

Ontology Reuse: the Real Test of Ontological Design

Piotr SOWIŃSKI ^{a,b}, Katarzyna WASIELEWSKA-MICHNIEWSKA ^b,
Maria GANZHA ^{a,b}, Marcin PAPRZYCKI ^b, Costin BĂDICĂ ^c

^a *Warsaw University of Technology, Warsaw, Poland*

^b *Systems Research Institute, Polish Academy of Sciences, Warsaw, Poland*

^c *University of Craiova, Romania*

Abstract.

Reusing ontologies in practice is still very challenging, especially when multiple ontologies are (jointly) involved. Moreover, despite recent advances, the realization of systematic ontology quality assurance remains a difficult problem. In this work, the quality of thirty biomedical ontologies, and the Computer Science Ontology are investigated, from the perspective of a practical use case. Special scrutiny is given to cross-ontology references, which are vital for combining ontologies. Diverse methods to detect potential issues are proposed, including natural language processing and network analysis. Moreover, several suggestions for improving ontologies and their quality assurance processes are presented. It is argued that while the advancing automatic tools for ontology quality assurance are crucial for ontology improvement, they will not solve the problem entirely. It is ontology reuse that is the ultimate method for continuously verifying and improving ontology quality, as well as for guiding its future development. Specifically, multiple issues can be found and fixed primarily through practical and diverse ontology reuse scenarios.

Keywords. Ontology evaluation, Ontology engineering, Semantic Web, Ontology reuse, Ontology alignment

1. Introduction

Ontologies are invaluable in representing domain knowledge; possessing a number of characteristics that make them a tempting choice when designing intelligent information systems [1]. They are very expressive and flexible, allowing one to describe a wide variety of subjects, in a precise manner. Nowadays, with the introduction of the Linked Open Data principles [2], accessing, reusing, and combining ontologies should have become easier than ever. However, in practice, ontology reuse can still be very challenging [3]. This is especially true when multiple ontologies are involved, where ambiguous or conflicting statements between ontologies often require human intervention in order to be resolved [4].

Obviously, high quality of an ontology is vital to ensure its potential reusability across a spectrum of applications. As ontologies model knowledge, any issues in them immediately reflect on the quality of their application. Many criteria for evaluating ontologies were proposed, and there are multiple approaches to their evaluation [5, 6].

However, individual ontologies use very different evaluation methods, or (in most cases) do not have any evaluation plan in place [7]. Some approaches are hard to implement in practice – such as the gold standard method, for which, often, there is no available gold standard to be used. Obviously, expert-based evaluation is resource-intensive and naturally limited in its scope. On the other hand, the automated, more universal methods are usually too general to identify many issues in the ontologies [8].

Overall, as of today, there do not seem to be any established, agreed upon, and widely used standards and methodologies for ontology evaluation. Instead, for years, there have been examples of evaluations that were *tailored to the specific ontology* [9, 10] or motivated by a *specific, practical use case* [8, 11, 12]. Lack of a commonly accepted evaluation methodology is one of problems negatively impacting potential reuse and adaptation of ontology-based solutions, as it was investigated, for instance, for the enterprise domain [13].

In earlier work [14] an unsupervised, ontology-based method for classifying scientific publications has been described (see Section 2). During research, when the classifier was to be adapted to different domains, numerous problems with the used ontologies have been noticed. Each of those issues either (i) increased the amount of work necessary to adapt the classifier to the domain, or (ii) decreased the accuracy of predictions. These observations motivated further research, which is reported in this contribution.

Specifically, the goal of this work is to demonstrate that the *varied and demanding reusability scenarios* provide the ultimate measure of the quality of an ontological design. Here, the real-world-anchored use case motivates exploring each issue and supports the conclusions drawn from the experiments.

In this context, in the following sections, ontologies from two different fields of science are investigated. To provide additional context to the study, Section 2 briefly describes the publication classification method that serves as the driving use case. Section 3 explores thirty diverse biomedical ontologies, related to the field of food safety. Section 4 examines a medium-sized Computer Science Ontology.

