent processing. Notably, since its publication Ligt has been adopted by third party users to model and annotate IGT from 280 endangered languages and their publication as Linked Open Data [89].

Although Ligt was designed for very specific set of domain requirements it can be considered a useful contribution to LLD vocabularies for textual annotation. This is because it provides data structures that are relevant for low-resource and morphologically rich languages but which had been neglected by earlier RDF vocabularies for linguistic annotation on the web, in particular, by NIF and Web Annotation.¹²⁴

Another domain specific RDF-based vocabulary which aims to provide a serialisation-independent way of dealing with textual annotations is **CoNLL-RDF** [90]. This latter vocabulary is based on the so-called “CoNLL formats”, a family of a tab-separated values (TSV) based-formalisms used to represent linguistically annotated natural language in fields such as NLP,¹²⁵ corpus linguistics, and more generally in the language sciences. CoNLL-RDF [90] provides a data model and a programming library that aim to facilitate the processing and transformation of such data regardless of the original order and number of columns, whether the source format used fixed-size tables (as for most CoNLL dialects) or variable size tables (such as all CoNLL formats that contain semantic role annotations). Sentence after sentence is converted to an RDF graph in accordance to the label information provided by the user. The listing below provides a slightly simplified annotation from the 2005 edition of the Shared

¹²³One should note that these tools are currently incompatible with each other and information can only exchanged between them if manual corrections are applied.

¹²⁴However it would be possible to encode Ligt information with a generic LAF-based vocabulary such as POWLA

¹²⁵Indeed in NLP the CoNLL formats have become de-facto standards for the most frequently used types of annotations having been popularised in a long-standing series of shared tasks over the last two decades

Task of the SIGNLL Conference on Computational Natural Language Learning (CoNLL-05):

```
# WORD          POS  PARSE
The             DT   (S (NP  *
spacecraft     NN
...
```

Here, the wordform is provided in the first column, the second column provides part-of-speech tag. The PARSE column contains a full parse in accordance with the Penn Treebank [91]. The CoNLL-RDF library reads such data as a continuous stream, every sequence of rows enclosed in empty lines will be processed as a block, assigned a URI and the type `nif: Sentence`, every row will be assigned a URI and the type `nif: Word`, and the annotation of every column will be stored as value of a property in the `conll: namespace` that is generated from the column label.¹²⁶ Links between and among sentences and words are encoded in accordance with NIF:

```
:s1_1 a nif:Word; nif:nextWord :s1_2;
conll:WORD "The"; conll:POS "DT";
conll:PARSE "(S (NP *".
```

Amongst other things, a CoNLL-RDF edition of the Universal Dependencies corpora¹²⁷ is available in the LLOD cloud diagram. The corpora are linked with the OLiA ontologies; further linking with additional LLOD resources, in particular, lexical resources, has not been explored so far. CoNLL-RDF has also been applied to the linking of corpora to dictionaries [92] and knowledge graphs [93]. It has also formed the basis of work on the syntactic parsing of historical languages [94, 95], the consolidation of syntactic and semantic annotations [96], corpus querying [97], and language contact studies [98]. In addition to the storing of syntactic parses as plain strings a further extension of CoNLL-RDF adds native support for tree structure [99], extending NIF/CoNLL-RDF data structures with POWLA [100]. As a result, the phrase structure of the example above can now be decoded:

```
:s1_1 a nif:Word; nif:nextWord :s1_2;
conll:WORD "The"; conll:POS "DT";
powla:hasParent _:np.
_:np a conll:PARSE; rdf:value "NP";
powla:next _:vp;
```

¹²⁶The columns HEAD (for dependency annotation) and PRED-ARGS (for semantic role annotations) are treated differently as they produce object properties, i.e., links, rather than datatype properties. Similarly, the column ID receives special handling. If provided as column label, as its value is used to overwrite the offsets that CoNLL-RDF normally adopts for creating word (row) URIs.

¹²⁷<https://universaldependencies.org/>


```

1   powla:hasParent _:s.
2   _:s a conll:PARSE; rdf:value "S".
3   ...

```

The CoNLL-RDF tree extension uses a minimal fragment of POWLA, the properties `powla:hasParent` (pointing to the parent node in a DAG) and `powla:next` (pointing to the following sibling in a tree). The class `powla:Node`, implicit in the listing above, can be RDFS-inferred from the use of these properties.

4.2.5. Towards a Convergence

The large number of vocabularies mentioned above already reveals something of a problem, that is, that applications and data providers may choose from a broad range of options, and depending on the expectations and requirements of their users, they may even need to support multiple, different output formats, protocols and service specifications that could potentially be mutually incompatible. So far, no clear consensus has emerged, albeit NIF and Web Annotation appear to enjoy relatively high popularity in their respective user communities. However, they are not compatible with each other and nor do they support linguistic annotation to the same or even a sufficient extent, thus motivating the continuous development of novel, more specialised vocabularies. Synergies between Web Annotation and NIF were explored relatively early on [58], and Cimiano et al. [101, p.89-122] describe how they can be used in combination with each other, more specialised vocabularies such as CoNLL-RDF, and more general vocabularies such as POWLA to model data in a way that suits the following criteria:

- it is applicable to any kind of primary data, including non-textual data (via Web Annotation selectors);
- it can also express reference to primary data in a compact fashion (via NIF String URIs);
- permits round-tripping between RDF graphs and conventional formats (via CoNLL-RDF and the CoNLL-RDF library);
- it supports generic linguistic data structures (via POWLA, resp., the underlying LAF model).

However, while the combination of these various components is possible and in principle operational, this also means that a user or provider of data needs to understand and develop a coherent vision of at least five different data models: Web Annotation, NIF, CoNLL-RDF, POWLA and the original or conventional structure of the data. Moreover, the data structures of these formats are parallel, in parts, and then, a principled and

consistent choice between, say, a `oa:Annotation` (from Web Annotation), a `powla:Node` (from POWLA), a `nif:String` and a `nif:Annotation`, has to be made.

Generally speaking, this situation is intractable, and thus, **the W3C Community Group Linked Data for Language Technology (LD4LT)** is currently engaged in a process to develop a harmonization of these vocabularies. While this has been under development since about mid-2018, regular discussions via LD4LT began in early 2020 only. Concrete results so far include a survey over requirements for any vocabulary for linguistic annotation on the web and on the degree to which NIF, Web Annotation and other vocabularies support these at the moment.¹²⁸ So far, 51 requirements have been identified, clustered in 6 groups:

1. LLOD compliancy (adherence to web standards, compatibility with community standards for linguistic annotation)
2. expressiveness (necessary data structures to represent and navigate linguistic annotations)
3. units of annotation (addressing primary data and annotations attached to it)
4. sequential data structures (preserving and navigating sequential order)
5. relations (annotated links between different units of annotation)
6. support for/requirements from specific applications and use cases (e.g., intertextual relations, linking with lexical resources, alignment, dialog annotation).

So far, this is still work in progress, but if indeed, these challenges can be resolved at some point in the future, and a coherent vocabulary for linguistic annotations emerges, we expect a similar rise in popularity for the adoption of the Linked Data paradigm for encoding linguistic annotations as we have seen in the last years for lexical resources. Also, this was largely driven by the existence of a coherent and generic vocabulary, and indeed, the drift in application that the OntoLex model has recently faced very much reflects the need for consistent, generic data models.

A question at this point may be what the general benefit of modelling annotations as linked data may be in comparison to other conventional solutions, and different user communities may have different answers

¹²⁸The survey can be accessed via <https://github.com/ld4lt/linguistic-annotation/blob/master/survey/required-features.md>, also compare the tabular view under <https://github.com/ld4lt/linguistic-annotation/blob/master/survey/required-features-tab.md>.

to that. It does seem, though, that one potential killer application can be seen in the capability to integrate, use and re-use pieces of information from different sources. In linguistic annotation, a still largely unsolved problem is how to efficiently process standoff annotation, and indeed, the application of RDF and/or Linked Data has long been suggested as a possible solution [56, 59, 61, 64], but only recently, have systems that support RDF as an output format emerged [62]. While it is clear that standoff is a solution it is also true that, the different communities involved have not agreed on commonly used standards to encode and exchange their respective data. In DH and BioNLP, Web Annotation and JSON-LD seems to dominate; in knowledge extraction and language technology, NIF (serialised in JSON-LD or Turtle) seem to be more popular; for digital humanities, the TEI is currently revising XML standoff specifications,¹²⁹ and support for RDF serializations (RDFa) or standoff (Web Annotation in JSON-LD) also seems to be growing, as mentioned above.

4.3. Metadata

Summary In the first section, Section 4.3.1, we give an introduction and overview of metadata trends in LLD and related areas. Next we give a detailed description of two important metadata resources for LLD. These are META-SHARE, described in Section 4.3.2, and the OntoLex-Lemon *lime* module, described in Section 4.3.3. The latter section also features a discussion of future metadata challenges for LLD language resources. Finally in Section 4.3.4 we address the important challenge of language identification which is an essential part of the metadata of a language resource.

4.3.1. Introduction

The rise of data-driven approaches that use Machine Learning, and in particular recent breakthroughs in the field of Deep Learning, have secured a central place for data in all scientific and technological areas. Cross-disciplinary research has also boosted the sharing of data arising across different communities. Thus, a huge volume of data has become available through various repositories, but also via aggregating catalogues, such as the European Open Science Cloud¹³⁰ and the

Google dataset search service¹³¹. Metadata play an instrumental role in the discovery, interoperability and hence (re-)use of digital objects, and indeed that they act as the intermediary between consumers (humans and machines) and digital objects. For this reason, the FAIR principles [1] include specific recommendations for metadata (see also section 1). Of particular relevance to this section is principle R1.3 which recommends that "(Meta)data meet domain-relevant community standards". According to this principle, the adoption of community standards or best practices for data archiving and sharing, including "documentation (metadata) following a common template and using common vocabulary" facilitates the re-use of data. We thus take a closer look at metadata models commonly used for language resources in the linguistics, digital humanities and language technology communities.

Although the focus of this section is on community models, we cannot leave the most popular general purpose models for dataset description out of this overview. Language is an essential part of human cognition and expression and thus present in all types of data; research on language and language-mediated research is carried out on data from all domains and human activities. All of this obviously extends the search space for data to catalogues other than the purely linguistic ones. The three models that currently dominate the description of datasets are DCAT¹³², schema.org¹³³ and DataCite¹³⁴. DCAT profiles are used in various open data catalogues, such as the EU Open Data portal¹³⁵, while schema.org is used for the Google dataset search engine; finally, DataCite, a leading provider of persistent identifiers (namely DOIs), has developed a schema with a small set of core properties which have been selected for the accurate and consistent identification of a resource for citation and retrieval purposes. There are various initiatives for the collection of crosswalks of community-specific metadata models with these models¹³⁶, as well as recommendations for extensions for specific data types (e.g., CodeMeta¹³⁷ and Bioschemas¹³⁸ for source code software and life science resources respectively). Of course, these mod-

¹²⁹See <https://github.com/TEIC/TEI/issues/1745> for pointers.

¹³⁰<https://www.eosc-portal.eu>

¹³¹<https://toolbox.google.com/datasetsearch>

¹³²<https://www.w3.org/TR/vocab-dcat-2/>

¹³³<https://schema.org/>

¹³⁴<https://schema.datacite.org/>

¹³⁵<https://data.europa.eu/euodp/en/data/>

¹³⁶See for instance, <https://rd-alliance.github.io/Research-Metadata-Schemas-WG/>

¹³⁷<https://codemeta.github.io/>

¹³⁸<https://bioschemas.org/>

els are not intended to capture all the specificities required for the description of linguistic features and, thus, we do not go into further details for them in this paper.

Among models for the description of language resources in general (and not just LLD resources), the Component Metadata Infrastructure (CMDI) profiles [102, 103], and the TEI guidelines (introduced above) stand out. CMDI is a framework designed to describe and re-use metadata; "profiles" can be constructed on the basis of building blocks ("components") that group semantically related metadata elements (e.g., address, identity, etc.) and can be used as ready-made templates catering for specific use cases (e.g., for lexica, for linguistic corpora, for audio corpora, etc.). CMDI profiles are used by various humanities and social sciences communities within the CLARIN¹³⁹ research infrastructure. The TEI standard specifies an encoding scheme for the representation of texts in digital form, chiefly in the humanities, social sciences and linguistics; it includes specific elements for the description of texts at the collection and individual text levels. Both models, however, are XSD-based¹⁴⁰, and therefore not discussed further in this section.

We should also mention the CLARIN Concept Registry (CCR)¹⁴¹, which is a collection of linguistic concepts [105, 106]. It is the successor to the ISOcat data category registry (described in Section 3.5) and is currently maintained by CLARIN. The CCR is implemented in SKOS and includes a concept scheme for metadata, but is a structured list without ontological relations, either internally or externally to other vocabularies. It mainly serves as the semantic interoperability layer of CLARIN; this is achieved through linking metadata fields included in CMDI profiles to concepts from the CCR.

4.3.2. Language Resource Metadata: The META-SHARE ontology

The META-SHARE¹⁴² (or MS-OWL in short) model [107], implemented as an OWL ontology, has been designed specifically for language resources, including data resources (structured or unstructured datasets, lexica, language models, etc.) and technologies used for lan-

guage processing [108]. The first version of MS-OWL was (semi-)automatically created from the META-SHARE XSD schema [108, 109] (which was designed to support the META-SHARE infrastructure [110]) and discussed in the framework of the LD4LT group. The second version, which is described here, has evolved from it, taking into account advancements in the Language Technology domain and related metadata requirements (such as the necessity for the description of workflows, interoperability issues between language processing tools and processing resources, etc.) as well as current trends in the overall metadata landscape [107].

MS-OWL has been constructed according to three key concepts: *resource type*, *media type* and *distribution*. These give rise to the following basic classes:

- LanguageResource, with four subclasses derived from the notion of resource type:
 - * Corpus: for structured collections of pieces of language data, typically of considerable size and which have been selected according to criteria external to the data (e.g., size, language, domain, etc.) with the aim of representing as comprehensively as possible a specific object of study;
 - * LexicalConceptualResource: covering resources such as term glossaries, word lists, semantic lexica, ontologies, etc., organised on the basis of lexical or conceptual units (lexical items, terms, concepts, phrases, etc.) along with supplementary information (e.g., grammatical, semantic, statistical information, etc.);
 - * LanguageDescription: for resources which are intended to model a language or some aspect(s) of a language via a systematic documentation of linguistic structures; members of this class are typically statistical and machine learning-computed language models and computational grammars;
 - * ToolService: for any type of software that performs language processing and/or related operations (e.g., annotation, machine translation, speech recognition, speech-to-text synthesis, visualization of annotated datasets, training of corpora, etc.);
- MediaPart: this is a parent class for a number of other subclasses, combining together the notions of resource and media type; it is not meant to be used directly in the description of language

¹³⁹<https://www.clarin.eu>

¹⁴⁰The conversion of CMDI metadata records offered in CLARIN into RDF [104] should not be confused with the construction of an RDF model for CMDI profiles

¹⁴¹<https://concepts.clarin.eu/ccr/browser/>

¹⁴²<http://w3id.org/meta-share/meta-share>

resources. The media type refers to the form/physical medium of a data resource (i.e., member of one of the first three subclasses above) and it can take the values text, audio, image, or video. To cater for multimedia/multimodal language resources (e.g. a corpus of videos and their subtitles, or corpus of audio recordings and their transcripts), language resources are represented as *consisting* of at least one media part: the *mediaPart* property is used to link an instance of the class *Corpus* to instances of *CorpusTextPart*, *CorpusAudioPart*, and so on; similarly, *LexicalConceptualResource* is linked to *LCRTextPart*, *LCRVideoPart*, etc.

- *DatasetDistribution* and *SoftwareDistribution*: these are conceived as subclasses of *dcate:Distribution*, which represents the accessible form(s) of a resource. For instance, software resources may be distributed as web services, executable files or source code files, while data resources as PDF, CSV or plain text files or through a user interface.

MS-OWL caters for the description of the full lifecycle of language resources, from conception and creation to integration in applications and usage in projects as well as recording relations with other resources (e.g., raw and annotated versions of corpora, tools used for their processing, models integrated in tools, etc.) and related/satellite entities¹⁴³.

The properties recommended for the description of language resources are assigned to the most relevant class. Thus, the *LanguageResource* class groups properties common to all resource/media types, such as those used for identification purposes (title, description, etc.), recording provenance (creation, publication dates, creators, providers, etc.), contact points, etc. More technical features and classification elements, that depend on resource/media types, as well as instances of *MediaPart* and *Distribution* are attached to the respective *LanguageResource* subclasses. Thus, properties for *LexicalConceptualResource* encode the subtype (e.g. computational lexicon, ontology, dictionary, etc.), and the contents of the resource (unit of description, types of accompanying linguistic and extralinguistic information, etc.); properties for *Corpus* include corpus subclass (raw, annotated corpus, anno-

¹⁴³The current work discusses only the core part of MS-OWL targeting the description of language resources and leaves aside the representation of satellite entities (persons, organizations, projects, etc.)

tations), and information on corpus contents. It should be noted that the language of the resource's contents, a piece of metadata of particular relevance to all language resources, is encoded in the media part subclasses rather than the top *LanguageResource* class; this is in line with the principles adopted for the representation of multimedia/multimodal resources consisting of parts with different languages (e.g. a corpus of video recordings in one language, its subtitles in the same language and their translations in another language). Finally, the two distribution classes (*DatasetDistribution* and *SoftwareDistribution*) provide information on how to access the resource (i.e., how and where it can be accessed), technical features of the physical files (such as size, format, character encoding) and licensing terms and conditions. A dedicated module has been devised for the structured representation of licenses commonly used for language resources, reusing existing vocabularies and extending the Open Digital Rights Language¹⁴⁴ core model [111].

To better illustrate the structure of the MS-OWL, figure 4 depicts a subset of the mandatory and recommended properties for the description of a corpus.

One of the additions made between the two versions of the MS ontology is the development of another vocabulary, again implemented as an OWL ontology, OMTD-SHARE¹⁴⁵ [112]. OMTD-SHARE can be considered as complementary to MS-OWL. It covers *functions* (tasks performed by software components), *annotation types* (types of information extracted or annotated by such software), *methods* (classification of the theoretical method used in the algorithm), and *data formats* of the resources that can be processed by such software. The ontology was begun within the framework of the OpenMinTeD project¹⁴⁶, which focused on Text and Data Mining resources, and has been enriched afterwards. The class *Operation* has been extended to cover Language Technology (LT) operations at large (now also referred to as "LT taxonomy"). Specific properties of MS-OWL make reference to the OMTD-SHARE classes. *Operation* is used for describing the function of tools/services as well as for applications for which a data resource can be used or has already been used. The *annotationType* for annotated corpora takes values from the *AnnotationType* class; linguistic annotation types are linked to the OLIA on-

¹⁴⁴<https://www.w3.org/ns/odrl/2/>

¹⁴⁵<http://w3id.org/meta-share/omtd-share/>

¹⁴⁶<https://www.openminded.eu>

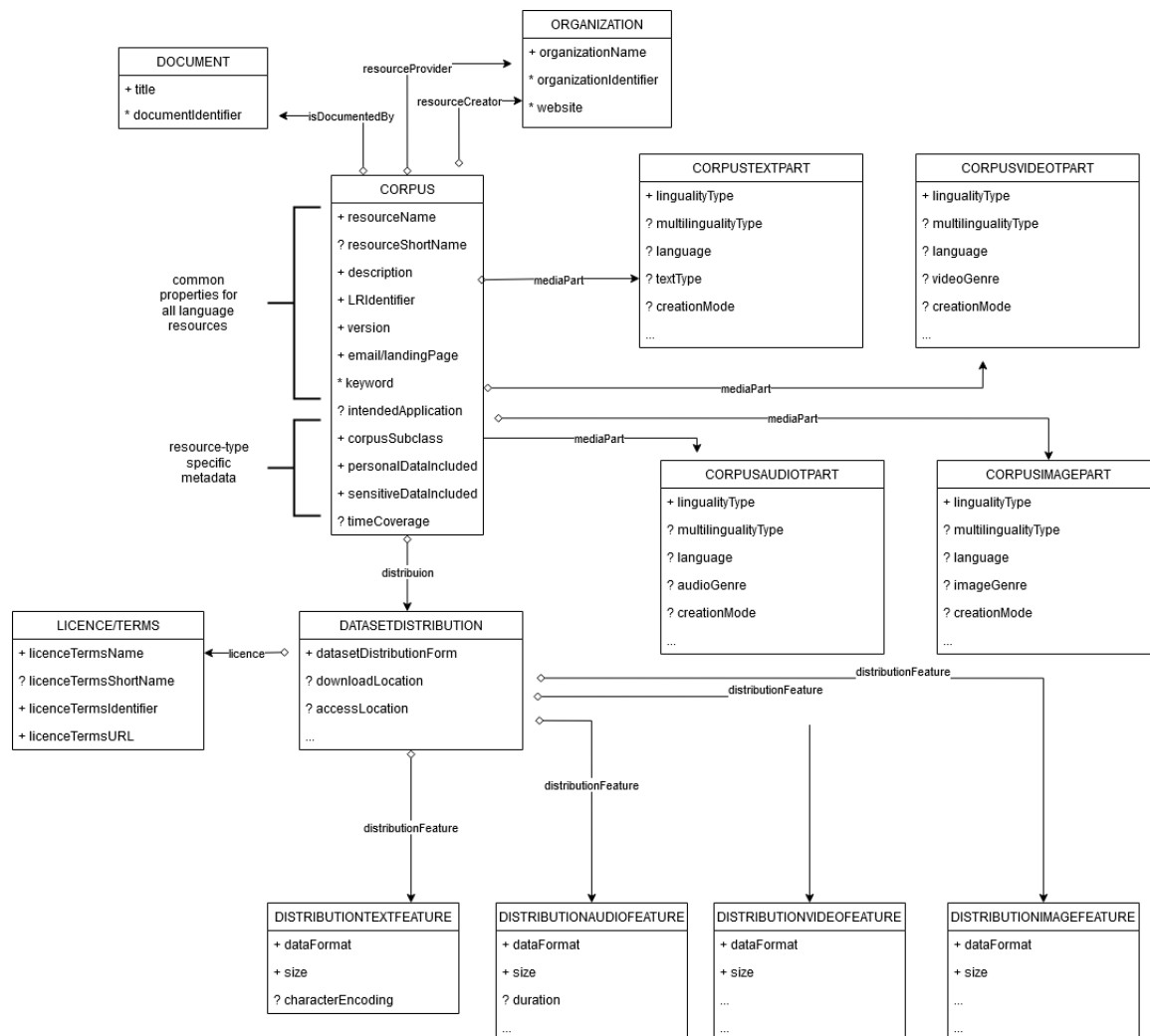


Fig. 4. Simplified subset of the MS-OWL for corpora.

tology (work in progress), while domain-specific annotation types for neighbouring domains are also foreseen (e.g., for elements in the document structure of publications, biomedical entities, etc.).