It should be stressed that all code used in the reported work, as well as additional results and materials, are available on GitHub, under an open-source license ¹. Detailed instructions for reproducing the experiments are also provided there.

2. Use Case Description

The investigation presented in this contribution is driven by a specific use case from our earlier work [14] – an unsupervised method for classifying scientific publications, that aims to be domain-agnostic. Here, all domain knowledge necessary to perform the classification is obtained from ontologies relevant to the pertinent field of science.

The considered method uses the publication’s title, abstract, and (optionally) unstructured keywords as an input. It then assigns to the article several entities from the loaded ontologies that best describe its topic. The classification pipeline consists of six steps. (1) Using a neural named entity recognition method, possible mentions of entities of interest are identified in the text. (2) The ontologies are searched for candidate entities that may match the mention. This is done using a text index. (3) The neighborhoods (most

¹<https://github.com/Ostrzyciel/ontology-quality-2022>

closely related entities) of the candidate entities are retrieved from the ontologies. (4) Using text embeddings and string matching, the candidate entities are compared to the mention in the text. Here, both the context of the mention (its sentence) and the context of the entity (its neighborhood) are taken into account. This two-way context sensitivity helps the method to correctly classify ambiguous mentions in the text. (5) The relations between the identified candidate entities are examined (entity-entity coherence). Then, the least consistent (least strongly connected) candidates are removed. (6) For the final result, the highest-scoring and most frequently occurring candidate entities are selected. The result is further enhanced with their parent entities. This allows for classification of publications with topics that are only implied and never mentioned explicitly.

The method relies solely on the ontologies as its source of domain knowledge, utilizing both textual labels (search index, text embeddings) and relations between entities (neighborhood estimation, entity-entity coherence, results enhancement). Unsurprisingly, the method's performance was observed empirically to be largely determined by the quality of the used ontologies. Following the ontology evaluation criteria as described by Raad & Cruz [6], accuracy, completeness, conciseness, consistency, and clarity of the ontologies had a direct, significant impact on the classifier's performance. The observed issues are discussed and investigated in detail in the following sections.

3. Evaluation of OBO Foundry Ontologies

Food safety was the first field, to which the publication classification method was applied. As a source of domain knowledge thirty ontologies from the Open Biological and Biomedical Ontology (OBO) Foundry [15] were used, along with SNOMED CT [16] (full list of used ontologies is available on GitHub²). Here, we focus only on the OBO ontologies, which range in size from 1 221 triples for BFO (Basic Formal Ontology) to 16 105 832 triples for NCBITaxon. When combined, the thirty ontologies contained 39 156 273 triples. These ontologies are highly heterogeneous, due to their different goals and being, typically, maintained by independent groups. Most considered ontologies are *developed* in the OBO format³, which is then *translated* into the standard Web Ontology Language (OWL), to facilitate reuse.

The issue of reliably maintaining quality over independently developed ontologies is an active area of research, with tools such as ROBOT [17] or the OBO Dashboard [18], attempting to tackle it. OBO Dashboard, in particular, performs a wide selection of tests on ontologies, essentially operationalizing the review activities that would otherwise be performed by a human. However, the set of tests is limited – for example the validity of cross-ontology references is not checked. Most of the issues have to be fixed manually, which proves difficult in practice, especially for the rarely-updated ontologies. Nevertheless, it can be stipulated that OBO Foundry is making significant efforts to assure the high quality of its ontologies. Keeping this in mind, let us report issues that have been found.

²<https://github.com/Ostrzyciel/ontology-quality-2022/blob/main/obo/ONTOLOGIES.md>

³https://owlcollab.github.io/oboformat/doc/G0.format.obo-1_4.html

3.1. *Ease of Access*

Before the ontologies could be used, they had to be first downloaded and loaded into a triple store. All ontology files downloaded from the OBO Foundry had the `.owl` extension, but their actual format varied. Most were distributed in the RDF/XML format, which is in line with the official OBO Foundry guidelines. However, the Common Anatomy Reference Ontology (CARO) used the OWL functional syntax instead, and had to be converted into RDF before being loaded. The Food-Biomarker Ontology (FOBI) was distributed in the OWL/XML format, which could not be converted to RDF/XML, due to invalid element errors. Thus, for FOBI, an earlier release was used that did not exhibit this issue and could be successfully converted to RDF.

It is worth pointing out that although these issues may appear as *minor annoyances*, they have significant implications. Inconsistency in file formats makes automatic reuse of these ontologies much harder or even impossible, due to the need for manual conversion of the files.

3.2. *Rarely Used Properties*

The publication classifier requires the user to assign weights to all object properties that may be useful for discovering how different entities are interrelated (steps 3 and 5 of the method). Similarly, weights are assigned to text properties, which describe the entities (step 2). To ascertain these weights, an obvious first step is to retrieve a list of all properties, sorted by the number of times they were used. Although OBO Foundry ontologies strive for reuse and have common upper ontologies defining relations, in the resulting lists many rarely used properties were found, often seeming like mistakes. Thus, those properties that had at most ten unique uses across all ontologies have been manually inspected.

It was found that out of 140 rarely-used properties, 74 were in some way erroneous, while being used across 278 unique triples. The most common type of error (51 properties, 212 occurrences) were undefined properties in the `oboInOwl` namespace [19]. These are, most likely, errors in translation from the OBO format to OWL. Other issues included mistaken ontology prefixes, typos, and invalid URIs. These issues were found in 24 out of thirty ontologies.

3.3. *Property Value Type Mismatch*

Next, the classifier requires dividing properties into URI-valued and literal-valued. This turned out to be a non-trivial task, as many properties were used inconsistently – sometimes their objects were URIs, and sometimes they were literals.

A list of all properties was created, with the number of times each was used in triples in which the object is an URI, a blank node, or a literal. For manual review, properties that exhibited conflicting usage patterns were selected. Most did not have an `rdfs:range` property, which would define formally their set of possible values. In fact, some properties were not defined at all, in any ontology. In other cases, the range was specified very broadly, allowing both URIs and literals. This presented a major challenge in distinguishing between erroneous and valid annotations. In such cases it had to be decided “intuitively” whether a given use is valid or not. The decision was made on the basis

of scarce documentation, or other occurrences of the same property. After establishing the valid usage patterns, lists of erroneous occurrences for each property were obtained, using SPARQL queries.

In the results, the `oboInOwl:hasDbXref` property was excluded, as it is discussed separately in the following subsection. Overall, 12 296 erroneous triples were identified in 20 ontologies. Often, in places where an URI-typed value was expected, a string literal representing that URI was present. Thus, the value itself was valid, but the datatype did not match. In other cases, various ontologies had different conventions for applying the same property. For some metadata properties, e.g., those for marking authorship, significant inconsistencies in their use were observed. However, it was decided not to mark them as errors, due to our lack of knowledge as to which convention should be applied, in a given case. Nonetheless, the lack of a common convention for representing metadata makes understanding and reusing ontologies considerably harder.

3.4. *Inconsistent Cross-ontology References*

Food safety integrates a wide variety of topics, from microbiology, to chemistry, and human diseases. Therefore, it is crucial to ensure that the imported ontologies are well-connected. Such connections, in OBO Foundry, can be either partial imports, (where an entity from one ontology is imported into another), or cross-ontology references, made with various types of properties. When examining inter-ontology references, it was observed that their usage patterns varied and their targets were hard to resolve automatically. Thus, all occurrences of cross-ontology references were retrieved and their targets analyzed. For these references, the dominant property is `oboInOwl:hasDbXref` (3 809 415 uses), but some SKOS [20] properties are also used. There were 259 occurrences of `skos:*Match` properties, of which only three had valid URIs as values. Others had literals of external identifiers, e.g., `MESH:C536189`.