Both the MS-OWL and OMTD-SHARE ontologies have been published and are currently undergoing evaluation and improvements. They are deployed in the description of language resources in catalogues of language resources. More specifically, the first version of MS-OWL is used in LingHub¹⁴⁷, a data portal aggregating metadata records for language resources hosted in various repositories and

catalogues [113, 114], while the second version, the one described here, is used in the European Language Grid¹⁴⁸, which is a platform for language resources with a focus on industry-relevant Language Technology in Europe [115]. Among the immediate plans, crosswalks with DCAT and schema.org are a priority, to ensure wider uptake and interoperability with (meta)data from other communities.

¹⁴⁷<http://linghub.org/>

¹⁴⁸<https://live.european-language-grid.eu/>

4.3.3. Linguistic Metadata for Lexical Resources: *lime*

Another metadata model that has relevance here is OntoLex-Lemon's¹⁴⁹ own dedicated metadata module that, in keeping with the overall citric theme, is called *lime* or the *LInguistic METadata* module [116].

The module is illustrated in Figure 5. Before we go onto describe the module in more detail, it is important to point out that *lime* focuses on providing metadata descriptions at the "level of lexicon-ontology interface"¹⁵⁰. That is, it concentrates on how ontological concepts in a so-called *reference dataset* (e.g. defined as an ontology that describes "the semantics of the domain") are *lexicalised* or given a linguistic grounding in a lexicon (viewed as a collection of lexical entries). The OntoLex-Lemon guidelines further refer to a *concept set* which is a set of individuals of class *Lexical Concept* (described as potentially "bearing a conceptual backbone to a lexicon"). The aim of the *lime* module then is to provide quantitative and qualitative (metadata) information about the relations between these (the links between them). In other words many, though not all as we will see below, of its classes and properties will not apply in cases where OntoLex-Lemon is *only* used to encode a lexicon, and where entries and their senses aren't linked to either *Lexical Concept* individuals or to ontology entities (such as is true of an increasing number of lexicon-centric use cases).

More generally useful classes and properties, however, include the *lime:Lexicon* class. This is defined as a subclass of *void:Dataset*¹⁵¹, and represents a set of individuals of the class *Lexical Entry* which are linked to *lime:Lexicon* via the property *lime:entry*. The whole lexicon, as well as individual entries, can be assigned to a certain language, as specified by the datatype property *lime:language* (the guidelines also recommend the use of the Dublin Core property and the use of either *lexVo* or Library of Congress language tags, see Section 4.3.4 for an extended discussion of language tags). The property *lime:linguisticCatalog* specifies the linguistic model, i.e. the catalogue of linguistic categories used for the annotation of the lexical entries. This could be, for instance *LexInfo*.

In order to show the use of these more general *lime* classes in practise we will look at a simple example

taken from the W3C guidelines¹⁵². This can be seen in diagrammatic form in Figure 6.

Figure 6 corresponds to the following listing (again based on that give in the W3C guidelines¹⁵³).

```
:lexicon a lime:Lexicon;
lime:language "en";
lime:entry :lex_high;
lime:entry :lex_cat;
lime:entry :lex_marry;
lime:entry :lex_intangible_assets.
]
```

Here we can see how the *lime* properties and class introduced above can help us to describe some of the most fundamental lexicon-specific metadata categories of a lexical resource. We can of course use Dublin Core properties such as *description* and *creator* to further flesh out this metadata description. The *lime:LexicalizationSet* class (again a subclass of *void:Dataset*) represents a collection of *lexicalizations*, defined as pairs consisting of a lexical entry and an associated entry in the reference dataset (this could be an OWL ontology but as the guidelines specify might be any "RDF dataset which contains references to objects of a domain of discourse"). The metadata properties on the *lime:LexicalizationSet* allow us to describe, among other things¹⁵⁴ how many entities have been lexicalised (by at least one entry), how many pairs of entries and ontology elements there are, as well as how many ontology elements have been lexicalised on average. *lime* also defines the class *Lexical Linkset* (*lime:LexicalLinkSet*, subclass of *void:Dataset*), individuals of which are links between a set of lexical concepts (i.e., members of the class *ontolex:LexicalConcept*) and the reference dataset (ontology). For this class, *lime* defines properties describing, for example, the number of links between the two resources in question. Last, the *Conceptualization Set* (*lime:ConceptualizationSet*) is analogous to the *lime:LexicalizationSet* but caters for the links between the the lexicon and the concept set.

Future Metadata Challenges for Lexical Resources

The *lime* classes which we have just looked at, along with several others, enable users to give a detailed metadata description of the lexicalisation relationship of a lexicon with a semantic resource. The META-SHARE ontology offers a wide number of classes and properties for encoding common metadata properties

¹⁴⁹This section assumes some familiarity with OntoLex-Lemon; an introduction to the model is given in Appendix A

¹⁵⁰See <https://www.w3.org/2016/05/ontolex/#metadata-lime>

¹⁵¹See <https://www.w3.org/TR/void/>

¹⁵²<https://www.w3.org/2016/05/ontolex/#metadata-lime>

¹⁵³<https://www.w3.org/2016/05/ontolex/#metadata-lime>

¹⁵⁴see <https://www.w3.org/2016/05/ontolex/> for a full description.

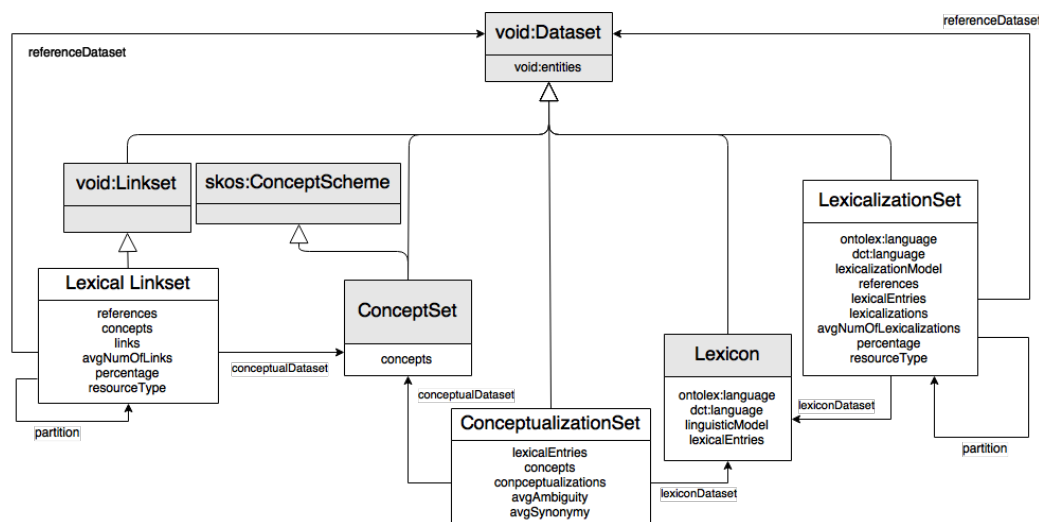


Fig. 5. The *lime* module.

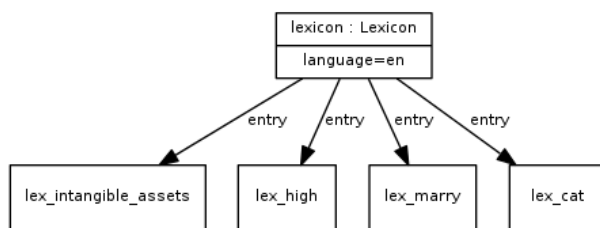


Fig. 6. *lime* Example (diagram taken from OntoLex-Lemon guidelines).

of all kinds of language resources. There are several other linked data vocabularies which are relevant for the description of the metadata of language resources. Aside from the well known Dublin Core vocabulary, as well as the PROV ontology for provenance information, there are more specialised vocabularies which can be applied in specific use cases.

For instance in the case of the the publication of retrodigitised dictionaries and the modelling of historical and scholarly lexical resources as LLD, there is a need for extensive metadata provision at both the lexical and at the individual entry level in order to reflect bibliographic and historic informatio and the representation of scientific-scholarly hypotheses. In the case of retrodigitised LLD dictionaries, metadata for the resources in question should contain information on who compiled the original ‘paper’ version of the dictionary and when, what edition of the dictionary the resource is based on, who digitised it, what tools were used, etc.

There already exist vocabularies for doing a lot of these things (vocabularies not specialised to the language resources domain), however. These include the Semantic Publishing and Referencing suite of ontologies for bibliographic information¹⁵⁵, which can be used in conjunction with META-SHARE and *lime* and others in creating metadata solutions, and potentially application profiles, for the cases which we have mentioned. They also include the CIDOC-CRM family of ontologies. It does however require that whatever solutions are proposed are then made accessible to a wider community of users via the use of guidelines containing typical examples as well as and/or design patterns (for a discussion of the potential use of ontology design patterns in the context of OntoLex-Lemon see Section 6.1.2).

¹⁵⁵<http://www.sparontologies.net/>

4.3.4. Language Identification

The reliable identification of languages and language varieties is of the utmost importance for language resources. For applications in linguistics and lexicography it defines the very scope of investigation and the data provided by a language resource; for applications in language technology and knowledge extraction, language identifiers define the suitability of training data or the applicability for a particular tool for the data at hand.

There are two different ways of encoding language identification information currently being used in RDF. The first is a URI-based mechanism that uses terminology repositories, the other is by attaching a language tag to a literal to indicate its language. In the latter case the language tag is treated similarly to a data type. Language information provided in this way does not entail an additional RDF statement over the literal, allowing a compact, readable and efficient identification with minimal overhead on data modelling. Note that the RDF specifications [117] already include provision for the use of language identification via the attachment of language tags to strings.

In the former case, there exist a number of RDF vocabularies which provide the means to mark the language of a resource explicitly using RDF triples, i.e., using properties such as `dc:language` (for language URIs or string representations) or `lime:language` (for string representations). We elaborate on the differences in practice below.

RDF language codes are defined by BCP47¹⁵⁶ and the IANA¹⁵⁷ registry on the basis of the ISO 639 standard for language codes¹⁵⁸.

For application to RDF data, ISO provides three relevant subsets of language tags: **ISO 639-1**, maintained by the Library of Congress and available as plain text

¹⁵⁶<https://tools.ietf.org/rfc/bcp/bcp47.txt>

¹⁵⁷<https://www.iana.org/>

¹⁵⁸The need for the provision of machine readable identifiers for single languages or language varieties is clear from instances where a language has more than one name. For instance, the Manding language *Bamanakan* (bm) which is also known as *Bambara*. It is also essential for dealing with cases where the same language name is used to refer to what are quite different varieties. Take, for instance, the case of *Saxon* which as well as being an English heavy metal band has also been used to designate both Old English (Anglo-Saxon, ISO 639-3 `ang`) as well as a number of varieties of *Low German*, both historical and modern (Old Saxon, `osx`; Low Saxon, `nds`), along with various different dialects of *High German* (Upper Saxon, `sxu`; Transylvanian Saxon [currently no ISO language code]).

or RDF data,¹⁵⁹ provides an extensive set of two letter codes for major languages that date back to the beginning of the modern-age of computing, but long before the emergence of the internet. ISO 639-1 codes are composed of two lower-case letters with values from *a* to *z* each. In theory, such a system is sufficient to identify up to 676 languages.

Yet, with language technology developing into a truly global phenomenon, it became clear that two-letter codes were not sufficient to reflect the linguistic diversity of the world both past and present – and in the present case this diversity is estimated to comprise more than 6,000 language varieties. As a response to this, **ISO 639-2** provides a set of three-letter codes for (theoretically) up to 17,576 languages. Again, the Library of Congress acts as maintainer and provides the data in human-readable form and as RDF¹⁶⁰. However, it has to be recognised that the primary use case of ISO 639-2 was library-based and focused on languages with an extensive literature, whereas the demands of linguistics and lexicography, especially historical linguistics and language documentation, exceed far beyond this. Indeed they comprise languages that are primarily spoken, not written, but for which field recordings, text books, grammars or word lists must nevertheless be identifiable in order to be retrieved from metadata portals such as, as an example, the Open Language Archives Community (OLAC)¹⁶¹.

For applications in linguistics, SIL International acts as maintainer of **ISO 639-3**, another, and more extensive set of three-letter codes. In distinction to ISO 639-1 and 2 codes, which are meant to be stable and develop at a slow pace, if at all,¹⁶² ISO 639-3 codes are actively maintained by the research community and a continuous process of monitoring, approving (or rejecting) updates, additions and deprecation requests is in place. At the moment, ISO 639-3 codes are published by means of human-readable code tables only,¹⁶³ along with their history and associated documentation, but not in any machine-readable form. Within the LLOD community, it is a common prac-

¹⁵⁹<https://id.loc.gov/vocabulary/iso639-1.html>

¹⁶⁰<https://id.loc.gov/vocabulary/iso639-2>

¹⁶¹<http://www.language-archives.org/>

¹⁶²Changes in ISO 639-1 and 639-2 codes are very rare and occur mostly as a result of political changes, e.g., after the split of Yugoslavia, Serbian (`sr`, `srp`) and Croatian (`hr`, `hrv`) were to be considered independent languages (with two tags) whereas they were previously considered dialects of a single language, Serbo-Croatian (language tag `sh`, deprecated in 2000).

¹⁶³<https://iso639-3.sil.org/>

1 tice to apply the ISO 639-3 codes provided as part of
 2 LexVo [118] whenever language URIs are required and
 3 ISO 639-3 codes are sufficient. However, it is to be
 4 noted that, unlike SIL code tables, LexVo identifiers
 5 are not authoritative and may not be up to date with the
 6 latest version of SIL.

7 But ISO 639-3 only represents the basis for lan-
 8 guage tags as specified by BCP47 [119, Best Common
 9 Practices 47, also referred to as IETF language tags or
 10 RFC 4646] that are incorporated into the RDF specifi-
 11 cation:

12 BCP47 defines how ISO 639 language tags can be
 13 extended with information regarding geographical use,
 14 script, among other variables as follows:

```
15 language(-script) (-region) (-variant) *  
16 (-extension) * (-x-privateuse)
```

17 where:

- 18 – **language**: this is an ISO 639-1 tag if this is avail-
 19 able or an ISO 639-3 tag otherwise;
- 20 – **Script** (optional): an **ISO 15924** 4-letter code, for
 21 instance the code for Latin is `Latn`;
- 22 – **region** (optional): this is an **ISO 3166** 2-letter re-
 23 gion code or a **UN M.49** 3-number) code, for in-
 24 stance either `US` or `840` for the United States of
 25 America
- 26 – **variant**: zero or more registered variants taken
 27 from the current list of registered variants pro-
 28 vided by IANA¹⁶⁴.
- 29 – **extension**: zero or more extensions in one or
 30 more custom schemes
- 31 – **private use** (optional): for use for internal notes
 32 about identification in a single application.

33 The W3C provides means for validating BCP 47
 34 language tags, part of the specification is also that lan-
 35 guage tags should be registered at the Internet As-
 36 signed Numbers Authority. The IANA language sub-
 37 tag registry¹⁶⁵ currently provides registered language
 38 tags in XML, HTML and plain text. As of 2020, dis-
 39 cussions around the provision of a machine-readable
 40 view in RDF and by means of resolvable URIs have
 41 undergoing and are expected to bear fruit in the coming
 42 years. We expect that, by then, the IANA registry will
 43 supersede LexVo as a default provider of ISO 639(-3)
 44 language URIs¹⁶⁶. However, it should be noted that the

45 ¹⁶⁴[https://www.iana.org/assignments/language-subtag-registry/](https://www.iana.org/assignments/language-subtag-registry/language-subtag-registry)
 46 [language-subtag-registry](https://www.iana.org/assignments/language-subtag-registry) (accessed 10-07-2019)

47 ¹⁶⁵[https://www.iana.org/assignments/lang-subtags-templates/](https://www.iana.org/assignments/lang-subtags-templates/lang-subtags-templates.xhtml)
 48 [lang-subtags-templates.xhtml](https://www.iana.org/assignments/lang-subtags-templates/lang-subtags-templates.xhtml)

49 ¹⁶⁶Cf. <https://github.com/w3c/i18n-discuss/issues/13>.

1 very notion of language tags has been criticised as be-
 2 ing both too inflexible as well as unable to address the
 3 needs of linguistics, e.g., recently by [120, 121], and
 4 alternatives are being explored [122].

5 URI-based language identification represents a nat-
 6 ural alternative in such cases, as these are not tied to
 7 any single standardization body or maintainer, but al-
 8 low the marking of both the respective organization or
 9 maintainer of the resource (as part of the namespace)
 10 and the individual language (in the local name). As a
 11 consequence they would naturally support to shift from
 12 one provider to another, if this were required for a par-
 13 ticular task.

14 Finally, another provider of language identifiers
 15 which is relevant for the current discussion is **Glottolog**
 16 [123],¹⁶⁷. This is a repository of identifiers for
 17 language varieties with a specific focus on (although
 18 by no means restricted to) low-resource languages and
 19 with an eye to applications in linguistic typology and
 20 language documentation. Glottolog maintains an inde-
 21 pendent set of language variety identifiers accessible
 22 in human- and machine-readable (RDF) form via re-
 23 solvable URIs, along with additional metadata, an as-
 24 sociated bibliography and a view on the phylogenetic
 25 structure of specific varieties¹⁶⁸. In order to avoid any
 26 of the unintended political connotations that inevitably
 27 arise from the use of the term ‘language’¹⁶⁹, Glottolog
 28 uses the more neutral (though rather uglier) term *lan-*
 29 *guoid*, where this latter is defined as a language variety
 30 about which, or in which, there exists some kind of
 31 literature.

32 A Glottolog ID for a languoid, then, consists of a
 33 4-letter alphabetic code followed by a 4-character nu-
 34 merical code; for instance the Glottolog ID for stan-
 35 dard English is `stan1293`. These are the basis
 36 of resolvable URIs, for instance [http://glottolog.org/](http://glottolog.org/resource/languoid/id/stan1293)
 37 [resource/languoid/id/stan1293](http://glottolog.org/resource/languoid/id/stan1293), which once resolved
 38 provide links to other relevant resources such as ISO
 39 639. Note that Glottolog maintains a certain bias to-
 40 wards endangered modern languages and therefore re-
 41 mains rather sketchy for what concerns the histori-

42 ¹⁶⁷<https://glottolog.org/>

43 ¹⁶⁸That is, Glottolog allows for the specification of the phy-
 44 logenetic relationships between different varieties, specifying En-
 45 glish, for instance, as a subconcept of the category ‘Macro-English’
 46 (`macr1271`), which groups together Modern Standard English and
 47 a number of English Pidgins; and relating it in its turn to narrower
 48 subconcepts such as Indian English (`indi1255`) and New Zealand
 49 English (`newz1240`).

50 ¹⁶⁹Recall Max Weinreich’s famous observation that “a language
 51 is a dialect with an army and a navy”.

cal dimension. Yet the popularity of Glottolog and fact that it has already had a wide uptake beyond the language documentation community, including in Wikipedia, would suggest that provision of identifiers for historical varieties is only a matter of time (pardon the pun).

5. Projects

Summary In this section we will give an overview of a range of different projects that have had an impact or which are currently having an impact on the use or definition of LLD vocabularies. Table 3 gives a summary of the projects discussed below. In Section 5.1 we give a detailed overview of this topic; this overview includes a subsection on recent projects which combine LLD and DH (Section 5.1.1) and an introduction and description of an LLD project matrix given as Figure 7 (Section 5.1.2). Next we describe a series of selected projects in detail. These are (in order of appearance):

- LiODi Section 5.2.1)
- POSTDATA (Section 5.2.2)
- ELEXIS (Section 5.2.3)
- Prêt-à-LLOD (Section 5.2.5)
- NexusLinguaram (Section 5.2.6)

5.1. An Overview

As we mentioned in the introduction, the funding of an ever increasing number of projects in which LLD plays a key role at the international, transnational (including European), national, and regional levels, is evidence of the success of LLD as a means of publishing language resources. These projects also provide us with an important picture of the use which is being made of LLD models and vocabularies across different disciplines and use cases as well as indicating where future challenges may lie. Therefore as part of a task currently being undertaken in the NexusLinguaram COST action (see Section 5.2.6) several of the authors of the current article decided to carry out a survey of research projects in which a significant part of each project was dedicated to making language resources available using linked data or which had LLD as one of its main themes.