The SKOS references were not included in further analysis. Instead, two standard OBO properties from the `oboInOwl` namespace were analyzed: `hasDbXref` and `hasAlternativeId`. Their values are supposed to be strings, having the form of `Namespace:Identifier`, with the namespace (usually) corresponding to an OBO Foundry ontology. All triples (3 908 752) containing these properties were retrieved. Then, the obtained reference values were processed and classified using a custom algorithm. The found namespace identifiers were cross-referenced with the publicly available list of all OBO Foundry ontologies⁴, to filter out valid references. Out of all references, 52 122 (1.3%) were URIs, not identifiers. Among them only 112 were valid URIs of entities in OBO Foundry. Others pointed to very diverse resources, e.g., Wikipedia, GitHub, online databases, scientific publications, structured vocabularies, and various other websites.

Among non-URI references, 187 167 (4.8%) referred to OBO Foundry ontologies. Others, usually, pointed to external biomedical databases and structured vocabularies, such as MeSH or SNOMED CT. However, multiple prefixes were used for some databases, (e.g., `SCTID`, `SMID`, `SNOMEDCT_US`, `SNOMEDCT_US_2021_09_1` for SNOMED CT), which complicates translating the identifiers into URIs. Using the recently created Bioregistry [21], it was possible to resolve prefixes for another 3 013 315 references (77.1%). Nonetheless, the identified databases form an eclectic collection, in-

⁴<https://obofoundry.org/registry/ontologies.yml>

cluding the Unified Medical Language System, PubMed, English Wikipedia, ISBN, ORCID, and Google Patents. This shows that the cross-reference properties are used for a wide variety of purposes. Obviously, this complicates their potential reuse. Moreover, the references did not always conform to the `Namespace:Identifier` format, with some using the underscore (`_`) as a separator instead of a colon. For OBO identifiers, there were 19 such cases, while for Bioregistry-resolved prefixes 26 572 were identified. Finally, 214 ($< 0.01\%$) references pointed to empty RDF blank nodes. These are, most likely, the results of errors during modeling, or translation into OWL. All these findings are summarized in Table 1.

Table 1. Cross-ontology references in OBO Foundry

Type	# references	% total
URI	52 122	1.33%
valid OBO URI	112	$< 0.01\%$
en.wikipedia.org	17 365	0.44%
orcid.org	8 085	0.21%
langual.org	5 592	0.14%
other	20 968	0.54%
Textual	3 856 416	98.66%
recognized OBO identifier	187 167	4.79%
recognized other Bioregistry prefix	3 013 315	77.09%
unknown	655 934	16.78%
Empty blank node	214	$< 0.01\%$
Total	3 908 752	100.00%

It is worth mentioning that the OBO ontology format has a number of header tags (`treat-xrefs-as-*`) that allow one to more precisely describe how to treat cross-references, on a per-ontology basis. This information, however, is lost in the translation to OWL, and thus largely inaccessible.

The (in)consistency of OBO cross-ontology references has also been explored by Laadhar et al. [10], who noted the overwhelming variety of uses for the cross-reference properties, which hampers reuse. Their work gives valuable suggestions of possible improvements, which augment the discussion presented here.

3.5. Summary

In summary, 720 844 issues were found across thirty OBO Foundry ontologies (treated jointly). These issues, divided into five categories, and broken down by the source ontology, have been presented in Table 2. Note that due to the use of partial imports (cf. MIREOT [22]), some ontologies overlap. Thus, the actual number of issues in a given category may be lower than the sum of issues encountered in individual ontologies. It is also worth noting that *only two ontologies* were found to be entirely problem-free, from the perspective of the considered use case. These were BFO and NCBITaxon.

4. Evaluation of the Computer Science Ontology

The same classifier was adapted to the Computer Science Ontology (CSO) [23], which models research topics in computer science. Its main areas of application include publi-

Table 2. Summary of issues found in OBO Foundry ontologies

Ontology	Rare prop.¹	Prop. obj.²	Xref: blank³	Xref: URI⁴	Xref: unk.⁵
AEO	4	0	0	10	136
AGRO	18	51	0	1 266	6 710
APOLLO-SV	4	308	214	2	21
BFO	0	0	0	0	0
BTO	3	0	0	0	3 479
CARO	1	6	0	380	1 800
CHEBI	12	0	0	0	313 736
CL	38	236	0	2 297	34 296
DOID	2	2	0	1	12 824
DRON	9	6	0	0	35 148
EHDAA2	3	0	0	5	67
ENVO	3	1 612	0	3 299	1 649
FOBI	5	0	0	0	0
FoodOn	0	5 702	0	8 416	6 329
GAZ	0	6	0	0	25 505
GO	1	2 536	0	354	118 473
HP	45	313	0	3 520	28 386
IAO	0	22	0	0	0
MP	47	388	0	15 253	37 229
NCBITaxon	0	0	0	0	0
OBI	0	1 295	0	0	0
PATO	13	96	0	3 485	17 144
PCO	3	19	0	9	41
PECO	2	0	0	0	685
PO	3	24	0	3	6 547
RO	2	35	0	0	15
SYMP	2	0	0	1	449
Uberon	87	375	0	23 845	14 627
UO	9	0	0	0	0
XCO	3	0	0	0	494
All	278	12 296	214	52 122	655 934

¹ Invalid occurrences of rarely-used properties.² Property object type mismatch (URI instead of literal or vice versa).³ Cross-references pointing to blank nodes.⁴ Cross-references pointing to URIs instead of identifiers.⁵ Non-resolvable cross-reference identifiers.

cation classification and research trend detection. It uses a small set of semantic properties, extending the SKOS data model [20]. It also contains references to entities in external knowledge bases (KBs) such as DBpedia [24], Wikidata [25], and YAGO [26], using the `owl:sameAs` property. The most recent version of CSO (3.3) contains 14 290 topics and 163 470 triples, and was used in all experiments.

The ontology was generated automatically, using the Klink-2 algorithm [27]. Relationships in some research areas (Semantic Web and Software Architecture) were manually reviewed by experts, but most of the burden of ontology quality control and main-

tenance was left to the community. To facilitate this, the authors of CSO built the *CSO Portal*⁵, which collects suggestions for improving the ontology.

In the following subsections, we examine various issues found in CSO. To assure that the results properly identify the magnitude of individual issues, each experiment was performed independently, starting from a fresh copy of the ontology.

4.1. Synonym Description Structure

In CSO, alternative names of one concept are treated in such a way that each synonym is its own entity of type `cso:Topic`. Clusters of equivalent topics are connected using the `cso:relatedEquivalent` and `cso:preferentialEquivalent` properties. There are 11 187 such clusters. This structure is in a stark contrast to most other ontologies, which usually attach alternative names to a single entity, using text properties (cf. Gene Ontology and other OBO ontologies [28], Open Research Knowledge Graph [29], YAGO [26]).

When processing the ontology, it was found to be laborious to work with this structure. For example, when displaying the results of a query searching for topics, one would have to filter out those that are not marked as *preferential*. Similarly, adapting the publication classification method to CSO would be impossible without writing additional logic to handle the atypical structure. Therefore, CSO was transformed using SPARQL 1.1 updates into an ontology with synonymous topics merged into a single entity. Here, synonyms were attached using the `skos:altLabel` property. The transformation changed the structure of the ontology, but it did not affect its information content. The original 163 470 triples were transformed into 88 227 triples, a 46% decrease. An example of the applied transformation can be found in Figure 1.