The survey has so far been carried out via queries on **CORDIS**¹⁷⁰ and the **OpenAIRE explorer site**¹⁷¹, as

¹⁷⁰<https://cordis.europa.eu/projects>

¹⁷¹<https://explore.openaire.eu/>

well as through a study of the literature and by requesting input from other participants of the NexusLinguaram COST action¹⁷². Our project survey also included an analysis of influential survey publications and anthologies in this area of linguistic linked data (such as [124, 125]) and an analysis of the programs of the major conferences in the sector of language resources¹⁷³. This may of course have, inadvertently, led us towards a natural selection bias in the project overview, namely, towards projects that tended to publish their results at these venues. In addition it should also be noted that since our most important sources of project information were the CORDIS and OpenAIRE project platforms, both of which have a severely limited coverage of national and non-European projects, we were at a disadvantage with respect to information with regards to the latter. We were however able to partially compensate for this by information retrieved via the active consultation of our respective networks.

Based on this exploratory work we were able to make a number of observations. Probably the most important of these is that the effort towards the definition of common models for linguistic linked data has never been dependent on any single, large-scale project, but was conducted within the confines of a much broader community, one which overlapped with a number of funded projects, often carried out in parallel. Over and above this, the community was also maintained by other kinds of networks and initiatives. What also came through quite strongly, however, both from the research carried out as part of the survey and the personal experience of the individual authors of the current work, is that international (and especially European level) projects played a crucial role in **supporting and sustaining** LLD models and vocabularies, *once they had been proposed*.

¹⁷²As part of the preparation for the survey, we set up a Wikipedia page on OntoLex, (<https://en.wikipedia.org/wiki/OntoLex>) and extended another Wikipedia page on Linguistic Linked Open Data (https://en.wikipedia.org/wiki/Linguistic_Linked_Open_Data). We also encouraged partners from our respective networks to contribute and extend those pages, especially with respect to applications of OntoLex and LLOD in general. Information retrieved as part of this process was used to complement the survey described above.

¹⁷³In particular, the Language Resource and Evaluation Conference (LREC) series and associated workshops as well as domain-specific events (workshops on Linked Data in Linguistics (LDL), conferences on Language, Data and Knowledge (LDK), lexicographic events such as EURALEX, ASIALEX, and GLOBALEX as well as the eLex series of electronic lexicography conferences, and associated workshops.

| Summary | | | |
|----------------|-----------|--|-----------------------------|
| Project Name | Duration | Type | Coverage in Current Article |
| EAGLES | 1993-1995 | European Project (FP3) | Section 5.1 |
| ISLE | 2000-2002 | European Project (FP5) | Section 5.1 |
| E-MELD | 2007-2012 | American National Project (NSF) | Section 5.1.2 |
| MONNET | 2010-2013 | European Project (FP7) | Section 5.1 |
| SemaGrow | 2012-2015 | European Project (FP7) | Section 5.1 |
| CLLD | 2013-2016 | German Project (Max Planck) | Section 4.2 |
| LIDER | 2013-2015 | European Project (FP7) | Section 5.1 |
| QTLep | 2013-2016 | European Project (H2020) | Section 5.1 |
| TDWM | 2014-2029 | German Regional Project | Section 5.1.1 |
| FREME | 2015-2017 | European Project (H2020) | Section 5.1 |
| LiODi | 2015-2022 | German Project | Section 5.2.1 |
| Lynx | 2017-2021 | European Project (H2020) | Section 5.1 |
| DiTMAO | 2016-2019 | German-Italian (funded by Deutsche Forschungsgemeinschaft (DFG)) | Section 5.1.1 |
| POSTDATA | 2016-2022 | European Project (H2020-ERC) | Section 5.2.2 |
| MTAAC | 2017-2020 | International (funding from DFG, SSHRC and NEH) | Section 5.1.1 |
| Nénufar | 2017- | French Project (mixed funds) | Section 5.1 |
| ELEXIS | 2018-2022 | European Project (H2020-ERC) | Section 5.2.3 |
| LiLa | 2018-2023 | European Project (H2020-ERC) | Section 5.2.4 |
| Prêt-à-LLOD | 2019-2022 | European Project (H2020-ERC) | Section 5.2.5 |
| NexusLinguaram | 2019-2023 | EU Cost Action | Section 5.2.6 |
| ItAnt | 2020-2023 | Italian National Project (PRIN) | Section 5.1.1 |
| MORdigital | 2021-2024 | Portuguese National Project | Section 5.1.1 |

Table 3

Projects Discussed in the Current Article

This can be demonstrated by the development history of OntoLex-Lemon, probably the most popular and well known of the LLD specific models featured in this article. The origins of this model ultimately go back to the Lexical Markup Framework (LMF) [3], a conceptual Uniform Markup Language (UML)-based model¹⁷⁴ for representing NLP-lexicons and machine-readable dictionaries which was developed over the course of a number of projects that covered lexical resources in NLP and related use cases, most notably **Expert Advisory Group on Language Engineering Standards (EAGLES, 1993-1995)**¹⁷⁵, and **International Standards for Language Engineering (ISLE,**

2000-2002)¹⁷⁶. LMF was subsequently further developed within ISO TC37. And it continues to be so: the latest version of LMF is a multi-part standard of which the first four parts have been all been published at the time of writing [126].

The project **Multilingual Ontologies for Networked Knowledge (MONNET, 2010-2013)**¹⁷⁷ subsequently developed the original *lemon* model on the basis of LMF. In 2011, MONNET project members initiated the formation of W3C Community Group Ontology-Lexica. OntoLex-Lemon was subsequently developed within the ambit of this community group as a revision of the original *lemon* model for the specific application use case of ontology lexicalization. OntoLex-Lemon was further developed in the subsequent **LIDER project (2013-2015)**¹⁷⁸. LIDER contributed to the

¹⁷⁴LMF also had an official XML serialization was included as part of the standard. Attempts towards a RDF/OWL serialization were made by Gil Francopoulo and can be found linked under <http://www.lexicalmarkupframework.org/>, but have not been otherwise published.

¹⁷⁵<http://www.ilc.cnr.it/EAGLES/home.html>

¹⁷⁶http://www.ilc.cnr.it/EAGLES96/isle/ISLE_Home_Page.htm

¹⁷⁷<https://cordis.europa.eu/project/id/248458>

¹⁷⁸<http://lider-project.eu/lider-project.eu/index.html>

1 formation of numerous W3C community groups as
 2 a means of providing a long-term perspective for
 3 its activities. As far as lexical resources are con-
 4 cerned, this included the a W3C Community Group on
 5 **Best Practices for Multilingual Linked Open Data**
 6 **(BP-MLOD)** which, among other contributions, de-
 7 veloped guidelines for the application of OntoLex-
 8 Lemon for modelling lexical resources (dictionar-
 9 ies and terminologies) independently from ontolo-
 10 gies. This represents the basis for most modern uses
 11 of OntoLex-Lemon, and its development towards a
 12 general-purpose community standard for publishing
 13 lexical resources on the Semantic Web.

14 Monnet and LIDER were seminal in their impact
 15 on the development of LLD models and vocabular-
 16 ies. Other important (European) projects in this re-
 17 gard include the FP7 project **Eurosentiment**¹⁷⁹ which
 18 leveraged *lemon* to model language resources for sen-
 19 timent analysis. They also include **FREME**¹⁸⁰ which
 20 explored the application of the NIF and *lemon*; and
 21 **SemaGrow**¹⁸¹ which, along with the LIDER project,
 22 helped to support the development of the *lime* meta-
 23 data module).

24 Other projects with a significant recent impact on
 25 the application of LLD vocabularies include the Hori-
 26 zon 2020 project **Lynx: Building the Legal Knowl-
 27 edge Graph for Smart Compliance Services in Mul-
 28 tilingual Europe** (2017-2021) [127] which has con-
 29 tributed to data modelling in the area of machine-
 30 readable licensing, a topic that is much broader than
 31 the area covered by our survey. They also include the
 32 project **Quality Translation by Deep Language En-
 33 gineering Approaches (QTLeap)**, 2013-2016 which
 34 has primarily focused on Natural Language Process-
 35 ing. Due to their more recent impact on the definition
 36 and use of LLD models and vocabularies we will dedi-
 37 cate specific sections to the following European H2020
 38 projects **ELEXIS** (Section 5.2.3), **Prêt-à-LLOD** (Sec-
 39 tion 5.2.5), and the ERC projects **LiLa** (Section 5.2.4)
 40 and **POSTDATA** (Section 5.2.2) below.

41 5.1.1. Recent Projects combining LLD and DH

42 As is typical of DH projects, those which we de-
 43 scribe in this section, along with ELEXIS, LiLa and
 44 POSTDATA (described in their own sections below),
 45 aim to engage with a wide and diverse scholarly com-
 46 munity, which includes linguists, philologists, histori-
 47

1 ans and archaeologists, as well as language learners. In
 2 the case of the classics (the case of LiLa in particular
 3 Section 5.2.4), there is also a reliance on a strong and
 4 lengthy tradition of previous scholarship.

5 However by making it easy to, for instance, empha-
 6 sise different kinds of connections both within and be-
 7 tween different past civilizations, their languages and
 8 cultures LLD offers a powerful and effective solu-
 9 tion to the challenges of modelling heterogeneous hu-
 10 manities data and making it findable and interopera-
 11 ble. In particular LLD is well placed to facilitate the
 12 integration of historical and geographical with lexi-
 13 cographic and linguistic information – using the ap-
 14 propriately defined classes and properties that is. The
 15 use of linked data made in DH projects such as Pела-
 16 gios [82], Mapping Manuscript Migrations [128] and
 17 in the Finnish Sampo datasets [129], among others
 18 very clearly demonstrates this.

19 In the rest of this section, then, we will provide sum-
 20 maries of a number of small and medium scale projects
 21 at the regional, national and trans-national levels which
 22 are notable for bringing together, or being at the over-
 23 lap of, LLD and DH.

24 At a national level, we can list the French project
 25 **Nénufar**, already mentioned above, which aims to-
 26 wards the creation of successive early 20th century edi-
 27 tions of the French language **Le Petit Larousse Il-
 28 lustré** dictionary in both TEI/XML and in RDF us-
 29 ing OntoLex-Lemon [130]¹⁸², along with the Ger-
 30 man project **Linked Open Dictionaries** which is de-
 31 scribed in detail in Section 5.2.1. In addition we can
 32 also mention the Italian project (part of the Progetti
 33 di Rilevante Interesse Nazionale or PRIN program)
 34 **Languages and Cultures of Ancient Italy. Histori-
 35 cal Linguistics and Digital Models (ItAnt)** (currently
 36 ongoing) which aims to publish a linked data lexicon
 37 of the ancient Italic languages¹⁸³ using the OntoLex-
 38 Lemon model and its extensions. Also relevant here
 39 is the Italo-German project **DiTMAO**, funded by
 40 the DFG (Deutsche Forschungsgemeinschaft), (com-
 41 pleted) which produced a lexicon of Old Occitano
 42 medical terminology and which also proposed an ex-
 43 tension of *lemon* to deal with the specifics of this use-
 44 case¹⁸⁴ [131]. Another national project worth men-
 45 tioning here is the recently initiated Portuguese na-
 46 tional project **MORdigital** [132] which has the aim of
 47

48 ¹⁷⁹<https://cordis.europa.eu/project/id/296277>

49 ¹⁸⁰<https://cordis.europa.eu/project/id/644771>

50 ¹⁸¹<https://cordis.europa.eu/project/id/318497>

51 ¹⁸²Despite the best of intentions however the RDF part isn't cur-
 52 rently very well developed.

53 ¹⁸³<https://www.prin-italia-antica.unifi.it/>

54 ¹⁸⁴<https://www.uni-goettingen.de/en/ditmao/487498.html>

1 digitising the historically significant 18th century Por-
 2 tuguese language dictionary *O Dicionario da Lingua*
 3 *Portuguesa* by António de Morais Silva with the in-
 4 tention of producing digital editions of this important
 5 lexicographic work both in TEI Lex-0 and OntoLex-
 6 Lemon. The MORdigital project will be an impor-
 7 tant test case for understanding both the coverage of
 8 already existing LLD vocabularies when it comes to
 9 retrodigitized dictionaries, and the advantages and dis-
 10 advantages of using linked data as a means of publish-
 11 ing such data.

12 Many of the projects we have mentioned have used
 13 OntoLex-Lemon or its predecessor *lemon*. However a
 14 lightweight alternative to these vocabularies, and one
 15 which enables the multilingual annotation of concep-
 16 tual hierarchies, is SKOS-XL. For instance SKOS-
 17 XL has been used in a number of related projects at
 18 the Computational Linguistics lab of the University of
 19 Saarland, part of a major effort towards the transforma-
 20 tion of several influential classification schemes in the
 21 field of folk literature¹⁸⁵ (including among others folk-
 22 tales, ballads, myths, fables) into Semantic Web repre-
 23 sentation languages in order to support interoperabil-
 24 ity between those schemes; see [136]. The terms used
 25 in the original classification schemes were transformed
 26 into (multilingual) SKOS-XL labels. These were used
 27 for encoding folktale text sequences, extracted from
 28 a manually annotated multilingual folktale corpus and
 29 which had been identified as representing motifs listed
 30 in [134]. The use of SKOS-XL meant that motifs could
 31 be annotated in different multilingual versions of tales.

32 Staying with the theme of the use of a range of
 33 different models/vocabularies for the encoding of lin-
 34 guistic data as LLD (and not necessarily the ones we
 35 have focused on in this survey) another project worth
 36 mentioning here is **Text Database and Dictionary of**
 37 **Classic Mayan (TDWM)**¹⁸⁶ (2014-2029). This latter
 38 aims to develop a corpus-based dictionary of Mayan
 39 hieroglyphic writing alongside a near-exhaustive cor-
 40 pus of Classic Mayan which would allow for the verifi-
 41 cation of different textual interpretations and aid in fin-
 42 ishing the decipherment of Maya writing. The project
 43 faces the problem, typical of ancient languages, of the
 44 need to represent multiple interpretations of characters
 45 and texts, in part due to damaged sources in concomi-
 46 tance with the need to update this data with the in-

185The classification schemes in question were those proposed by
 49 Vladimir Propp [133], Stith Thompson [134] and Anti Aarne, Stith
 50 Thompson and Hans-J. Uther [135].

186Based at the University of Bonn, Germany.

1 conclusion of new data during dictionary development. In
 2 the case of the Mayan the situation is even more dif-
 3 ficult due to the signs not yet having been fully deciphered.
 4 In order then to deal with the challenges which
 5 arise from the existence of different sign catalogues
 6 (which might cluster different signs into meanings dif-
 7 ferently) and the necessity of linking with other cat-
 8 alogues which have been developed in the field, the
 9 project's sign catalogue has been formalised in SKOS
 10 in addition to using properties and concepts from the
 11 CIDOC-CRM vocabulary¹⁸⁷ and GOLD (mentioned
 12 above in Section 3.5). The TDWM project is also de-
 13 veloping its own vocabulary for identifying signs, link-
 14 ing them to different sign catalogues, possible read-
 15 ings, graphical variants, etc. At the time of writing, nei-
 16 ther the sign catalog nor any texts are publicly avail-
 17 able, but Diehr et al. [137] provide a detailed descrip-
 18 tion¹⁸⁸.

19 Finally, another recent project which used a range of
 20 different LLD vocabularies is **Machine Translation**
 21 **and Automated Analysis of Cuneiform Languages**
 22 **(MTAAC)**. This was a Data international funded
 23 project which saw the collaboration of specialists of
 24 cuneiform languages and computational linguists in
 25 the development of cutting edge tools for the annota-
 26 tion and distribution of linguistic data of the cuneiform
 27 corpus. Although the project's overall objective was
 28 to open the way to the development of tools and the
 29 production of richer linguistic data for all cuneiform
 30 languages, its specific focus was on a group of unan-
 31 notated Sumerian texts issued from the bureaucratic
 32 apparatus of the Ur III period (21st century BC)¹⁸⁹;
 33 in addition, another corpus composed of royal inscrip-
 34 tions in the Sumerian language [138], annotated with
 35 morphology, was also employed. Amongst the several

187<http://www.cidoc-crm.org/>

188It is interesting to note that TDWM stands in a longer tra-
 38 dition of projects in the Digital Humanities that aim to comple-
 39 ment a TEI/XML edition with terminology management using an
 40 ontology. Similar ideas have already been driving force behind the
 41 project *Sharing Ancient Wisdoms (SAWS, 2010-2013)*(<http://www.ancientwisdoms.ac.uk/>), a joint project at King's College London,
 42 UK, the Newman Institute in Uppsala, Sweden, and the University
 43 of Vienna, Austria, funded in the context of the Humanities in the
 44 European Research Area (HERA) program to facilitate the study and
 45 electronic edition of ancient wisdom literature. Both projects employ
 46 resolvable URIs, but the linking is expressed by means of narrowly
 47 defined TEI/XML attributes rather in terms of RDF semantics. In
 48 that regard, the data published in accordance with these guidelines
 49 does not qualify as Linked Data, but can still be converted to Linked
 50 Data with moderate effort.

189These texts were extracted from CDLI

objectives of the project [139] was the aim of formalising the new data produced by the project by utilising Linked Open Data (LOD, including Linguistic LOD) vocabularies, and fostering the practices of standardisation, open data and LOD as integral to projects in digital humanities and computational philology. The project, which ended in 2020 successfully achieved these aims, including making new data in the form of linguistic annotations and translations available under open licenses¹⁹⁰; this will soon be accessible through the new web platform of the Cuneiform Digital Library Initiative (CDLI <https://cdli.ucla.edu>) in many forms, including (L)LOD.

CoNLL was chosen as a flexible and robust format for storing the multi-layer annotations which were produced and worked on as part of the project¹⁹¹. However CoNLL-RDF was also employed in the project in order to ensure integration with LLOD, as well as for easier querying and transformation and was used to link annotations, lexical information, and metadata. The ETCSRI morphological annotations¹⁹² were mapped to Unimorph¹⁹³ using Turtle-RDF¹⁹⁴, rendering Sumerian material accessible for cross-linguistic queries. SPARQL was leveraged through CoNLL-RDF for syntactic annotation which was mapped to Universal Dependencies for POS and dependency labels. Lexical data was linked to guide word entries through the employment of an OntoLex-Lemon compliant index. Metadata concerning the analysis of the medium of the text and other meta classifications of the texts were mapped to the CIDOC-CRM. Overall, MTAAC succeeded in preparing a (L)LOD edition and linking of Sumerian language corpora. The model can be extended in part to other cuneiform languages. Various Assyriological resources had been integrated using (L)LOD [98]: The CDLI data, (CoNLL-RDF plus CIDOC-CRM), ORACC:ETCSRI (by conversion; CoNLL-RDF), ePSD (by conversion and links to HTML; lemon) and ModRef & BM (by federation; CIDOC-CRM). Other vocabularies are planned

¹⁹⁰<https://gitlab.com/cdli/framework>; <https://github.com/cdli-gh>

¹⁹¹A derivative internal format, called CDLI-CoNLL is employed to store the data locally – this was an essential step to support the preservation of domain specific annotation which are richer than their counterparts found in linguistic all-encompassing models. But this can be exported in CoNLL-U format, as well in Brat Standalone format, for better compatibility.

¹⁹²<http://oracc.museum.upenn.edu/etcsri/parsing/index.html>

¹⁹³<http://unimorph.org/>.

¹⁹⁴https://github.com/cdli-gh/mtaac_work/blob/master/lod/annotations/um-link.ttl.

to be added in the future (Pleiades, perio.do, etc.). The model developed is currently being integrated into the CDLI platform.

5.1.2. An LLD Project Matrix; The Relationship between Projects and Community Initiatives

In Figure 7, we provide an overview in the form of a matrix of a number of selected projects on the basis of the kinds of contributions they have made to a number of LLD vocabularies. We distinguish three varieties of contribution: namely, a project is said to have

Developed (deep green) a vocabulary if the development of that vocabulary was a designated project goal, to have

Contributed (light green) to a standard if vocabulary development was not a designated project goal, but the project provided a use case or application that was discussed in the process of its development, or to have

Used (yellow) a vocabulary if they applied an existing vocabulary, worked with or produced data of that type

Note that this survey, and indeed any survey which focuses on projects, will provide a partial view only. In particular, contributions by community groups (Open Linguistics Working Group, W3C working groups, etc.) are not explicitly covered in this section (although they are described in some depth in Section 4 and their contribution is also discussed in Sections 2.2 and 5.1). For instance the reader will notice that very few of the projects in Fig. 7 address the area of LLD for linguistic typology. In fact the interaction between linguistic typology and language technology operates primarily on informal contacts on mailing lists and via workshops and less in terms of large-scale infrastructural projects, and that, thus, the development of standard (computational) models and vocabularies has only rarely a priority in typological projects¹⁹⁵. For such discussions,

¹⁹⁵There are notable exceptions here the **E-MELD project** (<http://emeld.org/>), for example, developed the GOLD ontology as part of an attempt to improve interoperability and sustainability of language documentation material. But while several typological projects developed ontologies and RDF vocabularies, and have been actively contributing to the community, esp., in the Open Linguistics working group, we see a very limited degree of linking between such resources. The **Cross-Linguistic Linked Data project (CLLD, <https://clld.org/>)**, for example, does provide an RDF view on their data, but linking is primarily internal, and neither complete data dumps nor a SPARQL end point or any form of an API is provided. Instead, their RDF data seems to be generated on the fly, without any links to external resources. We take this to reflect the fact that for



Fig. 7. Usage of and contribution to major LLOD vocabularies by selected research projects

more informal networks present great opportunities (and act as a driver) for experts to participate in the Linguistic Linked Open Data movement, whereas the chances for acquiring substantial funding directed towards vocabulary development and community participation are rather unreliable (if past experience is anything to go).