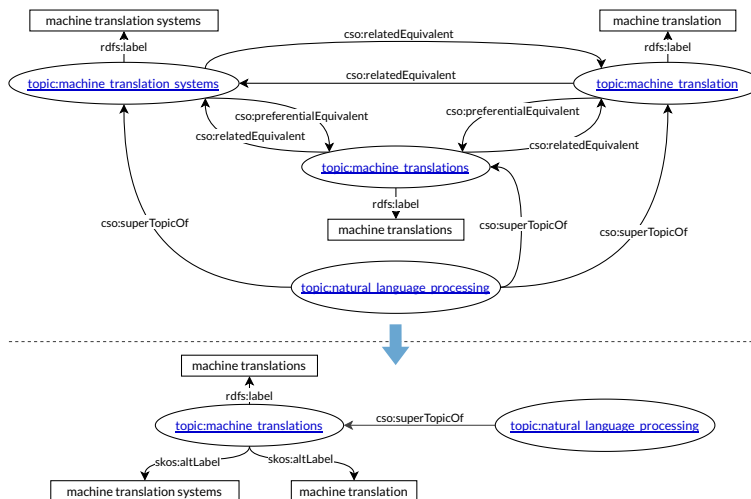


Figure 1. An example of the changes made in CSO

⁵<https://cso.kmi.open.ac.uk/>

4.2. Externally Inconsistent References

Next, CSO’s references to external KBs were examined. The goal was to be able to use federation to provide additional information to the classifier. However, multiple references were found to point to entities very distant from computer science, such as musicians or rivers. Assuming that in the external KB there should not be many connections between, e.g., popular music and computer science, one could attempt to find the “outlier” entities that are not well connected with the rest. These entities will, most likely, be invalid references.

A list of all DBpedia entities referenced from CSO was obtained. Then, for each entity, a list of links to other DBpedia entities was retrieved. This data was used to construct an undirected wikilink graph that was analyzed using the NetworkX library [30]. Three approaches to identifying outliers were used in a sequence:

- T1** Find all connected components (CCs) in the graph. Add the vertices from all except the largest CC to the set of candidate outliers (result set). Afterwards, remove these vertices from the graph, to leave only one CC.
- T2** Identify and remove bridges in the graph, increasing the number of CCs [31]. Add the vertices from the newly created CCs to the result set, and then remove them from the graph.
- T3** Find “communities” in the graph, using the label propagation algorithm [32]. Add the vertices from these communities to the result set.

The outlier detection method identified 140 potentially erroneous alignments, which were then manually evaluated by a reviewer, using a custom web application. As the errors were usually easily discernible and there was little room for ambiguity, only one reviewer was included in this experiment. Here, each alignment has been marked as *valid*, *invalid* or having a different issue, such as being out of scope of computer science entirely. Table 3 presents the results of the review, broken down by the tactic that led to marking the alignment as “suspect”. It can be observed that tactics **T1** and **T2** were highly effective in discovering potentially erroneous alignments. The false positive rate was relatively low, at only 25% of all suspected entities not having any issues. Most occurrences in the *other issues* category referred to out-of-scope entities, which explains why their external references were inconsistent with others.

Table 3. DBpedia alignments manual review results

Verdict	# total	# T1	# T2	# T3
invalid	74	47	27	0
other issues	31	14	15	2
valid	35	15	13	7
Total	140	76	55	9

4.3. Missing References to the Corresponding Knowledge Bases

CSO contains references to several external KBs, most of which have corresponding entities for each other. For example: each DBpedia entity refers to an English Wikipedia article, and each such article has a corresponding entity in Wikidata. Thus, one may

expect there to be a 1:1 mapping between CSO’s references to DBpedia and Wikidata. To examine this, a series of SPARQL queries was performed, looking for instances where this 1:1 mapping was not the case.

The method did not identify any missing DBpedia references. However, it did find 31 topics with missing references to 13 unique Wikidata items. Fixing the issue is trivial, and a patch in the form of a Turtle file⁶ was prepared.