Also note that in this section, we have concentrated on research projects with a specific focus on linguistic linked (open) data – many of them, indeed, with industry partners involved – but which do not, for the most part, directly target industrial applications. More industry-focused LLD projects do exist, however, and are the basis for businesses specialising in text analytics [140], terminology and knowledge management [141] or lexicography [142]. But linked data in these contexts tends to be viewed as a technical facet that has an impact on interoperability, (re)usability and information aggregation rather than being fundamental for any existing business model. With the increasing maturity of the technology, however, this may change over the longer term, especially in the area of establishing interoperability between AI platforms [143], their providers and users and data provided and exchanged between them [144].

To conclude then, it has really been the combination of open community initiatives and projects that has de-

termined the success and then the subsequent maintenance of the LLDs models and vocabularies. The importance of funded projects is clear for the development of tools and hosting solutions for Linguistic Linked (Open) Data which are not yet in place; open community initiatives have also proven themselves vital for dissemination and wider community engagement. With the increasing maturity of OntoLex-Lemon and the convergence between existing solutions in linguistic annotation, the necessary requirements for developing large-scale Linguistic Linked (Open) Data infrastructures and their respective linking are in place, now. In Section 6.1.1 below we take a brief look at the prospects for the involvement of research infrastructures in the kinds of initiatives mentioned in this section.

In what follows we will give extended descriptions of six ongoing projects. We have chosen these projects the basis of their impact and their scale, as well as for their importance in the development of LLD models and vocabularies and/or in their innovative use of such. These are **LiODi** in Section 5.2.1; **POST-DATA**, in Section 5.2.2; **Prêt-à-LLOD** in Section 5.2.5; **ELEXIS** in Section 5.2.3 and finally **NexusLinguarum** in Section 5.2.6. Please note that the length of the following project descriptions will vary on the basis of their relevance to some of the models and vocabularies discussed in the rest of this paper.

this community, interoperability is a priority, but also, to maintain control over internal data and independence from external contributions.

5.2. Innovative Projects

5.2.1. LiODi (2015-2022)

The **Linked Open Dictionaries project (LiODi)**¹⁹⁶ aims to develop LLOD-enabled methodologies and infrastructures to facilitate language research for low-resource languages, validating these developments mostly on the languages of the Caucasus. Within the project, a set of loosely connected tools are being created with the aim of facilitating language contact studies over lexical and corpus data. One of the primary development goals of the project is an environment for detecting semantically and phonologically similar words across different languages in order to facilitate the detection of possible cognates. Other tools include interfaces for converting, validating, and exploring linguistic data to aid in linguistic research both within and outside of the project. Tool development and linguistic research are both integral parts of LiODi and the tools and pipelines implemented are tested on the data generated and used in the project [98, 145].

The most important contributions of LiODi from a modelling perspective relate to the fact that its members have developed, and are in the course of developing, LLD vocabularies for a wide-range of applications in the language sciences: in particular, vocabularies with an emphasis on the requirements of low-resource languages and especially morphologically rich languages which are not well served by existing formats. These vocabularies include individual, task-specific vocabularies such as *Ligt* and *CoNLL-RDF* (see 4.2.4), but also an extension of *OntoLex* for diachronic relations (cognate and loan relations) [44]. In addition to that, the LiODi project (along with *Prêt-à-LLOD*, see 5.2.5) is the main contributor to the **ACoLi Dictionary Graph** [34]¹⁹⁷. To the best of our knowledge, the ACoLi Dictionary Graph currently represents the most extensive collection of machine-readable bilingual open source dictionaries available, with currently more than 3000 substantial data sets for more than 430 ISO 639-3 languages (including full *OntoLex-Lemon* editions of **PanLex**¹⁹⁸, **Apertium**¹⁹⁹, **FreeDict**²⁰⁰, **MUSE**²⁰¹, **Wikidata**²⁰², the **Open Mul-**

tilingual WordNets²⁰³, the **Intercontinental Dictionary Series**, **XDXF**²⁰⁴ and **StarDict**²⁰⁵ – the latter only to the extent that the copyright could be clarified and an open license was confirmed).

More significant than lexical resources and novel vocabularies, however, are the contributions of LiODi to the development of community standards for LLD vocabularies. This includes, among other aspects, significant contributions to the emerging *OntoLex Morphology* module (Section 4.1.2), initiating and moderating the development of the *OntoLex FrAC* module (Section 4.1.3) and the *LD4LT* initiative on harmonizing vocabularies for linguistic annotation on the web.

Furthermore, LiODi has a strong dedication to the dissemination and promotion of linked data approaches for linguistics. The project co-organised two summer schools, **SD-LLOD 2017** and **SD-LLOD 2019**; two conferences **LDK 2017** and **LDK 2019**; three workshops **LDL 2016**, **LDL 2018**, and **LDL 2020**; and collaborated with international partners and the *Prêt-à-LLOD* project (see Section 5.2.5) in the publication of the first monograph on the topic [124] along with a number of edited volumes (not counting the five volumes of proceedings which resulted from the aforementioned events, including a collection on linked data for collaborative, data-intense research in the language sciences [146]).

Outside of conjoined activities at summer schools and datathons, the project supports numerous external partners in expertise with data modelling and language resource management. Indeed LiODi has close ties with most of the projects listed here. To mention one notable example here, a collaboration with the *POSTDATA* project (see next section) and the *Academy of Sciences in Heidelberg*, Germany, led to the first practical applications of RDFa within TEI editions in the *Digital Humanities* [67, 147], and ultimately to the development of an official TEI+RDFa customization (see above).

5.2.2. POSTDATA (2016-2021)

The **Poetry Standardization and Linked Open Data (POSTDATA)** project²⁰⁶, seeks to bridge the digital gap between traditional cultural assets and the growing sophistication of data modelling and publication practices in the field of the *Digital Humanities*.

¹⁹⁶<https://acoli-repo.github.io/liodi/>

¹⁹⁷<https://github.com/acoli-repo/acoli-dicts>

¹⁹⁸<https://panlex.org/>

¹⁹⁹<https://www.apertium.org/>

²⁰⁰<https://freedict.org>

²⁰¹<https://github.com/facebookresearch/MUSE>

²⁰²<https://www.wikidata.org/>

²⁰³<http://compling.hss.ntu.edu.sg/omw/>

²⁰⁴<https://sourceforge.net/projects/xdxf/>

²⁰⁵<http://stardict.sourceforge.net/>

²⁰⁶<http://postdata.linhd.uned.es>

1 It focuses on poetry analysis, bringing Semantic Web
 2 standards and technologies for enhanced interoperabil-
 3 ity to bear on numerous different poetry-related re-
 4 sources. The project is founded upon two central pil-
 5 lars: the use of linked open data and the implementa-
 6 tion and use of a set of dedicated Natural Language
 7 Processing (NLP) tools, **PoetryLAB**. In particular one
 8 of the aims of the project is to share scholarly knowl-
 9 edge about the domain of poetry and publish literary
 10 works on the linked open data cloud.

11 In order to fulfil this aim, POSTDATA is developing
 12 a poetry ontology. This ontology is based on the anal-
 13 ysis and comparison of different data structures and
 14 metadata arising from eighteen projects and databases
 15 devoted to poetry in different languages at the Euro-
 16 pean level, [148–151]. The POSTDATA ontology is an
 17 *encapsulated ontology model*, where domain knowl-
 18 edge is implemented in 3 layers: **Postdata-core**, **Post-**
 19 **data metrical and literary analysis** and **Postdata-**
 20 **transmission**. This layered ontology is based on the
 21 re-use of other ontologies relevant to the project’s do-
 22 main of interest and covers different levels of descrip-
 23 tion from the abstract concept of the poetry work to its
 24 bibliographic representation [152–156]. The model is
 25 intended to support tasks associated with the analysis
 26 of poetry such as close reading, distant reading or crit-
 27 ical analysis. All of these ontologies are modelled in
 28 OWL and will be exposed via SPARQL endpoints.

29 The POSTDATA metrical layer represents knowl-
 30 edge pertaining to the poetical structure and prosody of
 31 a poem and contains salient (general) linguistic, pho-
 32 netic and metrical concepts. From the metrical point of
 33 view, a poem is formed by *stanzas* that contain *lines*,
 34 where the latter is understood as a list of *words*. The
 35 concept, *word*, is present in OntoLex-Lemon and in
 36 NIF. In both cases, the definition of the concept is
 37 not sufficient to capture all the knowledge needed for
 38 the analysis and description of a word from a metri-
 39 cal point of view. From this latter point of view the
 40 concept *word* is associated both with more general lin-
 41 guistic information (such as its *lemma* or headword)
 42 as well as more specific phonetic features such as *syl-*
 43 *lable*, *foot*, *feet type onset or coda* as well as other
 44 types of metrical information. However, the intention
 45 is to link the *Word* concept in the POSTDATA metri-
 46 cal ontology with the OntoLex-Lemon concept *Word*
 47 through the property *wordsense*, allowing us to cap-
 48 ture the range of meanings of the concept. Moreover,
 49 the POSTDATA *Word* class will also be linked to the
 50 NIF *Word* class due to the shared relationship of both
 51 of them to NLP operations.

1 The second pillar of POSTDATA the use of NLP
 2 tools, represented by PoetryLab, encompasses the sev-
 3 eral different levels of poetry scholarship, from the
 4 most formal analyses relating to scansion, to more
 5 cognitive levels which concern the understanding of
 6 metaphor as well as others related to knowledge and
 7 subjective perception involving AI techniques. POST-
 8 DATA has already implemented the first level of NLP
 9 algorithms for poem analysis. These allow for the au-
 10 tomated extraction of information from poems at differ-
 11 ent levels of description and include an Name-Entity
 12 Recognition system (NER) for medieval place names
 13 and organizations, [157] as well as automatic enjamb-
 14 ment analysis and basic metrical scansion tools (which
 15 allow for lexical syllabification and the recognition of
 16 stressed and unstressed syllables) testing different ap-
 17 proaches. These latter range from traditional ruled-
 18 based systems to the latest deep learning based tech-
 19 niques, [158–160]. The goal in this case is to use the
 20 results of these tools in order to build an RDF knowl-
 21 edge graph that is compliant with Postdada ontology.

5.2.3. ELEXIS (2018-2022)

22 Following on from the European Network for e-
 23 Lexicography COST Action²⁰⁷, the ELEXIS project
 24 is currently undertaking the construction of a Euro-
 25 pean infrastructure for electronic lexicography [161].
 26 LLD will play a key role in the ELEXIS infrastruc-
 27 ture, as means of connecting dictionaries and other lex-
 28 icographic resources both within and across language
 29 boundaries. Indeed, the idea of ELEXIS is to eventu-
 30 ally construct a network of interlinked electronic lex-
 31 ica and other lexicographic and language resources in
 32 several different languages, a network that the project
 33 calls a **Matrix Dictionary**. Another relevant aspect of
 34 the project concerns the conversion of legacy lexico-
 35 graphic resources into structured data, and potentially,
 36 linked data in order to feed into this Matrix Dictionary.

37 The main models being used in the project are
 38 OntoLex-Lemon and the TEI Lex-0 model mentioned
 39 above [162]. Here it will perhaps be useful to give a
 40 brief description of the latter.

41 *TEI Lex-0 and ELEXIS* TEI Lex-0 is a customiza-
 42 tion of the TEI schema²⁰⁸ that is adapted to the encod-
 43 ing of lexical resources. In particular it was designed
 44 to enhance the interoperability of such datasets by lim-
 45 iting the range of encoding possibilities (offered by the

²⁰⁷<https://www.elexicography.eu/>

²⁰⁸<https://dariah-eric.github.io/lexicalresources/pages/TEILex0/TEILex0.html>

current TEI guidelines) in the representation of lexical content (for instance TEI Lex-0 has deprecated elements such as `superEntry` or `entryFree`). This makes the possibility of a crosswalk from (at least a subset of) TEI Lex-0 to OntoLex-Lemon more feasible than, say, a crosswalk from say, a minimal customisation of TEI based on the TEI dictionary guidelines to OntoLex-Lemon. TEI Lex-0 is being developed by a special working group which (pre-Covid) organised regular in-person training schools with support from ELEXIS too. Both OntoLex-Lemon and TEI Lex-0 have been previously used for smaller lexicography projects, but never in a project with such wide coverage in terms of the languages and kinds of lexicographic resource under consideration. ELEXIS has provided support to the development of both OntoLex-Lemon as well as TEI Lex-0 and a joint workshop was held between these projects at the 2019 edition of the e-lexicography convention eLex. Work is also underway on a crosswalk between TEI Lex-0 and OntoLex-Lemon. The latest version of a proposed TEI Lex-0 to OntoLex converter can be found at <https://github.com/elexis-eu/tei2ontolex>.

The project is also promoting the standardisation of OntoLex-Lemon and TEI Lex-0 through the OASIS working group on Lexicographic Infrastructure Data Model and API (LEXIDMA)²⁰⁹, which will lead to a new unifying standard for lexicographic data that will be serialised in both OntoLex-Lemon and TEI Lex-0.

The Impact of ELEXIS on the Use of OntoLex-Lemon
ELEXIS aims to provide support for the creation and editing of dictionary resources using both models. To this end extensive teaching materials are also being developed as part of the project with the aim of introducing lexicographers to linked data and the OntoLex-Lemon model. It should be noted that the availability of manuals and targeted teaching materials plays an important factor in increasing the uptake of models such as OntoLex-Lemon and technologies such as linked data, (as of course is the case with new technologies and new technological approaches in general), especially amongst users who haven't had much previous exposure to linked data or conceptual modelling. The original designers of such models are usually unable to take into consideration of every kind of use-case for which the model might be used. The gap between a general purpose model as it is presented in

²⁰⁹https://www.oasis-open.org/committees/tc_home.php?wg_abbrev=lexidma

some final set of guidelines, and its use or appropriation (along with other pertinent models and vocabularies) in a specialist domain or task can be bridged by such targeted training materials (and also by design patterns see Section 6.1.2. This is also one of the motivations behind the strong emphasis on training in NexusLinguaram (see Section 5.2.6).

Both the production of training materials and the push to promote OntoLex-Lemon as a common serialisation format for a standard for e-lexicography seems to promise much in terms of the future use of linked data in the context of electronic lexicography. It is inevitable that the experiences of lexicographers and linguists in using OntoLex-Lemon (and its lexicographic extension, see Section 4.1.1) both within and outside of the ELEXIS project to create and edit lexicographic resources will have an important impact on the use of the model and also, potentially, on future extensions and/or versions of OntoLex-Lemon.

5.2.4. LiLa (2018-2023)

The **LiLa: Linking Latin** ERC project²¹⁰ aims to connect language resources developed for the study of Latin, bringing the worlds of textual corpora, digital libraries, lexica and tools for Natural Language Processing together. To this end, LiLa makes use of the LOD paradigm and of a set of ontologies from the LLD cloud to build an interoperable network of resources. LiLa's ambition is to create an infrastructure where researchers and students of Latin can find answers to complex questions that involve multiple layers of linguistic annotations and knowledge, such as: *what subjects are constructed with verbs formed with a certain prefix [163]? What WordNet synsets do they belong to?*

As Latin is characterised by a very rich morphology (where, for instance, a single verb can potentially yield more than 100 forms, excluding the nominal inflection of participles), LiLa focuses on lemmatization as the key task that allows for a meaningful and functional connection between the different layers of annotation and information involved in the project. Indeed, while lemmas are used by lexica to label entries, lemmatization is often performed in digital libraries of Latin texts to index words and is included in most NLP pipelines (like e.g. UDPipe)²¹¹ as a preliminary step for more advanced forms of analysis²¹².

²¹⁰<http://lila-erc.eu>

²¹¹<https://ufal.mff.cuni.cz/udpipe>.

²¹²For the state of the art in automatic lemmatization and PoS tagging for Latin, see the results of the first edition of **EvaLatin**, a cam-

LLD standards such as OntoLex-Lemon (see Section 4.1) provide an adequate framework to model the relations between the different classes of resources via lemmatization, while also offering a robust solution for modelling most lexica. The central component in LiLa's framework, the gateway between the different projects, is the collection of canonical forms that are used to lemmatize texts (called **lemma bank**). This collection was created starting from the lexical database of the morphological analyzer **Lemlat** [165], and currently includes a set of about 190,000 forms that can potentially be used as lemmas in corpora or lexica²¹³.

The forms in the lemma bank are described in an OWL ontology that reuses several concepts from the LLD standards discussed in the previous sections. The canonical forms are instances of the class **Lemma**, which is defined as a subclass of the **Form** from the OntoLex-Lemon vocabulary. The part-of-speech and morphological annotations in the Lemlat database have been included in the ontology and linked to the OLiA reference model (see Section 3.5). For a selection of circa 36,000 lemmas, the lemma bank also includes derivational information, listing the morphemes (i.e. the prefixes, affixes and lexical bases) that can be identified in each lemma [166].

The fact that OntoLex-Lemon forms are allowed to have multiple written representations is a particularly helpful feature for a language which has been attested across circa 25 centuries and a wide spectrum of genres, and which is, moreover, characterised by a substantial amount of spelling variation. Harmonizing different lemmatization solutions adopted by corpora and NLP tools, however, requires practitioners to deal with other kinds of variation as well [167]. In the case of words with multiple inflectional paradigms or forms which may be interpreted as either autonomous words or inflected forms of a main lemma (such as participles, or adverbs built from adjectives: see e.g. English "quickly" from "quick"), projects may vary considerably in the adopted strategies. For this reasons, the LiLa ontology introduces one sub-class of the **Lemma** and two object properties that connect forms to forms.

paign devoted to the evaluation of NLP tools for Latin [164]. The first edition of *EvaLatin* focused on two shared tasks (i.e. lemmatization and PoS tagging), each featuring three sub-tasks (i.e. Classical, Cross-Genre, Cross-Time). These sub-tasks were specifically designed to measure the impact of genre variation and diachrony on NLP tool performances.

²¹³The lemma bank can be queried using the lemmaBank SPARQL endpoint of the project: <https://lila-erc.eu/sparql/>.

The property lemma variant connects two lemmas that can be alternatively used to lemmatize forms of the same words. Hypolemma is a sub-class of the **Lemma** that groups forms (e.g. participles) that can be either promoted to canonical or be lemmatised under a hyperlemma (e.g. the main verb); hypolemmas are connected to their hyperlemma via the *is hypolemma* property.

Currently, the canonical forms in the LiLa lemma bank connect lexical entries of four lexical resources. Two lexica provide etymological information, which was modelled using the OntoLex-Lemon extension *lemonEty* [54], respectively on the lexicon inherited from proto-Indo-european²¹⁴ [168] and loans from Greek²¹⁵ [169]. The polarity lexicon *LatinAffectus* connects a polarity value (expressed using the **Marl** ontology²¹⁶) to a general sense for 1,998 entries²¹⁷ [170]. Finally, 1,421 verbs from the *Latin WordNet* have been manually revised and published as *LOD*²¹⁸ [171].

In addition to lexica, two annotated corpora are currently linked to the LiLa lemma bank. The *Index Thomisticus* Treebank²¹⁹ provides morpho-syntactic annotation for 375,000 tokens from the Latin works of Thomas Aquinas (13th century CE), while the *Dante Search* corpus²²⁰ includes the lemmatised text of four Latin works of Dante Alighieri (14th century), which are currently undergoing a process of syntactic annotation following the Universal Dependencies annotation style [172]²²¹. The POWLA ontology was used to represent texts and annotations for both corpora. However, the link between a corpus token and a lemma of the LiLa collection was expressed using a custom property has lemma defined in the LiLa ontology²²², which takes an instance of the **Lemma** class as its range, since no existing vocabulary provided a suitable way to express this relation.

5.2.5. Prêt-à-LLOD (2019-2022)

The goal of the **Prêt-à-LLOD** project is to make linguistic linked open data 'ready-to-use' and part of

²¹⁴<https://lila-erc.eu/data/lexicalResources/BrillEDL/Lexicon>

²¹⁵<https://lila-erc.eu/data/lexicalResources/IGVLL/Lexicon>

²¹⁶<http://www.gsi.upm.es:9080/ontologies/marl/>

²¹⁷<https://lila-erc.eu/data/lexicalResources/LatinAffectus/>

Lexicon

²¹⁸<http://lila-erc.eu/data/lexicalResources/LatinWordNet/Lexicon>

con

²¹⁹<http://lila-erc.eu/data/corpora/ITTB/id/corpus>

²²⁰<http://lila-erc.eu/data/corpora/DanteSearch/id/corpus>

²²¹<https://universaldependencies.org/guidelines.html>

²²²<https://lila-erc.eu/lodview/ontologies/lila/>

1 this mission is to contribute to the development of
 2 new vocabularies for linguistic linked data in appli-
 3 cation scenarios that facilitate the development of a
 4 next-generation multilingual internet. Several aspects
 5 of linked data technology are being pursued in this
 6 context. This includes, without being restricted to

7
 8 **linking** In its linking aspect, Prêt-à-LLOD explores
 9 technologies to facilitate the linking between and
 10 among lexical, terminological and ontological re-
 11 sources. In this context, it has provided signifi-
 12 cant support to the development of the OntoLex-
 13 Lemon, including the development of a module
 14 for lexicography, a module for morphology, and
 15 corpus information (all of which are discussed in
 16 Section 4.1). Further extensions for terminologies
 17 and linking metadata (Fuzzy Lemon) have been
 18 proposed in the context of the project, as well.
 19 In addition, the project is contributing models for
 20 dataset linking to the Naisc project²²³ that pro-
 21 vides a toolkit for generic dataset linking.

22 **transformation** Prêt-à-LLOD provides a generic frame-
 23 work for transforming, enriching and manipulat-
 24 ing language resources by means of RDF tech-
 25 nology [173]. The idea here is to transform a lan-
 26 guage resource into an equivalent RDF represen-
 27 tation, to manipulate and enrich it with SPARQL
 28 transformation and external knowledge, and to se-
 29 rialize the result in RDF or non-RDF formats. To
 30 the extent that different formats can be mapped
 31 to or generated from the same RDF representa-
 32 tion, they can be transformed one into another.
 33 For lexical data, the OntoLex model and its afore-
 34 mentioned extensions represent a de facto stan-
 35 dard and are being used as such. For linguistic an-
 36 notations, several competing standards exist, and
 37 Prêt-à-LLOD contributes to on-going consolida-
 38 tion efforts within the W3C CG Linked Data for
 39 Language Technology with case studies on and
 40 support for CoNLL-RDF, NIF, Ligt, POWLA,
 41 and OLiA (see Sect. 4.2).

42 **metadata** Prêt-à-LLOD provides a workflow manage-
 43 ment system, a metadata repository for language
 44 resources, and machine-readable license informa-
 45 tion. In that regard, it also contributes to the devel-
 46 opment of metadata standards. This work is lead-
 47 ing to a new version of the Linghub site [113]²²⁴,
 48 that is based around the DSpace open source soft-

1 ware repository as well as the linking technolo-
 2 gies to provide a single authoritative source of in-
 3 formation about language resources across a wide
 4 range of languages.

5
 6 The key priority of Prêt-à-LLOD, however, is less to
 7 develop novel vocabularies, than to develop technical
 8 solutions on that basis. Accordingly, Prêt-à-LLOD in-
 9 volves four industry-led pilot projects that are designed
 10 to demonstrate the relevance, transferability and appli-
 11 cability of the methods and techniques under develop-
 12 ment in the project to concrete problems in the lan-
 13 guage technology industry. The pilots showcase po-
 14 tentials in the context of various sectors: technology
 15 companies, open government services, pharmaceutical
 16 industry, and finance, details of which are described
 17 in [174]. As overarching challenges, all pilots are ad-
 18 dressing facets of *cross-language transfer* or *domain*
 19 *adaptation* to varying degrees. Particularly relevant to
 20 LLOD, the project is developing tools that are help-
 21 ful to practical lexicographic applications including for
 22 the Oxford Dictionaries [175].

23 Notable project results in the context of this pa-
 24 per are a **Report on Vocabularies for Interoperable**
 25 **Language Resources and Services** that gives a brief
 26 overview over standards for language resources as of
 27 2019²²⁵ and the publication of the first monograph on
 28 LLOD technologies [124]. Whereas the latter builds
 29 on long-standing collaborations between its authors in
 30 previous projects and community groups, it was final-
 31 ized with support from the Prêt-à-LLOD project.

32 5.2.6. *NexusLinguarum* (2019-2023)

33 The **European network for Web-centred linguis-
 34 tic data science (NexusLinguarum)**²²⁶ is a COST
 35 Action project that involves researchers from 42 coun-
 36 tries. The network started in October 2019 and will
 37 continue its activities for four years. The COST Ac-
 38 tion promotes synergies across Europe between lin-
 39 guists, computer scientists, terminology experts, lan-
 40 guage professionals, and other stakeholders from both
 41 industry and society, in order to investigate into and
 42 to extend the areas of applicability of **linguistic data**
 43 **science** in a Web-centred context. Linguistic data sci-

44
 45
 46
 47
 48
 49
 50
 51
 223 <https://github.com/insight-centre/naisc>

224 <https://linghub.org>

225 Christian Chiarcos, Philipp Cimiano, Julia Bosque-Gil, Thierry Declerck, Christian Fäth, Jorge Gracia, Maxim Ionov, John McCrae, Elena Montiel-Ponsoda, Maria Pia di Buono, Roser Saurí, Fernando Bobillo, Mohammad Fazleh Elahi (2020), Report on Vocabularies for Interoperable Language Resources and Services, available from <https://cordis.europa.eu/project/id/825182/results>

226 <https://nexuslinguarum.eu/>

ence is concerned with providing a formal basis for the analysis, representation, integration and exploitation of linguistic data for language analysis (e.g. syntax, morphology, terminology, etc.) and language applications (e.g. machine translation, speech recognition, sentiment analysis, etc.). NexusLinguarum seeks to identify several key technologies to support such a study, including language resources, data analysis, NLP, and LLD. The latter is considered to be a cornerstone for the building of an ecosystem of multilingual and semantically interoperable linguistic data technologies and resources at a Web scale. Such an ecosystem is needed to foster the systematic cross-lingual discovery, exploitation, extension, curation and quality control of linguistic data.

One of the main research coordination objectives of NexusLinguarum is to propose, agree upon and disseminate best practices and standards for linking data and services across languages. In that regard, an active collaboration has been established with W3C community groups for the extension of existing standards such as OntoLex-Lemon as well as for the convergence of standards in language annotation (see Section 4). Several surveys of the state of the art are also being drafted by the NexusLinguarum community covering different salient aspects of the domain (e.g., multilingual linking across different linguistic description levels). A number of activities organised by NexusLinguarum have been planned with the aim of fostering collaboration and communication across communities. These include scientific conferences (e.g., LDK 2021²²⁷), and training schools (e.g., EuroLAN 2021²²⁸), where linguistic linked data will take a central role. Finally, NexusLinguarum is also devoted to the collection and analysis of relevant use cases for linguistic data science and to developing prototypes and demonstrators that will address a selection of prototypical cases. In an initial phase, the definition of use cases will cover Humanities and Social Sciences, Linguistics (Media and Social Media, and Language Acquisition), Life Sciences, and Technology (Cybersecurity and FinTech). NexusLinguarum also places a strong emphasis on lesser resourced languages.

A NexusLinguarum Use Case: ReTeRom

As an example of the kinds of complex, heterogeneous resources which have been proposed by consortium members as candidates for publication as linked

data with the support of members of the COST action, we will look at the corpora being produced in a Romanian language project.

The ReTeRom (*Resources and Technologies for Developing Human-Machine Interfaces in Romanian*) project²²⁹ is working towards adding the Romanian language to the multilingual Linguistic Linked Open Data cloud²³⁰. There are four different ReTeRom components. These are CoBiLiRo, SINTERO²³¹, TEPROLIN²³² and TADARAV²³³. We will focus on CoBiLiRo,

CoBiLiRo (*Bimodal Corpus for Romanian Language*), coordinated by the “Alexandru Ioan Cuza” University from Iai (UAIC), is working with a large collection of parallel speech/text data [179]. This collection is annotated on different levels on both the acoustic and the linguistic components [180], which

²²⁹https://www.racai.ro/p/reterom/index_en.html/

²³⁰Note that several Romanian language resources (e.g. Romanian WordNet (RoWN), Romanian Reference Treebank (RoRefTrees or RRT), Corpus-driven linguistic data, etc.) are currently in the process of conversion to LLD. The converter implementation is open source (<https://github.com/racai-ai/RoLLOD/>)

²³¹SINTERO (Technologies for the Realization of Human-Machine Interfaces for Text-to-Speech Synthesis with Expressivity), coordinated by Technical University of Cluj-Napoca (UTCN), primarily aims to implement a text-speech synthesis system in Romanian that allows the modelling and control of prosody (intonation in speech) in an appropriate way of natural speech. Secondly, SINTERO aims to create as many voices synthesised in Romanian as possible (in this project at least 10 voices), so that they too can be used by an extended community, including in commercial applications [176]

²³²TEPROLIN (*Technologies for Processing Natural Language - Text*) which is coordinated by the Research Institute for Artificial Intelligence Mircea Drăgulescu (ICIA), aims to create Romanian text processing technologies that can be readily used by the other component-projects of ReTeRom. For instance, higher layers of annotation may be performed using TEPROLIN services: on the speech component - the prosodic annotation (e.g. decrease of the fundamental frequency) and on the textual component - sub-syntactic (e.g. clauses) and syntactic annotation (e.g. parsing trees). TEPROLIN works inside a major language processing and text mining platform such as UIMA, GATE or TextFlows [177]

²³³TADARAV (*Technologies for automatic annotation of audio data and for the creation of automatic speech recognition interfaces*), coordinated by the University Politehnica of Bucharest (UPB), primarily aims to develop a set of advanced technologies for generating transcripts aligned correctly with the voice signal from the body collected in the CoBiLiRo component project. Secondly, TADARAV aims to increase the accuracy of the current Speed automatic speech recognition system [178] by requalifying its acoustic model based on the entire body of speech collected and using more powerful language models generated in the TEPROLIN component project.

²²⁷<http://2021.ldk-conf.org/>

²²⁸<http://eurolan.info.uaic.ro/2021>

facilitates searching, editing and statistical analysis operations over it. Three types of formats pairing speech and text components were identified in the building of the CoBiLiRo repository: (1) PHS/LAB, a format which separates text, speech and alignment in different files; (2) MULTTEXT/TEI, a format described initially in the MULTTEXT project and later used by various language resource builders; (3) TEXTGRID, a format supported by a large community of European developers and used in a large set of existing resources. In order to share and distribute these bimodal resources, a standard format for CoBiLiRo has been proposed, inspired by the TEI-P5.10 standard [181] and based on the idea of alignment between speech and the text components, taking into consideration several annotation conventions proposed in 2007 by Li and Zhi-gang [182]. At present, the header of this format includes the following metadata: *source of the object stored*; *speakers gender*; *speakers identity (if she/he agreed to this)*; *vocal type (spontaneous or in-reading)*; *recording conditions*; *duration*; *speech file type*; *speech-text alignment level*, etc. Moreover, the CoBiLiRo format allows for three types of segmentation ("file - adequate for resources held in multiple files, "startstop - adequate for resources that include only one speech file, and "file-start-stop a combination of the two types described before) and speech-text alignment, marked using <unit> tags. A <unit> tag includes two child nodes: the <speech> that names the file containing the speech component and the <text> that points to the corresponding textual transcription file.

As the preceding example (one of many within the project, falling within several different disciplines or technical domain) demonstrates NexusLinguaram's potential as a testing ground for many of the new vocabularies and modules mentioned above (as well as for the potential of linked data as a paradigm for the modelling and publication of language data) through the analysis of complex multifaceted use cases involving several different types of language resources.

6. Conclusions and Discussions of Future Challenges

We have attempted in the present article to give a comprehensive survey and a near-exhaustive²³⁴ de-

scription of the current state of affairs with respect to the use, definition and availability of models and vocabularies for Linguistic Linked (Open) Data. We have also gone into some detail as to the role of these models in various important initiatives, both past and present.

Moreover, as we hope that the article has demonstrated, this is an extremely active and dynamic area of research, with numerous projects and initiatives underway, or due to commence in the short term, which promise to bring further updates and improvements in expressivity and coverage in addition to the changes discussed here. For this reason, and in a vain attempt to stave off the dangers of rapid obsolescence, we have tried to situate our descriptions of recent advances in the field within a discussion of more general, ongoing trends. This was indeed our specific intention with Section 2 and in many other parts of the article. We hope that this survey will give the reader a good idea both of the future challenges which have yet to be fully confronted as well as the areas of immense opportunity which currently remain untapped.

In this rest of this section we will summarise the future prospects/challenges described in this paper. In the next and final subsection, Section 6.1, we focus on two particular areas and suggest a possible future trend and a proposal for a further direction of research.

In Section 3, we gave an overview of the most well known and widely used models for linguistic linked data, emphasising their FAIR-ness, and in particular their accessibility via ontology search engines, whether and how licensing information is made available, and how versioning is handled; we saw how, in many cases, there still remained work to be done in these areas. We classified these models into different subject areas based on the LLOD cloud, which helped us to describe how different areas of linguistics were covered by the models with some areas clearly better served than others. We also briefly discussed the provision of dedicated tools for LLD models; again this is an area which is still very much under development.

Next, in Section 4 we looked at the latest developments in LLD community standards. This section was divided into a subsection discussing OntoLex-Lemon related developments (Section 4.1), a section on the latest developments regarding LLD models for annotation (Section 4.2), and a section on metadata (Section 4.3). Each of these sections features a detailed description of different initiatives in their respective areas (including those still in progress), including in the case of Section 4.2 and Section 4.3 discussions of future trends and prospects (Section 4.2.5 and Section 4.3.3 respec-

²³⁴We were certainly exhausted after writing it.

tively). The main challenge in the case of LLD vocabularies for annotation is to respond to the need for a convergence of vocabularies. In the case of metadata vocabularies we looked at coverage issues, especially with regard to language identification.

Then in Section 5 we presented an overview of the impact of projects on the definition and use of LLD models and vocabularies. We focused on a number of ongoing projects and looked at their current and potential future contributions to LLD models and vocabularies. In the rest of this concluding section we will look at one important potential future trend, the involvement of research infrastructures alongside community groups and projects in the definition and ongoing development of models and vocabularies (Section 6.1.1). We will also make a proposal for handling the increasing complexity of LLD vocabularies (especially in the domain of language resources), namely, the recourse to ontology design patterns (Section 6.1.2).

6.1. Discussion of Future Trends and Challenges

6.1.1. Linguistic Linked Data, Projects, and Research Infrastructures

In this article we have tried to underline the role of research projects alongside that of community groups such as the Open Linguistics Working Group or the W3C Ontology-Lexicon Community Group in driving the development of LLD vocabularies and models. Going ahead however, the role of SSH research infrastructures could begin to play an important role too in order, that is, to ensure longer term hosting solutions and the greater sustainability of resources and tools based on these models.

Research infrastructures could also help to give long term support to the community groups which are developing such models and vocabularies that is in addition to and in a complimentary way to the support received from projects and COST actions in the short term. That is, in a similar way to which the TEI Lex-0 initiative (described in Section 5.2.3) has been supported both by a number of funded projects and COST actions as well as by the DARIAH "Lexical Resources" Working Group²³⁵.

At the same time research infrastructures could also help in the dissemination of LLD vocabularies and models, making them more accessible to ever wider numbers of users and cutting across different dis-

ciplinary boundaries via training and teaching initiatives. In other words, the Linguistic Linked Data community could exploit both the technical and the knowledge infrastructures provided by such **European Resource Infrastructure Consortia** (ERICs) as CLARIN, DARIAH in order to further sustain the work carried out in individual research projects and via open community groups.

A recent CLARIN event brought members of these two communities together in order to initiate a dialogue on future collaboration between the two²³⁶. The event was well received, being a promising start for future collaborations.

Note that although we will not discuss it here (as it would shift us too far into the realms of research policy), the role played by the European Open Science Cloud²³⁷ will also be crucial here (at the very least for projects and initiatives taking place in Europe) and especially for its consistent promotion of FAIR data.

6.1.2. A Proposal? The Use of Design Patterns

The OntoLex-Lemon model has come to be used in (or has at least been proposed for) a wide range of use cases pertaining to a range of different disciplines and kind of resources. As we have seen, the original model been extended to cover new kinds of use cases (or is planned to do so) by the W3C OntoLex group through the definition and publication of new extensions with their own separate guidelines. In the long term, however, this has the potential to become very complicated very quickly.

Take for instance the creation of specialised language resources for such areas as morpho-syntax or historical linguistics (in the former case these are dealt with in part in the original guidelines and in the new morphology module). In both these cases, there are so many different types (and sub-types) of resource as well as different theoretical approaches and schools of thought (not to mention language-specific modelling requirements) that it would be difficult to produce guidelines with detailed enough provision for any and all of the exigencies that might potentially arise. Or instead take the modelling of lexicographic resources (something which falls within the compass of the lexicog extension, Section 4.1.1). This could encompass numerous different kinds of sub-cases – e.g., etymo-

²³⁵See <https://dariah-eric.github.io/lexicalresources/pages/TEILex0/TEILex0.html>.

²³⁶A textual summary of the virtual event and recordings of the presentations and discussion can be found here <https://www.clarin.eu/event/2021/clarin-cafe-linguistic-linked-data>

²³⁷<https://eos-portal.eu/>

logical dictionaries, philological dictionaries, rhyming dictionaries – each of which brings its own specific varieties of modelling challenges. Furthermore, there can often also differing technical solutions to given modelling problems without a strong enough consensus on any single one of these. Such is the case with modelling ordered sequences in RDF.

One way of handling this potential modelling complexity and that avoids the drafting of ever more elaborate guidelines and the definition of ever more specialised modules is by through the publication and maintenance of a repository of **ontology design patterns** (ODP). ODP's are modelling solutions for recurring problems in the field of ontology design and are intended as a means of enhancing resuability in knowledge base design. In this they are based on previous work on patterns in software engineering. ODPs are arranged in six types [183]. These range from so called **Logical ODPs**, i.e., patterns that deal with problems in expressivity of formal knowledge engineering languages such as OWL (such as the representation of n-ary relations), and **Architectural ODPs** which are compositions of Logical ODPs, to **Reasoning ODPs** which propose procedures for automatic inference (for a full list see [183]). In our case the most relevant of these types are the **Content ODPs**, which are described as solving domain specific problems.

The idea would be to define, promote, and collect OntoLex-Lemon specific design patterns (as well as those pertaining to other vocabularies) within the LLD community and beyond. This is not a completely new idea and design patterns had been proposed for OntoLex-Lemon's predecessor *lemon* in the past. These are currently available on github²³⁸ and offer templates for the creation of nominal, verbal and adjectival lexical entries as well as more specific kinds of the former categories such as Relational Nouns, State Verbs and Intersective Adjectives. These patterns are fairly limited in scope and so our more specific proposal would be for patterns covering a variety of different areas/kinds of use cases. These patterns would initially correspond to the different sections of the W3C Ontolex guidelines, such as for example the syntax and semantics and the decomposition sections, as well as the lexicography module and the forthcoming Morphology and Frequency Attestation and Corpus (FrAC) modules. Each of these patterns would follow the set of criteria proposed in for instance [183] for Content

ODPs and would be based on competency questions, e.g., potential SPARQL queries.

Any resulting OntoLex-Lemon ODPs could then either be hosted on the ontology design patterns site²³⁹, or a special repository, or both. ODPs would provide a bridge between OntoLex guidelines and concrete applications; they would help to prevent those guidelines from becoming overly-complicated and unwieldy and would keep the extensions themselves as simple (and hopefully uncontroversial) as possible²⁴⁰. They would make models such as OntoLex-Lemon, and indeed several of the other models featured in this article, more accessible. Furthermore, they would allow us to recommend the re-use of other vocabularies without having to include them 'officially' within the OntoLex-Lemon guidelines themselves, ensuring the decoupling of the OntoLex-Lemon guidelines from these other vocabularies.

Acknowledgments

This article is based upon work from COST Action NexusLinguarum (CA18209), supported by COST (European Cooperation in Science and Technology). The authors thank Milan Dojchinovski and Francesca Frontini for several very helpful suggestions.

References

- [1] M.D. Wilkinson, M. Dumontier, I.J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L.B. da Silva Santos, P.E. Bourne, J. Bouwman, A.J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C.T. Evelo, R. Finkers, A. Gonzalez-Beltran, A.J.G. Gray, P. Groth, C. Goble, J.S. Grethe, J. Heringa, P.A.C. 't Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S.J. Lusher, M.E. Martone, A. Mons, A.L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S.-A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M.A. Swertz, M. Thompson, J. van der Lei, E. van Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao and B. Mons, The FAIR Guiding Principles for scientific data management and stewardship, *Scientific Data* 3(1) (2016), 160018. doi:10.1038/sdata.2016.18. <https://www.nature.com/articles/sdata201618>.
- [2] TEI Consortium, TEI P5: Guidelines for Electronic Text Encoding and Interchange, Zenodo, 2020. doi:10.5281/zenodo.3992514.

²³⁹http://ontologydesignpatterns.org/wiki/Main_Page

²⁴⁰Although of course the original modules would still need to be revised and extended on the basis of new kinds of use-cases/modelling needs; ODPs would help to keep these to a minimal.

²³⁸<https://github.com/jmccrae/lemon.patterns>

- [3] G. Francopoulo, M. George, N. Calzolari, M. Monachini, N. Bel, M. Pet and C. Soria, Lexical markup framework (LMF), 2006.
- [4] E. de la Clergerie and L. Clément, MAF: a morphosyntactic annotation framework, *Actes de LTC* (2005), 90–94.
- [5] N. Guarino and C.A. Welty, An Overview of OntoClean, in: *Handbook on Ontologies*, Springer Berlin Heidelberg, 2004, pp. 151–171.
- [6] B. Mons, FAIR Science for Social Machines: Let’s Share Metadata Knowlets in the Internet of FAIR Data and Services, *Data Intelligence* 1(1) (2019), 22–42.
- [7] J. Bosque-Gil, J. Gracia, E. Montiel-Ponsoda and A. Gómez-Pérez, Models to represent linguistic linked data, *Natural Language Engineering* 24(6) (2018), 811–859. doi:10.1017/S1351324918000347. <https://www.cambridge.org/core/journals/natural-language-engineering/article/models-to-represent-linguistic-linked-data/805F3E46882414B9144E43E34E89457D>.
- [8] H. Bohbot, F. Frontini, F. Khan, M. Khemakhem and L. Romary, Nénufar: Modelling a Diachronic Collection of Dictionary Editions as a Computational Lexical Resource, in: *The sixth biennial conference on electronic lexicography, eLex 2019*, 2019.
- [9] M. Passarotti, F. Mambrini, G. Franzini, F.M. Cecchini, E. Litta, G. Moretti, P. Ruffolo and R. Sprugnoli, Interlinking through Lemmas. The Lexical Collection of the LiLa Knowledge Base of Linguistic Resources for Latin, *Studi e Saggi Linguistici* 58(1) (2020), 177–212.
- [10] P. Cimiano, C. Chiarcos, J.P. McCrae and J. Gracia, *Linguistic Linked Data: Representation, Generation and Applications*, Springer International Publishing, Cham, 2020.
- [11] A. Pareja-Lora, M. Blume, B.C. Lust and C. Chiarcos (eds), *Development of linguistic linked open data resources for collaborative data-intensive research in the language sciences*, MIT Press, Cambridge, 2019. ISBN 978-0-262-53625-7.
- [12] C. Chiarcos, S. Hellmann and S. Nordhoff, Linking linguistic resources: Examples from the open linguistics working group, in: *Linked Data in Linguistics*, Springer, 2012, pp. 201–216.
- [13] Y. Le Franc, J. Parland-von Essen, L. Bonino, H. Lehvälaiho, G. Coen and C. Staiger, D2.2 FAIR Semantics: First recommendations (2020). doi:10.5281/ZENODO.3707985. <https://zenodo.org/record/3707985>.
- [14] P.-Y. Vandenbussche and B. Vatant, Metadata recommendations for linked open data vocabularies, *Version 1* (2011), 2011–12.
- [15] P.-Y. Vandenbussche, G.A. Atemezing, M. Poveda-Villalón and B. Vatant, Linked Open Vocabularies (LOV): a gateway to reusable semantic vocabularies on the Web, *Semantic Web* 8(3) (2017), 437–452.
- [16] J. McCrae, G. Aguado-de-Cea, P. Buitelaar, P. Cimiano, T. Declerck, A. Gómez-Pérez, J. Gracia, L. Hollink, E. Montiel-Ponsoda and D. Spohr, Interchanging lexical resources on the semantic web, *Language Resources and Evaluation* 46(4) (2012), 701–719, Publisher: Springer.
- [17] P. Cimiano, J.P. McCrae and P. Buitelaar, Lexicon Model for Ontologies: Community Report, W3C, 2016. <https://www.w3.org/2016/05/ontolex/>.
- [18] G. Sérasset, DBnary: Wiktionary as a Lemon-based multilingual lexical resource in RDF, *Semantic Web* 6(4) (2015), 355–361, Publisher: IOS Press.
- [19] M. Ehrmann, F. Cecconi, D. Vannella, J.P. McCrae, P. Cimiano and R. Navigli, Representing Multilingual Data as Linked Data: the Case of BabelNet 2.0, in: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, N.C.C. Chair, K. Choukri, T. Declerck, H. Loftsson, B. Maegaard, J. Mariani, A. Moreno, J. Odijk and S. Piperidis, eds, European Language Resources Association (ELRA), Reykjavik, Iceland, 2014. ISBN 978-2-9517408-8-4.
- [20] B. Klimek, N. Arndt, S. Krause and T. Arndt, Creating Linked Data morphological language resources with MMoOn (2016).
- [21] R. Forkel, The cross-linguistic linked data project, in: *3rd Workshop on Linked Data in Linguistics: Multilingual Knowledge Resources and Natural Language Processing*, 2014, p. 61.
- [22] M. Kemps-Snijders, M. Windhouwer, P. Wittenburg and S.E. Wright, ISOcat: Corraling data categories in the wild, in: *6th International Conference on Language Resources and Evaluation (LREC 2008)*, 2008.
- [23] S. Farrar and D.T. Langendoen, A linguistic ontology for the semantic web, *GLoT international* 7(3) (2003), 97–100.
- [24] H. Aristar-Dry, S. Drude, M. Windhouwer, J. Gippert and I. Nevskaya, Rendering endangered lexicons interoperable through standards harmonization: the relish project, in: *LREC 2012: 8th International Conference on Language Resources and Evaluation*, European Language Resources Association (ELRA), 2012, pp. 766–770.
- [25] D.T. Langendoen, Whither GOLD?, in: *Development of Linguistic Linked Open Data Resources for Collaborative Data-Intensive Research in the Language Sciences*, A. Pareja-Lora, B. Lust, M. Blume and C. Chiarcos, eds, MIT Press, 2019.
- [26] C. Chiarcos and M. Sukhareva, OLiA – ontologies of linguistic annotation, *Semantic Web* 6(4) (2015), 379–386.
- [27] C. Chiarcos, C. Fäth and F. Abromeit, Annotation Interoperability for the Post-ISOcat Era, in: *Proceedings of The 12th Language Resources and Evaluation Conference*, 2020, pp. 5668–5677.
- [28] C.A. Ferguson, Diglossia, *WORD* 15(2) (1959), 325–340. doi:10.1080/00437956.1959.11659702.
- [29] D.G. Martin Haspelmath Matthew S. Dryer and B. Comrie, *The world atlas of language structures*, Oxford University Press, 2005.
- [30] M.S. Dryer and M. Haspelmath (eds), *WALS Online*, Max Planck Institute for Evolutionary Anthropology, Leipzig, 2013. <https://wals.info/>.
- [31] P. Monachesi, A. Dimitriadis, R. Goedemans, A.-M. Mineur and M. Pinto, The typological database system, in: *Proceedings of the IRCS workshop on linguistic databases*, 2001, pp. 181–186.
- [32] A. Dimitriadis, M. Windhouwer, A. Saulwick, R. Goedemans and T. Bíró, How to integrate databases without starting a typology war: The Typological Database System, *The Use of Databases in Cross-Linguistic Studies*, Mouton de Gruyter, Berlin (2009), 155–207.
- [33] P. Westphal, C. Stadler and J. Pool, Countering language attrition with PanLex and the Web of Data, *Semantic Web* 6(4) (2015), 347–353.
- [34] C. Chiarcos, C. Fäth and M. Ionov, The ACoLi dictionary graph, in: *Proceedings of The 12th Language Resources and Evaluation Conference*, 2020, pp. 3281–3290.

- [35] G. De Melo, Lexvo.org: Language-related information for the linguistic linked data cloud, *Semantic Web* 6(4) (2015), 393–400.
- [36] A. Stellato, M. Fiorelli, A. Turbati, T. Lorenzetti, W. van Gemert, D. Dechandon, C. Laaboudi-Spoiden, A. Gerencsér, A. Waniart, E. Costetchi and et al., VocBench 3: A collaborative Semantic Web editor for ontologies, thesauri and lexicons, *Semantic Web* 11(5) (2020), 855881–. doi:10.3233/SW-200370.
- [37] A. Bellandi, E. Giovannetti and A. Weingart, Multilingual and multiword phenomena in a lemon old occitan medicobotanical lexicon, *Information* 9(3) (2018), 52.
- [38] C. Chiarcos, S. Hellmann and S. Nordhoff, Towards a Linguistic Linked Open Data cloud: The Open Linguistics Working Group, *TAL Traitement Automatique des Langues* 52(3) (2011), 245–275.
- [39] C. Chiarcos, S. Nordhoff and S. Hellmann, *Linked Data in Linguistics*, Springer, 2012.
- [40] J.P. McCrae, S. Moran, S. Hellmann and M. Brümmer (eds), *Semantic Web 6(4), Special Issue on Multilingual Linked Open Data*, IOS Press, 2015, pp. 313–400. <https://content.iiospress.com/journals/semantic-web/6/4>.
- [41] F. Khan, F. Boschetti and F. Frontini, Using lemon to model lexical semantic shift in diachronic lexical resources, in: *Proceedings of the 3rd Workshop on Linked Data in Linguistics (LDL-2014): Multilingual Knowledge Resources and Natural Language Processing*, 2014, pp. 50–54.
- [42] B. Klimek and M. Brümmer, Enhancing lexicography with semantic language databases, *Kernerman Dictionary News* 23 (2015), 5–10.
- [43] J. Bosque-Gil, J. Gracia, E. Montiel-Ponsoda and G. Aguado-de-Cea, Modelling multilingual lexicographic resources for the Web of Data: The K Dictionaries case, in: *GLOBALEX 2016 Lexicographic Resources for Human Language Technology Workshop Programme*, 2016, p. 65.
- [44] F. Abromeit, C. Chiarcos, C. Fäth and M. Ionov, Linking the Tower of Babel: Modelling a Massive Set of Etymological Dictionaries as RDF, in: *Proc. of the 5th Workshop on Linked Data in Linguistics: Managing, Building and Using Linked Language Resources*, 2016, p. 11.
- [45] J. Gracia, M. Villegas, A. Gómez-Pérez and N. Bel, The aperitium bilingual dictionaries on the web of data, *Semantic Web* 9(2) (2018), 231–240. doi:10.3233/SW-170258.
- [46] J. Bosque-Gil, J. Gracia and E. Montiel-Ponsoda, Towards a Module for Lexicography in OntoLex, in: *Proc. of the LDK workshops: OntoLex, TIAD and Challenges for Wordnets at 1st Language Data and Knowledge conference (LDK 2017), Galway, Ireland*, Vol. 1899, CEUR-WS, Galway (Ireland), 2017, pp. 74–84. ISSN 1613-0073. <http://ceur-ws.org/Vol-1899/OntoLex{ }2017{ }paper{ }5.pdf>.
- [47] J. Bosque-Gil, D. Lonke, J. Gracia and I. Kernerman, Validating the OntoLex-lemon lexicography module with K Dictionaries’ multilingual data, in: *Electronic lexicography in the 21st century. Proceedings of the eLex 2019 conference*, 2019, pp. 726–746.
- [48] B. Klimek, J.P. McCrae, M. Ionov, J.K. Tauber, C. Chiarcos, J. Bosque-Gil and P. Buitelaar, Challenges for the Representations for Morphology in Ontology Lexicons, in: *Proceedings of Sixth Biennial Conference on Electronic Lexicography, eLex 2019*, 2019. https://elex.link/elex2019/wp-content/uploads/2019/09/eLex_2019_33.pdf.
- [49] C. Chiarcos, M. Ionov, J. de Does, K. Depuydt, A.F. Khan, S. Stolk, T. Declerck and J.P. McCrae, Modelling Frequency and Attestations for OntoLex-Lemon, in: *Proceedings of the Globalex Workshop on Linked Lexicography (@LREC 2020)*, 2020, pp. 1–9. <https://lrec2020.lrec-conf.org/media/proceedings/Workshops/Books/GLOBALEX2020book.pdf#page=19>.
- [50] S. Peroni and D. Shotton, FaBiO and CiTO: ontologies for describing bibliographic resources and citations, *Web Semantics: Science, Services and Agents on the World Wide Web* 17 (2012), 33–43.
- [51] C. Chiarcos, T. Declerck and M. Ionov, Embeddings for the Lexicon: Modelling and Representation, in: *Proceedings of the 6th Workshop on Semantic Deep Learning (SemDeep-6), held virtually in January 2021, co-located with IJCAI-PRICAI 2020, Japan*, 2021.
- [52] S. Stolk, lemon-tree: Representing Topical Thesauri on the Semantic Web, in: *2nd Conference on Language, Data and Knowledge (LDK 2019)*, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.
- [53] S. Stolk, A Thesaurus of Old English as linguistic linked data: Using OntoLex, SKOS and lemon-tree to bring topical thesauri to the Semantic Web, in: *Proceedings of the eLex 2019 conference*, 2019, pp. 223–247.
- [54] A.F. Khan, Towards the Representation of Etymological Data on the Semantic Web, *Information* 9(12) (2018), 304.
- [55] F. Khan, Representing Temporal Information in Lexical Linked Data Resources, in: *Proceedings of the 7th Workshop on Linked Data in Linguistics (LDL-2020)*, European Language Resources Association, Marseille, France, 2020, pp. 15–22. ISBN 979-10-95546-36-8. <https://www.aclweb.org/anthology/2020.lidl-1.3>.
- [56] A. Burchardt, S. Padó, D. Spohr, A. Frank and U. Heid, Formalising Multi-layer Corpora in OWL/DL – Lexicon Modelling, Querying and Consistency Control, in: *Proc. of the 3rd International Joint Conference on NLP (IJCNLP)*, Hyderabad, India, 2008, pp. 389–396.
- [57] K. Verspoor and K. Livingston, Towards Adaptation of Linguistic Annotations to Scholarly Annotation Formalisms on the Semantic Web, in: *Proc. of the 6th Linguistic Annotation Workshop*, Association for Computational Linguistics, Jeju, Republic of Korea, 2012, pp. 75–84.
- [58] S. Hellmann, J. Lehmann, S. Auer and M. Brümmer, Integrating NLP using Linked Data, in: *Proc. 12th International Semantic Web Conference, 21-25 October 2013, Sydney, Australia*, 2013, also see <http://persistence.uni-leipzig.org/nlp2rdf/>.
- [59] N. Mazziotta, Building the syntactic reference corpus of medieval French using notabene rdf annotation tool, in: *Proc. of the Fourth Linguistic Annotation Workshop*, Association for Computational Linguistics, 2010, pp. 142–146.
- [60] B. Almas, H. Cayless, T. Clérice, Z. Fletcher, V. Jolivet, P. Liuzzo, E. Morlock, J. Robie, M. Romanello, J. Tauber and J. Witt, Distributed Text Services (DTS). First Public Working Draft, Technical Report, Github, 2019, version of May 23, 2019.
- [61] S. Cassidy, An RDF realisation of LAF in the DADA annotation server, in: *Proc. of the 5th Joint ISO-ACL/SIGSEM Workshop on Interoperable Semantic Annotation (ISA-5)*, Hong Kong, 2010.

- [62] N. Diewald, M. Hanl, E. Margaretha, J. Bingel, M. Kupietz, P. Bański and A. Witt, KorAP Architecture – Diving in the Deep Sea of Corpus Data, in: *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, 2016, pp. 3586–3591.
- [63] ISO, ISO 24612:2012. Language Resource Management - Linguistic Annotation Framework, Technical Report, ISO/TC 37/SC 4, Language resource management, 2012. <https://www.iso.org/standard/37326.html>.
- [64] C. Chiarcos, POWLA: Modeling Linguistic Corpora in OWL/DL, in: *The Semantic Web: Research and Applications*, E. Simperl, P. Cimiano, A. Polleres, O. Corcho and V. Presutti, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2012, pp. 225–239. ISBN 978-3-642-30284-8.
- [65] N. Ide and L. Romary, International Standard for a Linguistic Annotation Framework, *Natural language engineering* **10**(3–4) (2004), 211–225.
- [66] S. Tittel, H. Bermúdez-Sabel and C. Chiarcos, Using RDFa to Link Text and Dictionary Data for Medieval French, in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, J.P. McCrae, C. Chiarcos, T. Declerck, J. Gracia and B. Klimek, eds, European Language Resources Association (ELRA), Paris, France, 2018. ISBN 979-10-95546-19-1.
- [67] P. Ruiz Fabo, H. Bermúdez Sabel, C. Martínez Cantón and E. González-Blanco, The Diachronic Spanish Sonnet Corpus: TEI and linked open data encoding, data distribution, and metrical findings, *Digital Scholarship in the Humanities* (2020).
- [68] A. Gangemi, V. Presutti, D. Reforgiato Recupero, A.G. Nuzzolese, F. Draicchio and M. Mongiovi, Semantic Web Machine Reading with FRED, *Semantic Web* **8**(6) (2017), 873–893.
- [69] P. Vossen, R. Agerrí, I. Aldabe, A. Cybulska, M. van Erp, A. Fokkens, E. Laparra, A.-L. Minard, A.P. Aprosio, G. Rigau et al., Newsreader: Using knowledge resources in a cross-lingual reading machine to generate more knowledge from massive streams of news, *Knowledge-Based Systems* **110** (2016), 60–85.
- [70] N. Ide, J. Pustejovsky, C. Cieri, E. Nyberg, D. DiPersio, C. Shi, K. Suderman, M. Verhagen, D. Wang and J. Wright, The language application grid, in: *International Workshop on Worldwide Language Service Infrastructure*, Springer, 2015, pp. 51–70.
- [71] S. Peroni, A. Gangemi and F. Vitali, Dealing with markup semantics, in: *Proceedings of the 7th International Conference on Semantic Systems*, 2011, pp. 111–118.
- [72] A. Gangemi, N. Guarino, C. Masolo, A. Oltramari and L. Schneider, Sweetening ontologies with DOLCE, in: *International Conference on Knowledge Engineering and Knowledge Management*, Springer, 2002, pp. 166–181.
- [73] A. Fokkens, A. Soroa, Z. Beloki, N. Ockeloen, G. Rigau, W.R. Van Hage and P. Vossen, NAF and GAF: Linking linguistic annotations, in: *Proceedings 10th Joint ISO-ACL SIGSEM Workshop on Interoperable Semantic Annotation*, 2014, pp. 9–16.
- [74] M. Verhagen, K. Suderman, D. Wang, N. Ide, C. Shi, J. Wright and J. Pustejovsky, The LAPPS interchange format, in: *International Workshop on Worldwide Language Service Infrastructure*, Springer, 2015, pp. 33–47.
- [75] N. Ide, K. Suderman, J. Pustejovsky, M. Verhagen and C. Cieri, The language application grid and galaxy, in: *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, 2016, pp. 457–462.
- [76] E. Hinrichs, N. Ide, J. Pustejovsky, J. Hajic, M. Hinrichs, M.F. Elahi, K. Suderman, M. Verhagen, K. Rim, P. Stranák et al., Bridging the LAPPS Grid and CLARIN, in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*, 2018.
- [77] E. Wilde and M. Duerst, RFC 5147 – URI Fragment Identifiers for the text/plain Media Type, Technical Report, Internet Engineering Task Force (IETF), Network Working Group, 2008.
- [78] D. Filip, S. McCance, D. Lewis, C. Lieske, A. Lommel, J. Kosek, F. Sasaki and Y. Savourel, Internationalization Tag Set (ITS) Version 2.0, Technical Report, W3C Recommendation 29 October 2013, 2013.
- [79] J. Frey, M. Hofer, D. Obraczka, J. Lehmann and S. Hellmann, DBpedia FlexiFusion the best of Wikipedia> Wikidata> your data, in: *International Semantic Web Conference*, Springer, 2019, pp. 96–112.
- [80] R. Sanderson, P. Ciccarese and B. Young, Web Annotation Data Model, Technical Report, W3C Recommendation, 2017. <https://www.w3.org/TR/annotation-model/>.
- [81] R. Sanderson, P. Ciccarese and B. Young, Web Annotation Vocabulary, Technical Report, W3C Recommendation, 2017. <https://www.w3.org/TR/annotation-vocab/>.
- [82] L. Isaksen, R. Simon, E.T. Barker and P. de Soto Cañamares, Pelagios and the emerging graph of ancient world data, in: *Proceedings of the 2014 ACM conference on Web science*, 2014, pp. 197–201.
- [83] R. Simon, E. Barker, L. Isaksen and P. de Soto Cañamares, Linked Data Annotation Without the Pointy Brackets: Introducing Recogito 2, *Journal of Map & Geography Libraries* **13**(1) (2017), 111–132.
- [84] P. Cimiano, C. Chiarcos, J.P. McCrae and J. Gracia, *Linguistic Linked Data in Digital Humanities*, in: *Linguistic Linked Data*, Springer, 2020, pp. 229–262.
- [85] C. Chiarcos and M. Ionov, Ligt: An LLOD-Native Vocabulary for Representing Interlinear Glossed Text as RDF, in: *2nd Conference on Language, Data and Knowledge (LDK 2019)*, M. Eskevich, G. de Melo, C. Fäth, J.P. McCrae, P. Buitelaar, C. Chiarcos, B. Klimek and M. Dojchinovski, eds, OpenAccess Series in Informatics (OASISs), Vol. 70, Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany, 2019, pp. 3:1–3:15. ISSN 2190-6807. ISBN 978-3-95977-105-4. doi:10.4230/OASISs.LDK.2019.3. <http://drops.dagstuhl.de/opus/volltexte/2019/10367>.
- [86] S. Robinson, G. Aumann and S. Bird, Managing fieldwork data with Toolbox and the Natural Language Toolkit, *Language Documentation & Conservation* **1**(1) (2007), 44–57.
- [87] L. Butler and H. Van Volkinburg, Fieldworks Language Explorer (FLEx), *Technology Review* **1**(1) (2007), 1.
- [88] M.W. Goodman, J. Crowgey, F. Xia and E.M. Bender, Xigt: extensible interlinear glossed text for natural language processing, *Language Resources and Evaluation* **49**(2) (2015), 455–485.
- [89] S. Nordhoff, Modelling and Annotating Interlinear Glossed Text from 280 Different Endangered Languages as Linked Data with LIGT, in: *Proceedings of the 14th Linguistic Annotation Workshop*, 2020, pp. 93–104.

- [90] C. Chiarcos and C. Fäth, CoNLL-RDF: Linked Corpora Done in an NLP-Friendly Way, in: *Language, Data, and Knowledge*, J. Gracia, F. Bond, J.P. McCrae, P. Buitelaar, C. Chiarcos and S. Hellmann, eds, Springer, Cham, Switzerland, 2017, pp. 74–88. ISBN 978-3-319-59888-8.
- [91] M. Marcus, B. Santorini and M.A. Marcinkiewicz, Building a Large Annotated Corpus of English: The Penn Treebank, *Computational Linguistics* 19(2) (1993), 313–330.
- [92] F. Mambri and M. Passarotti, Linked Open Treebanks. Interlinking Syntactically Annotated Corpora in the LiLa Knowledge Base of Linguistic Resources for Latin, in: *Proceedings of the 18th International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2019)*, 2019, pp. 74–81.
- [93] M. Tamper, P. Leskinen, K. Apajalahti and E. Hyvönen, Using biographical texts as linked data for prosopographical research and applications, in: *Euro-Mediterranean Conference*, Springer, 2018, pp. 125–137.
- [94] C. Chiarcos, B. Kosmehl, C. Fäth and M. Sukhareva, Analyzing Middle High German Syntax with RDF and SPARQL, in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*, Miyazaki, Japan, 2018, pp. 4525–4534.
- [95] C. Chiarcos, I. Khaït, É. Pagé-Perron, N. Schenk, C. Fäth, J. Steuer, W. Mcgrath, J. Wang et al., Annotating a low-resource language with LLOD technology: Sumerian morphology and syntax, *Information* 9(11) (2018), 290.
- [96] C. Chiarcos and C. Fäth, Graph-based annotation engineering: towards a gold corpus for Role and Reference Grammar, in: *2nd Conference on Language, Data and Knowledge (LDK 2019)*, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.
- [97] M. Ionov, F. Stein, S. Sehgal and C. Chiarcos, cqp4rdf: Towards a Suite for RDF-Based Corpus Linguistics, in: *European Semantic Web Conference*, Springer, 2020, pp. 115–121.
- [98] C. Chiarcos, K. Donandt, H. Sargsian, M. Ionov and J.W. Schreur, Towards LLOD-based language contact studies. A case study in interoperability, in: *Proceedings of the 6th Workshop on Linked Data in Linguistics (LDL)*, 2018.
- [99] C. Chiarcos and L. Glaser, A Tree Extension for CoNLL-RDF, in: *Proceedings of the Twelfth International Conference on Language Resources and Evaluation (LREC-2020)*, ELRA, Marseille, France, 2020, pp. 7161–7169.
- [100] C. Chiarcos, POWLA: Modeling Linguistic Corpora in OWL/DL, in: *The Semantic Web: Research and Applications*, Vol. 7295, D. Hutchison, T. Kanade, J. Kittler, J.M. Kleinberg, F. Mattern, J.C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M.Y. Vardi, G. Weikum, E. Simperl, P. Cimiano, A. Polleres, O. Corcho and V. Presutti, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2012, pp. 225–239, Series Title: Lecture Notes in Computer Science. ISBN 978-3-642-30283-1 978-3-642-30284-8.
- [101] P. Cimiano, C. Chiarcos, J.P. McCrae and J. Gracia, *Modelling Linguistic Annotations*, in: *Linguistic Linked Data*, Springer, 2020, pp. 89–122.
- [102] D. Broeder, M. Kemps-Snijders, D. Van Uytvanck, M. Windhouwer, P. Withers, P. Wittenburg and C. Zinn, A Data Category Registry- and Component-based Metadata Framework, in: *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, European Language Resources Association (ELRA), Malta, 2010, pp. 43–47. http://www.lrec-conf.org/proceedings/lrec2010/pdf/163_Paper.pdf.
- [103] D. Broeder, D. van Uytvanck, M. Gavrilidou, T. Trippel and M. Windhouwer, Standardizing a Component Metadata Infrastructure, in: *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, N.C.C. Chair, K. Choukri, T. Declerck, M.U. Doğan, B. Maegaard, J. Mariani, A. Moreno, J. Odijk and S. Piperidis, eds, European Language Resources Association (ELRA), Istanbul, Turkey, 2012. ISBN 978-2-9517408-7-7.
- [104] M. Windhouwer, E. Indarto and D. Broeder, CMD2RDF: Building a Bridge from CLARIN to Linked Open Data, *Ubiquity Press* (2017), Publisher: Ubiquity Press. doi:10.5334/bbi.8.
- [105] D. Broeder, I. Schuurman and M. Windhouwer, Experiences with the ISOcat Data Category Registry, in: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, N.C.C. Chair, K. Choukri, T. Declerck, H. Loftsson, B. Maegaard, J. Mariani, A. Moreno, J. Odijk and S. Piperidis, eds, European Language Resources Association (ELRA), Reykjavik, Iceland, 2014. ISBN 978-2-9517408-8-4.
- [106] I. Schuurman, M. Windhouwer, O. Ohren and D. Zeman, CLARIN Concept Registry: The New Semantic Registry, in: *Selected Papers from the CLARIN Annual Conference 2015, October 14–16, 2015, Wrocław, Poland*, Linköping University Electronic Press, 2016, pp. 62–70.
- [107] P. Labropoulou, K. Gkirtzou, M. Gavrilidou, M. Deligiannis, D. Galanis, S. Piperidis, G. Rehm, M. Berger, V. Mapelli, M. Rigault, V. Arranz, K. Choukri, G. Backfried, J.M.G. Peñez and A. Garcia-Silva, Making Metadata Fit for Next Generation Language Technology Platforms: The Metadata Schema of the European Language Grid, in: *Proceedings of the 12th Language Resources and Evaluation Conference (LREC 2020)*, N. Calzolari, F. Bchet, P. Blache, C. Cieri, K. Choukri, T. Declerck, H. Isahara, B. Maegaard, J. Mariani, A. Moreno, J. Odijk and S. Piperidis, eds, European Language Resources Association (ELRA), Marseille, France, 2020, pp. 3421–3430.
- [108] M. Gavrilidou, P. Labropoulou, E. Desipri, S. Piperidis, H. Papageorgiou, M. Monachini, F. Frontini, T. Declerck, G. Francopoulo, V. Arranz and V. Mapelli, The META-SHARE Metadata Schema for the Description of Language Resources, in: *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC 2012)*, European Language Resources Association (ELRA), 2012. http://www.lrec-conf.org/proceedings/lrec2012/pdf/998_Paper.pdf.
- [109] J.P. McCrae, P. Labropoulou, J. Gracia, M. Villegas, V. Rodríguez-Doncel and P. Cimiano, One Ontology to Bind Them All: The META-SHARE OWL Ontology for the Interoperability of Linguistic Datasets on the Web, in: *The Semantic Web: ESWC 2015 Satellite Events*, F. Gandon, C. Guéret, S. Villata, J. Breslin, C. Faron-Zucker and A. Zimmermann, eds, Lecture Notes in Computer Science, Springer International Publishing, 2015, pp. 271–282. ISBN 978-3-319-25639-9. https://link.springer.com/chapter/10.1007/978-3-319-25639-9_42.

- [110] S. Piperidis, The META-SHARE Language Resources Sharing Infrastructure: Principles, Challenges, Solutions, in: *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, N.C.C. Chair, K. Choukri, T. Declerck, M.U. Doğan, B. Maegaard, J. Mariani, A. Moreno, J. Odijk and S. Piperidis, eds, European Language Resources Association (ELRA), Istanbul, Turkey, 2012. ISBN 978-2-9517408-7-7.
- [111] V. Rodriguez-Doncel and P. Labropoulou, Digital Representation of Licenses for Language Resources, in: *Proceedings of the 4th Workshop on Linked Data in Linguistics: Resources and Applications*, Association for Computational Linguistics, Beijing, China, 2015, pp. 49–58. doi:10.18653/v1/W15-4206. <http://aclweb.org/anthology/W15-4206>.
- [112] P. Labropoulou, D. Galanis, A. Lempesis, M. Greenwood, P. Knoth, R. Eckart de Castilho, S. Sachtouris, B. Georgantopoulos, S. Martziou, L. Anastasiou, K. Gkirtzou, N. Manola and S. Piperidis, OpenMinTeD: A Platform Facilitating Text Mining of Scholarly Content, in: *WOSP 2018 Workshop Proceedings, Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, European Language Resources Association (ELRA), Miyazaki, Japan, 2018, pp. 7–12. ISBN 979-10-95546-20-7. http://lrec-conf.org/workshops/lrec2018/W24/pdf/13_W24.pdf.
- [113] J.P. McCrae and P. Cimiano, Linghub: a Linked Data based portal supporting the discovery of language resources., *SEMANTiCS (Posters & Demos)* **1481** (2015), 88–91.
- [114] J.P. McCrae, P. Cimiano, V. Rodríguez Doncel, D. Vila-Suero, J. Gracia, L. Matteis, R. Navigli, A. Abele, G. Vulcu and P. Buitelaar, Reconciling Heterogeneous Descriptions of Language Resources, in: *Proceedings of the 4th Workshop on Linked Data in Linguistics: Resources and Applications*, Association for Computational Linguistics, Beijing, China, 2015, pp. 39–48. doi:10.18653/v1/W15-4205. <https://www.aclweb.org/anthology/W15-4205>.
- [115] G. Rehm, M. Berger, E. Elsholz, S. Hegele, F. Kintzel, K. Marheinecke, S. Piperidis, M. Deligiannis, D. Galanis, K. Gkirtzou, P. Labropoulou, K. Bontcheva, D. Jones, I. Roberts, J. Hajič, J. Hamrlová, L. Kačena, K. Choukri, V. Arranz, A. Vasiljevs, O. Anvari, A. Lagzdinš, J. Melņika, G. Backfried, E. Dikici, M. Janosik, K. Prinz, C. Prinz, S. Stampfer, D. Thomas-Aniola, J.M. Gómez-Pérez, A. Garcia Silva, C. Berrío, U. Germann, S. Renals and O. Klejch, European Language Grid: An Overview, in: *Proceedings of the 12th Language Resources and Evaluation Conference*, European Language Resources Association, Marseille, France, 2020, pp. 3366–3380. ISBN 979-10-95546-34-4. <https://www.aclweb.org/anthology/2020.lrec-1.413>.
- [116] M. Fiorelli, A. Stellato, J.P. McCrae, P. Cimiano and M.T. Paziienza, LIME: the metadata module for OntoLex, in: *European Semantic Web Conference*, Springer, 2015, pp. 321–336.
- [117] R. Cyganiak, D. Wood and M. Lanthaler, RDF 1.1 Concepts and Abstract Syntax, Technical Report, W3C Recommendation 25 February 2014, 2014.
- [118] G. De Melo, Lexvo. org: Language-related information for the linguistic linked data cloud, *Semantic Web* **6**(4) (2015), 393–400. Publisher: IOS Press.
- [119] A. Phillips and M. Davis, BCP 47 – Tags for Identifying Languages, Technical Report, Internet Engineering Task Force, 2006. <http://www.rfc-editor.org/rfc/bcp/bcp47.txt>.
- [120] F. Gillis-Webber and S. Tittel, The shortcomings of language tags for linked data when modeling lesser-known languages, in: *2nd Conference on Language, Data and Knowledge (LDK 2019)*, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.
- [121] S. Tittel and F. Gillis-Webber, Identification of Languages in Linked Data: A Diachronic-Diatopic Case Study of French, in: *Electronic lexicography in the 21st century. Proceedings of the eLex 2019 conference*, 2019, pp. 1–3.
- [122] F. Gillis-Webber and S. Tittel, A Framework for Shared Agreement of Language Tags beyond ISO 639, in: *Proceedings of The 12th Language Resources and Evaluation Conference*, 2020, pp. 3333–3339.
- [123] S. Nordhoff, Linked data for linguistic diversity research: Glottolog/langdoc and asjp online, in: *Linked Data in Linguistics*, Springer, 2012, pp. 191–200.
- [124] P. Cimiano, C. Chiarcos, J.P. McCrae and J. Gracia, *Linguistic Linked Data*, Springer, 2020.
- [125] A. Pareja-Lora, M. Blume, B.C. Lust and C. Chiarcos, *Development of Linguistic Linked Open Data Resources for Collaborative Data-Intensive Research in the Language Sciences*, MIT Press, 2020.
- [126] L. Romary, M. Khemakhem, M. George, J. Bowers, F. Khan, M. Pet, S. Lewis, N. Calzolari and P. Banski, LMF Reloaded., in: *Asialex 2019*, 2019.
- [127] V.R. Doncel and E.M. Ponsoda, LYNX: Towards a Legal Knowledge Graph for Multilingual Europe, *Law in Context. A Socio-legal Journal* **37**(1) (2020), 1–4.
- [128] T. Burrows, E. Hyvönen, L. Ransom and H. Wijsman, Mapping Manuscript Migrations: Digging into Data for the History and Provenance of Medieval and Renaissance Manuscripts, *Manuscript Studies: A Journal of the Schoenberg Institute for Manuscript Studies* **3**(1) (2018), 249–252.
- [129] E. Hyvönen, “Sampo” Model and Semantic Portals for Digital Humanities on the Semantic Web., in: *DHN*, 2020, pp. 373–378.
- [130] A. Khan, H. Bohbot, F. Frontini, M. Khemakhem and L. Romary, Historical Dictionaries as Digital Editions and Connected Graphs: the Example of Le Petit Larousse Illustré, in: *Digital Humanities 2019*, 2019.
- [131] A. Weingart and E. Giovannetti, A Lexicon for Old Occitan Medico-Botanical Terminology in Lemon., in: *SWASH@ESWC*, 2016, pp. 25–36.
- [132] R. Costa, A. Salgado, A.F. Khan, S. Carvalho, L. Romary, B. Almeida, M. Ramos, M. Khemakhem, R. Silva and T. Tasovac, MORDigital: The Advent of a New Lexicographical Portuguese Project, in: *eLex 2021 - Seventh biennial conference on electronic lexicography*, Brno, Czech Republic, 2021. <https://hal.inria.fr/hal-03195362>.
- [133] V. Propp, *Morphology of the folktale*, Trans., Laurence Scott. 2nd ed., University of Texas Press, 1968.
- [134] S. Thompson, *Motif-index of folk-literature: A classification of narrative elements in folktales, ballads, myths, fables, medieval romances, exempla, fabliaux, jest-books, and local legends*, Revised and enlarged edition (1955–1958), Indiana University Press, 1958.
- [135] H.-J. Uther, *The Types of International Folktales: A Classification and Bibliography. Based on the system of Antti Aarne and Stith Thompson*, Suomalainen Tiedekatemia, 2004.

- [136] T. Declerck, A. Kostová and L. Schäfer, Towards a Linked Data Access to Folktales classified by Thompsons Motifs and Aarne-Thompson-Uthers Types, in: *Proceedings of Digital Humanities 2017*, ADHO, 2017.
- [137] F. Diehr, M. Brodhun, S. Gronemeyer, K. Diederichs, C. Prager, E. Wagner and N. Grube, Modellierung eines digitalen Zeichenkatalogs für die Hieroglyphen des Klassischen Maya, in: *47. Jahrestagung der Gesellschaft für Informatik, Digitale Kulturen, INFORMATIK 2017, Chemnitz, Germany, September 25-29, 2017*, M. Eibl and M. Gaedke, eds, LNI, Vol. P-275, GI, 2017, pp. 1185–1196. doi:10.18420/in2017_120.
- [138] G. Zólyomi, B. Tanos and S. Sövegjártó, The Electronic Text Corpus of Sumerian Royal Inscriptions, 2008.
- [139] É. Pagé-Perron, M. Sukhareva, I. Khait and C. Chiarcos, Machine translation and automated analysis of the sumerian language, in: *Proceedings of the Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, 2017, pp. 10–16.
- [140] M. Hartung, M. Orlikowski and S. Veríssimo, Evaluating the Impact of Bilingual Lexical Resources on Cross-lingual Sentiment Projection in the Pharmaceutical Domain, 2020.
- [141] B.-P. Ivanschitz, T.J. Lampoltshammer, V. Mireles, A. Revenko, S. Schlarb and L. Thurnay, A Semantic Catalogue for the Data Market Austria., in: *SEMANTICS Posters&Demos*, 2018.
- [142] D. Lonke and J. Bosque Gil, Applying the OntoLex-lemon lexicography module to K Dictionaries' multilingual data, *K Lexical News (KLN)* (2019). <https://kln.lexicala.com/kln28/lonke-bosque-gil-ontolex-lemon-lexicog/>.
- [143] G. Rehm, D. Galanis, P. Labropoulou, S. Piperidis, M. Weiß, R. Usbeck, J. Köhler, M. Deligiannis, K. Gkirtzou, J. Fischer, C. Chiarcos, N. Feldhus, J. Moreno-Schneider, F. Kintzel, E. Montiel, V. Rodríguez Doncel, J.P. McCrae, D. Laqua, I.P. Theile, C. Dittmar, K. Bontcheva, I. Roberts, A. Vasiljevs and A. Lagzdīņš, Towards an Interoperable Ecosystem of AI and LT Platforms: A Roadmap for the Implementation of Different Levels of Interoperability, in: *Proceedings of the 1st International Workshop on Language Technology Platforms*, European Language Resources Association, Marseille, France, 2020, pp. 96–107. ISBN 979-10-95546-64-1. <https://www.aclweb.org/anthology/2020.iwltpl-1.15>.
- [144] Sator, Sator 2021 Data-for-AI Market Report, Technical Report, Sator, 2021.
- [145] C. Chiarcos, M. Ionov, M. Rind-Pawłowski, C. Fäth, J.W. Schreur and I. Nevskaya, LLODifying linguistic glosses, in: *Proceedings of Language, Data and Knowledge (LDK-2017)*, Galway, Ireland, 2017.
- [146] A. Pareja-Lora, B. Lust, M. Blume and C. Chiarcos, *Development of Linguistic Linked Open Data Resources for Collaborative Data-Intensive Research in the Language Sciences*, The MIT Press, 2019.
- [147] H.B.-S. Sabine Tittel and C. Chiarcos, Using RDFa to Link Text and Dictionary Data for Medieval French, in: *Proc. of the 6th Workshop on Linked Data in Linguistics (LDL-2018): Towards Linguistic Data Science*, European Language Resources Association (ELRA), Paris, France, 2018. ISBN 979-10-95546-19-1.
- [148] M. Curado Malta, P. Centenera and E. González-Blanco García, POSTDATA – Towards publishing European Poetry as Linked Open Data, *International Conference on Dublin Core and Metadata Applications* 16 (2016), 19–20.
- [149] M. Curado Malta, P. Centenera and E. Gonzalez-Blanco, Using Reverse Engineering to Define a Domain Model: The Case of the Development of a Metadata Application Profile for European Poetry, in: *Developing Metadata Application Profiles*, IGI Global, 2017, pp. 146–180. doi:10.4018/978-1-5225-2221-8. <http://e-spacio.uned.es/fez/view/bibliuned:365-Egonzalez9>.
- [150] M. Curado Malta, H. Bermúdez-Sabel, A.A. Baptista and E. Gonzalez-Blanco, Validation of a metadata application profile domain model, *International Conference on Dublin Core and Metadata Applications* (2018), 65–75.
- [151] M. Curado Malta, Modelação de dados poéticos: Uma perspectiva desde os dados abertos e ligados, in: *Humanidades Digitales. Miradas hacia la Edad Media*, D. González and H. Bermudez Sabel, eds, De Gruyter, Berlin, 2019, pp. 24–48. ISBN 978-3-11-058542-1. <https://doi.org/10.1515/9783110585421-004>.
- [152] E. González-Blanco, S. Ros Muñoz, M.L. Díez Platas, J. De la Rosa, H. Bermúdez-Sabel, A. Pérez Pozo, L. Ayrcirix and B. Sartini, Towards an Ontology for European Poetry, DARIAH Annual Event 2019, Warsaw, Poland, 2019. doi:10.5281/zenodo.3458772. https://zenodo.org/record/3458772#.Xhw_YOhKjIV.
- [153] P.E. project, Network of ontologies - POSTDATA, Postdata ERC project, [Online; accessed 2021-01-17]. <http://http://postdata.linhd.uned.es/results/>.
- [154] P.E. project, Postdata-core ontology, Postdata ERC project, [Online; accessed 2021-01-17]. <http://http://postdata.linhd.uned.es/results/>.
- [155] P.E. project, Postdata-prosodic ontology, Postdata ERC project, [Online; accessed 2021-01-17]. <http://http://postdata.linhd.uned.es/results/>.
- [156] P.E. project, Postdata-structural ontology, Postdata ERC project, [Online; 2021-01-17]. <http://http://postdata.linhd.uned.es/results/>.
- [157] M.L. Díez Platas, S. Ros Muñoz, E. González-Blanco, P. Ruiz Fabo and E. Álvarez Mellado, Medieval Spanish (12th-15th centuries) Named Entity Recognition and Attribute Annotation System based on contextual information, *JASIST (Journal of the Association for Information Science and Technology)* (2020). doi:<https://doi.org/10.1002/asi.24399>.
- [158] J. De la Rosa, S. Ros Muñoz, E. González-Blanco, Á. Pérez Pozo, L. Hernández and A. Díaz Medina, Bertsification: Language modeling fine-tuning for Spanish scansion, 4th International Conference on Science and Literature (postponed due to COVID-19 crisis), Girona, 2020.
- [159] J. De La Rosa, S. Ros Muñoz, E. González-Blanco and Á. Pérez Pozo, PoetryLab: An Open Source Toolkit for the Analysis of Spanish Poetry Corpora, Carleton University and the University of Ottawa, Virtual Conference, 2020, DH2020. doi:<http://dx.doi.org/10.17613/rsd8-we57>. <https://hcommons.org/deposits/item/hc:31763/>.
- [160] J. de la Rosa, S. Ros and E. González-Blanco, Predicting metrical patterns in Spanish poetry with language models, *arXiv preprint arXiv:2011.09567* (2020).
- [161] S. Krek, I. Kosem, J.P. McCrae, R. Navigli, B.S. Pedersen, C. Tiberius and T. Wissik, European lexicographic infrastructure (elexis), in: *Proceedings of the XVIII EURALEX Interna-*

- 1 *tional Congress on Lexicography in Global Contexts*, 2018,
2 pp. 881–892.
- 3 [162] P. Bański, J. Bowers and T. Erjavec, TEI-Lex0 guidelines for
4 the encoding of dictionary information on written and spoken
5 forms, in: *Electronic Lexicography in the 21st Century: Proceedings of ELex 2017 Conference*, 2017.
- 6 [163] F. Mambrini and M. Passarotti, Linked Open Treebanks.
7 Interlinking Syntactically Annotated Corpora in the LiLa
8 Knowledge Base of Linguistic Resources for Latin, in: *Proceedings of the 18th International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2019)*, Association for Computational Linguistics, Paris, France, 2019,
9 pp. 74–81. doi:10.18653/v1/W19-7808. <https://www.aclweb.org/anthology/W19-7808>.
- 10 [164] R. Sprugnoli, M. Passarotti, F.M. Cecchini and M. Pellegrini,
11 Overview of the Evalatin 2020 Evaluation Campaign, in: *Proceedings of LT4HALA 2020 - 1st Workshop on Language Technologies for Historical and Ancient Languages*, European Language Resources Association (ELRA), Marseille, France, 2020, pp. 105–110. ISBN 979-10-95546-53-5.
12 <https://www.aclweb.org/anthology/2020.lt4hala-1.16>.
- 13 [165] M. Passarotti, M. Budassi, E. Litta and P. Ruffolo, The Lemlat 3.0 Package for Morphological Analysis of Latin, in: *Proceedings of the NoDaLiDa 2017 Workshop on Processing Historical Language*, Linköping University Electronic Press, 2017, pp. 24–31.
- 14 [166] E. Litta, M. Passarotti and F. Mambrini, The Treatment of
15 Word Formation in the LiLa Knowledge Base of Linguistic Resources for Latin, in: *Proceedings of the Second International Workshop on Resources and Tools for Derivational Morphology*, Charles University, Faculty of Mathematics and Physics, Institute of Formal and Applied Linguistics, Prague, Czechia, 2019, pp. 35–43. <https://www.aclweb.org/anthology/W19-8505>.
- 16 [167] F. Mambrini and M. Passarotti, Harmonizing Different
17 Lemmatization Strategies for Building a Knowledge Base of Linguistic Resources for Latin, in: *Proceedings of the 13th Linguistic Annotation Workshop*, Association for Computational Linguistics, Florence, Italy, 2019,
18 pp. 71–80. doi:10.18653/v1/W19-4009. <https://www.aclweb.org/anthology/W19-4009>.
- 19 [168] F. Mambrini and M. Passarotti, Representing Etymology in
20 the LiLa Knowledge Base of Linguistic Resources for Latin, in: *Proceedings of the 2020 Globalex Workshop on Linked Lexicography*, European Language Resources Association, Marseille, France, 2020, pp. 20–28. ISBN 979-10-95546-46-7. <https://www.aclweb.org/anthology/2020.globalex-1.3>.
- 21 [169] G. Franzini, F. Zampedri, M. Passarotti, F. Mambrini and
22 G. Moretti, Græcissære: Ancient Greek Loanwords in the LiLa Knowledge Base of Linguistic Resources for Latin., in: *Seventh Italian Conference on Computational Linguistics*, J. Monti, F. Dell’Orletta and F. Tamburini, eds, CEUR-WS.org, Bologna, 2020, pp. 1–6. http://ceur-ws.org/Vol-2769/paper_06.pdf.
- 23 [170] A. Westerski and J.F. Sánchez-Rada, Marl Ontology Specification, V1.1 8 March 2016, 2016. <http://www.gsi.dit.upm.es/ontologies/marl/>.
- 24 [171] G. Franzini, A. Peverelli, P. Ruffolo, M. Passarotti, H. Sanna,
25 E. Signoroni, V. Ventura and F. Zampedri, Nunc Est Aestimandum. Towards an evaluation of the Latin WordNet, in: *Sixth Italian Conference on Computational Linguistics (CLiC-it 2019)*, R. Bernardi, R. Navigli and G. Semeraro, eds, CEUR-WS.org, Bari, Italy, 2019, pp. 1–8.
- 26 [172] F.M. Cecchini, R. Sprugnoli, G. Moretti and M. Passarotti,
27 UDante: First Steps Towards the Universal Dependencies Treebank of Dante’s Latin Works, in: *Seventh Italian Conference on Computational Linguistics*, CEUR-WS.org, 2020,
28 pp. 1–7.
- 29 [173] C. Fäth, C. Chiarcos, B. Ebbrecht and M. Ionov, Fintan –
30 Flexible, Integrated Transformation and Annotation eEngineering, in: *Proceedings of the Twelfth International Conference on Language Resources and Evaluation (LREC-2020)*, ELRA, Marseille, France, 2020, pp. 7212–7221.
- 31 [174] T. Declerck, J. McCrae, M. Hartung, J. Gracia, C. Chiarcos,
32 E. Montiel, P. Cimiano, A. Revenko, R. Sauri, D. Lee, S. Racioppa, J. Nasir, M. Orlikowski, M. Lanau-Coronas, C. Fäth, M. Rico, M.F. Elahi, M. Khvalchik, M. Gonzalez and K. Cooney, Recent Developments for the Linguistic Linked Open Data Infrastructure, in: *Proceedings of the Twelfth International Conference on Language Resources and Evaluation (LREC 2020)*, N. Calzolari, F. Béchet, P. Blache, C. Cieri, K. Choukri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk and S. Piperidis, eds, ELRA, 2020, pp. 5660–5667, ELRA.
- 33 [175] R. Sauri, L. Mahon, I. Russo and M. Bitinis, Cross-Dictionary
34 Linking at Sense Level with a Double-Layer Classifier, in: *2nd Conference on Language, Data and Knowledge (LDK 2019)*, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019.
- 35 [176] B. Lorincz, M. Nutu, A. Stan and G. Mircea, An Evaluation
36 of Postfiltering for Deep Learning Based Speech Synthesis with Limited Data, in: *IEEE 10th International Conference on Intelligent Systems (IS)*, 2020.
- 37 [177] R. Ion, Teprolin: an Extensible, Online Text Preprocessing
38 Platform for Romanian, in: *Proceedings of the ConsILR-2018*, 2018, pp. 69–76.
- 39 [178] A.L. Georgescu, H. Cucu, A. Buzo and C. Burileanu, RSC: A
40 Romanian Read Speech Corpus for Automatic Speech Recognition, in: *Proceedings of the 12th Language Resources and Evaluation Conference (LREC)*, 2020, pp. 6606–6612.
- 41 [179] D. Cristea, I. Pistol, S. Boghiu, A. Bibiri, D. Gifu, A. Scutelnicu,
42 M. Onofrei, D. Trandabat and G. Bugeag, CoBiLiRo: a Research Platform for Bimodal Corpora, in: *Proceedings of the 1st International Workshop on Language Technology Platforms (IWLTP 2020)*, European Language Resources Association, 2020, pp. 22–27.
- 43 [180] D. Gifu, A. Moruz, C. Bolea, A. Bibiri and M. Mitrofan, The
44 Methodology of Building CoRoLa, in: *Revue Roumaine de Linguistique (Romanian Review of Linguistics)/ On design, creation and use of of the Reference Corpus of Contemporary Romanian and its analysis tools. CoRoLa, KorAP, DRuKoLa and EuReCo / Conception, création et utilisation du Corpus de Référence du Roumain Contemporain et de ses outils d’analyse. CoRoLa, KorAP, DRuKoLa et EuReCo*, Vol. 64, 2019, pp. 241–253.
- 45 [181] C.M. Sperberg-McQueen and L. Burnard, Original editors,
46 revised and expanded under the supervision of the Technical Council of the TEI Consortium. TEI P5: Guidelines for Electronic Text Encoding and Interchange, 2018.
- 47 [182] A. Li and Z. Yin, Standardization of Speech Corpus, in: *Data Science Journal*, Vol. 6, 2007.

[183] V. Presutti, E. Blomqvist, E. Daga and A. Gangemi, Pattern-based ontology design, in: *Ontology Engineering in a Networked World*, Springer, 2012, pp. 35–64.

Appendix A. The OntoLex-Lemon Model

In order to make this paper as self-contained as possible we will introduce the OntoLex-Lemon model in the appendix. We will describe the core module with an example and briefly describe its different sub-modules. The full guidelines can be found at the following URL: <https://www.w3.org/2016/05/ontolex/>. This appendix is stand-alone and can be skipped by those who are already familiar with the OntoLex-Lemon model.

A.1. Introduction and the Core Module

As has been mentioned several times in this article, OntoLex-Lemon, is an RDF-native model for the modelling and publication of ontology-lexicons (that is language resources that consist both of a lexical and an ontological component) on the Semantic Web. OntoLex-Lemon is an update of the original *lemon* model and as with the latter model it has the aim of enriching or grounding linked data ontologies with linguistic information. However it has increasingly come to be used for the representation of linked data lexicons without any ontological component (and indeed extensions of OntoLex-Lemon such as lexicog 4.1.1 and FrAC 4.1.3 are strongly motivated by lexical use-cases rather than ontological ones).

OntoLex consists of a *core module* and four additional thematic modules three of which we will describe in the following sections (the Metadata module is described in Section 4.3.3). In the current section we will look at the core module, presented in Figure 8, with the help of an example. As the figure shows, the core of OntoLex-Lemon is essentially based around the class Lexical Entry and the different properties and relationships that it can have. Lexical Entry has the subclasses Word, Multiword Expression and Affix. We can associate information about the different forms that a Lexical Entry can have via the Form class. The Form class can be associated with written representations via the representation properties and its sub-properties. The semantics of a Lexical Entry can be described using the Lexical Sense and the Lexical Concept classes (a fuller ontological based description of the semantics of an entry can be found in the Syntax and Semantics

module which is described in Section A.2 . ontoLexical Entry is related to Lexical Sense via the sense property (and its inverse isSenseOf. Members of the Lexical Sense class represent

the lexical meaning of a lexical entry when *interpreted as referring to the corresponding ontology element...*A link between a lexical entry and an ontology entity via a Lexical Sense object *implies that the lexical entry can be used to refer to the ontology entity in question*²⁴¹. [Emphasis ours]

The Lexical Sense object relates to a corresponding ontology entity via the reference property; Lexical Entry individuals can be related directly to these ontology entities via the property denotes²⁴². A second class which is used to describe the semantics of an entry is the Lexical Concept class. Members of the latter class represent "a mental abstraction, concept or unit of thought that can be lexicalized by a given collection of senses"²⁴³; Lexical Concept is a subclass of the SKOS class Concept.

We will show how some of these classes and properties are used in practise with an example entry from Wiktionary, the open-source online multilingual dictionary. We will take the example of the Urdu word زبان (*zuban*) which means both ‘tongue’ and ‘language’²⁴⁴; a screenshot of the entry can be seen in Figure 9.

The entry contains some standard morphosyntactic and semantic information about the word, quite a lot of which can already be encoded using the OntoLex-Lemon core, along with a select number of properties from the LexInfo vocabulary (described above in Section 3.5).

Figure 10 shows the encoding of زبان as a OntoLex-Lemon Word. Here we can note the use of the LexInfo vocabulary to specify the part of speech and gender of the word and the use of the lime vocabulary to specify the language of the entry. We can also see the relation of the entry to the different forms of the word via the two object properties canonicalForm (linking it up to its headword form) and otherForm (linking the word up to its different morphological variants). In addition

²⁴¹Definition taken from <https://www.w3.org/2016/05/ontolex/#lexical-sense-reference>.

²⁴²The property denotes is equivalent to the property chain sense o reference.

²⁴³Definition taken from <https://www.w3.org/2016/05/ontolex/#lexical-concept>.

²⁴⁴The entry can be found here <https://en.wiktionary.org/wiki/%D8%B2%D8%A8%D8%A7%D9%86#Urdu>.

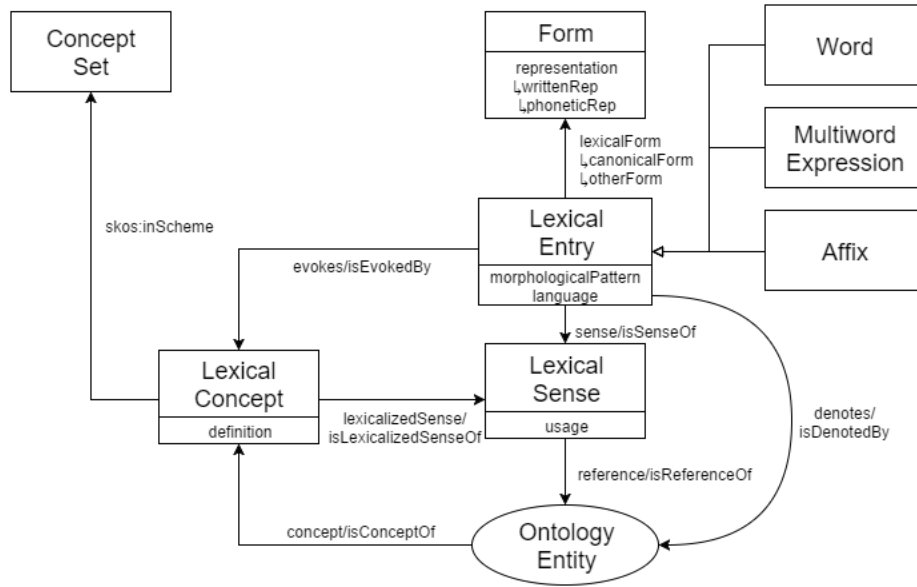


Fig. 8. The OntoLex-Lemon Core.

Urdu [\[edit\]](#)

Etymology [\[edit\]](#)

From Persian زبان (*zobān, zabān*).

Pronunciation [\[edit\]](#)

- IPA^(key): /zʊ.bɑːn/, /zə.bɑːn/

Noun [\[edit\]](#)

زبان • (*zubān, zabān*) *f* (*Hindi spelling ज़बान*)

- tongue (body part)

Synonyms: جیبھ (jībh), لسان (lisān)

- language

Synonyms: لسان (lisān), بھاشا (bhāshā), بولی (bolī)

Declension [\[edit\]](#)

| Declension of زبان [less ▲] | | |
|--|--------------|-----------------|
| | singular | plural |
| direct | زبان (zubān) | زبانیں (zubānē) |
| oblique | زبان (zubān) | زبانوں (zubānō) |
| vocative | زبان (zubān) | زبانو (zubāno) |

Derived terms [\[edit\]](#)

- مادری زبان (mādrī zabān)

Fig. 9. زبان (*zuban*)

```

1      :زبان a ontolex:Word;
2      lexinfo:gender lexinfo:feminine;
3      lexinfo:partOfSpeech lexinfo:noun;
4      lime:language "ur"^^xsd:language;
5      ontolex:canonicalForm :زبان_lemma;
6      ontolex:denotes <http://dbpedia.org/resource/Language>,
7                      <http://dbpedia.org/resource/Tongue>;
8      ontolex:otherForm :زبان_dir_pl,
9                        :زبان_obl_pl,
10                       :زبان_obl_sing,
11                       :زبان_voc_pl,
12                       :زبان_voc_sing;
13      ontolex:sense :زبان_sense1, :زبان_sense2 .
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

```

Fig. 10. The Word زبان.

we see the link to the two senses from the Wiktionary entry encoded as `زبان_sense1` and `زبان_sense2` respectively, both of which senses link to DBpedia pages and the `denotes` property which links the entry straight to these two. We show a form, the direct plural form of the noun, and a sense, the second sense of the word meaning ‘language’. Both of these are presented in Figure 11. Here we can see the use of the `writtenRep` property to associate both the roman and arabic script versions of the word’s orthography, and the use of reference to link the sense of the word to its ontological reference.

A.2. Syntax and Semantics

The next module we will look at is the *Syntax and Semantics* module of OntoLex-Lemon, shortened as *synsem*. It contains a number of classes and properties for modelling the syntactic behaviour of words and the relationship(s) between this syntactic behaviour and the semantic properties of those words. Figure 12 presents these classes and properties. The two classes used to model syntactic behaviour are *Syntactic Frame*²⁴⁵ and *Syntactic Argument*²⁴⁶. These can be related to corresponding ontological predicates via the *OntoMap* class and various object properties. To describe these mechanisms here would take us too far beyond the purpose of this brief introduction. However

²⁴⁵This "represents the syntactic behavior of an open class word in terms of the (syntactic) arguments it requires". See <https://www.w3.org/2016/05/ontolex/#syntactic-frames>

²⁴⁶This "represents a slot that needs to be filled for a certain syntactic frame to be complete." See <https://www.w3.org/2016/05/ontolex/#syntactic-frames>.

the interested reader is invited to consult the OntoLex-Lemon guidelines.

A.3. Decomposition

The OntoLex-Lemon *Decomposition* module (*decomp* for short), represented in Figure 13, is designed to indicate "which elements constitute a multiword or compound lexical entry"²⁴⁷. It does this via the *Component* class and the properties *subterm*, *constituent*, and *correspondsTo*.

For instance, to return to our earlier example, take the derived term which was listed in the Wiktionary entry for زبان, namely *مادر زبان* *madri zuban*. This expression literally means ‘mother tongue’ or ‘native language’ and can be decomposed into two subterms using the *decomp* property *subterm*. That is it can be decomposed into the two entries زبان and *مادر*, the latter meaning ‘maternal’, as in Figure 14 (note that we class *madri zuban* as a *MultiWordExpression*).

Other examples and details of the other properties can be found, as always, in the OntoLex-Lemon guidelines.

A.4. Variation and Translation

Finally in this appendix we will look at the *Variation and Translation* module (also known as *vartrans*). This module is concerned with representing "relations between lexical entries and lexical senses that are variants of each other." In this context these

²⁴⁷<https://www.w3.org/2016/05/ontolex/#decomposition-decomp>

```

1  :زبان_dir_pl a ontolex:Form;
2  lexinfo:number lexinfo:plural;
3  ontolex:writtenRep "zubānē", "زبان"@ur .
4
5  :زبان_sense2 a ontolex:LexicalSense;
6  ontolex:reference <http://dbpedia.org/resource/Language> .
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

```

Fig. 11. A form and a sense belonging to زبان.

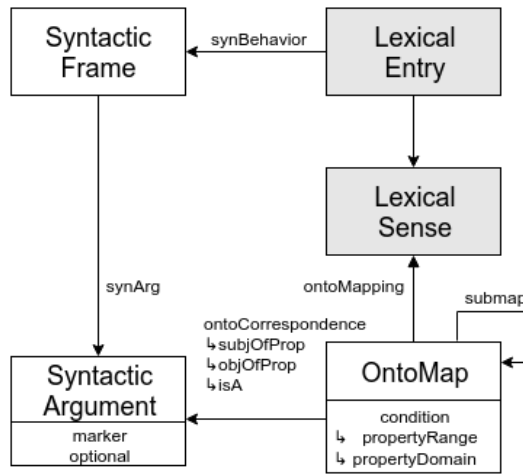


Fig. 12. The OntoLex-Lemon Synsem Module.

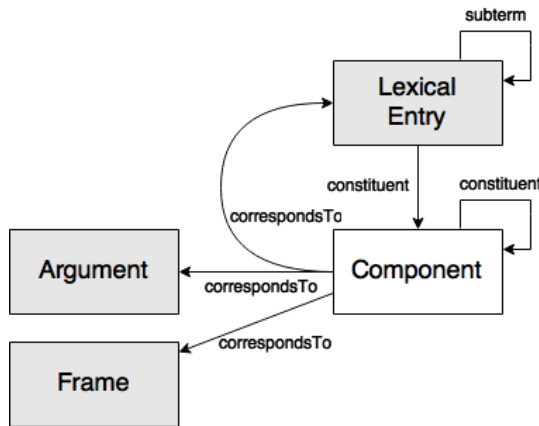


Fig. 13. The OntoLex-Lemon Decomposition Module.

entries can be variants via lexico-semantic relations such as hyponymy and synonymy or because they are translations. The classes and properties of the vartrans module are shown in Figure A as can be seen the vartrans module reifies lexico-semantic relations between

words, defining a class Lexico-semantic Relation with subclasses, Lexical Relation and Sense Relation .

We will illustrate the use of this module with two examples from our running زبان example .In the first case we relate the word زبان to one of its synonyms,

```

1      :MWE_مادری_زبان a ontolex:MultiWordExpression;
2      lime:language "ur"^^xsd:language ;
3      decomp:subterm :زبان,
4      :مادری .

```

Fig. 14. Multi-word Expression Example.

9 باشا *bhasa* ‘language’ – or rather we relate the sense
10 of the word meaning ‘language’ to a sense of the word
11 which means the same thing. This is represented in Fig
12 15 using vartrans properties and classes. We specify
13 the relation as a Sense Relation, use the LexInfo class
14 synonym to categorise it and specify its source and tar-
15 get using the vartrans source and target objects respec-

9 tively. In Fig 16 on the other hand we represent a trans-
10 lation relation between and the first sense of the En-
11 glish word *language*. This is done in much the same
12 way as in the previous example, except that we spec-
13 ify that the relation is one of direct equivalence using
14 the property category and a term from a terminological
15 vocabulary.

```

1      :senseRelation a vartrans:SenseRelation;
2      vartrans:category lexinfo:synonym;
3      vartrans:source :پهاتسا_sense1;
4      vartrans:target :زبان_sense2 .
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
    
```

Fig. 15. Example of the use of vartrans.

```

16      :translationRelation a vartrans:Translation;
17      vartrans:category <http://purl.org/net/translation-categories#directEquivalent>;
18      vartrans:source :زبان_sense2;
19      vartrans:target :language_sense_1 .
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
    
```

Fig. 16. Example of the use of vartrans.

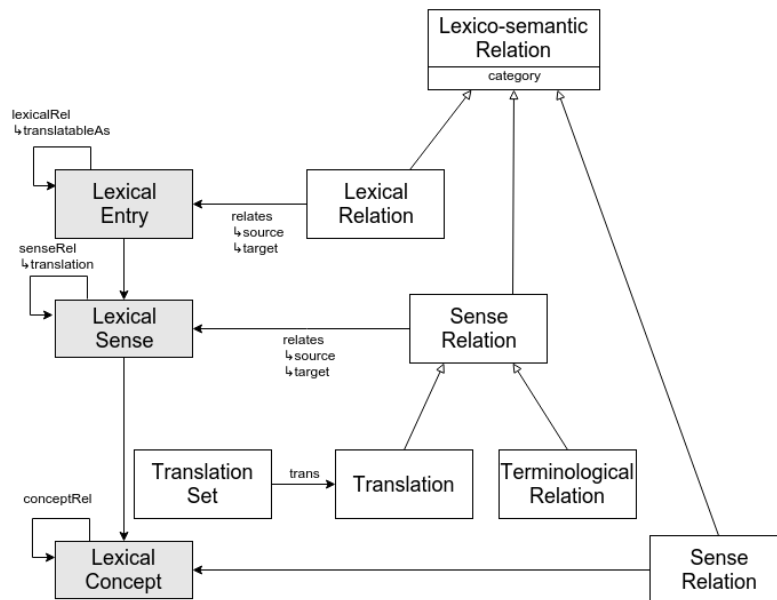


Fig. 17. The OntoLex-Lemon Vartrans Module.