4.4. Logically Invalid References

The publication classifier uses relationships contained in the ontology to determine the degree of similarity between entities. Each property type is assigned a weight, which should be low for very close conceptual connections (e.g., identity), and high for the more “remote” relationships. When assigning the weights, one can assume that the `owl:sameAs` relation, according to its definition, should represent identity (and thus, have zero weight). However, CSO uses the `owl:sameAs` property for referencing external knowledge bases. This property is defined as reflexive, symmetrical, and transitive [33], which implies that if two different concepts were to be aligned to the same external entity, they could be deduced to be identical, which is a logic violation. Multiple examples of such problems have been noticed (see, Figure 2). Here, for example, *malicious software* and *malware detection* can be reasoned to be identical, even though in the ontology they are (as expected) completely separate concepts.

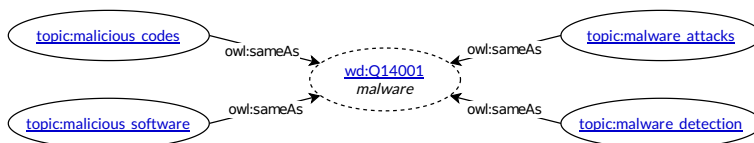


Figure 2. References to external KBs in CSO

malicious software and *malware detection* can be reasoned to be identical, even though in the ontology they are (as expected) completely separate concepts.

To find all such occurrences, the set of `owl:sameAs` relations in CSO was expanded, using a reasoner that exploits the reflexive, symmetrical, and transitive characteristics of this property. A SPARQL query was used to obtain a list of CSO topic pairs, marked as identical (`owl:sameAs`), but not belonging to the same cluster of synonyms. A total of 7 457 unique URI pairs was found to be incorrectly marked as identical. This issue concerned 3 250 unique topic URIs, out of 14 290 CSO topics in total.

For comparison, in the BioPortal, which hosts multiple biomedical ontologies, the alignments that were constructed with the LOOM lexical matching algorithm [34] are marked using the `skos:closeMatch` property, chosen due to its lack of transitivity, which helps avoid similar issues.

4.5. Intra-cluster Alignment Inconsistencies

CSO’s approach to synonym clusters requires caution to avoid inconsistencies in references to external KBs. For example, having topics *A* and *B* from a single synonym cluster, one would expect that if *A* is aligned to an external entity *X*, then also *B* is aligned

⁶https://github.com/Ostrzyciel/ontology-quality-2022/blob/main/cso/3_missing_refs/results/suggestion.ttl

to X . Moreover, CSO topics from one cluster should not be aligned to multiple different external entities from a single external KB.

To validate these assumptions, SPARQL queries were performed to find DBpedia alignments for each synonym cluster. The information was then processed with a script to find clusters that violate either assumption. The method identified 130 synonym clusters that contained more than one unique DBpedia reference. Of those, 124 clusters had two unique references, and six clusters had three. Moreover, there were 962 entities with a DBpedia reference missing, spread over 752 synonym clusters. Adding these missing references is easy and an appropriate ontology patch was prepared⁷. In total, 882 out of 2 225 (39.6%) larger-than-one synonym clusters exhibited some external reference consistency issues.

4.6. Term Conflation

When debugging the outputs of the publication classifier, several synonym clusters were noticed to include conflated and erroneous terms. Finding more such issues systematically would require the intervention of human experts. To narrow down the set of “suspected” clusters, natural language processing techniques can be used to determine the degree of similarity between terms. Then, the least coherent synonyms could be inspected manually by human reviewers.

Hence, using SPARQL queries, a list of all synonym clusters was obtained, along with the corresponding topic labels. These labels were then encoded into vectors using a sentence embedding model `all-mpnet-base-v2`⁸, based on the MPNet transformer [35]. Next, in each cluster, an all-pairs similarity matrix was computed between the encoded labels (using cosine similarity, however, any other similarity measure could have been used). For each topic, the mean and the standard deviation of its similarity to other topics were computed. Finally, clusters that contained at least three topics and had low mean and standard deviation of similarity to other topics were selected for review. This approach was designed to find clusters with topics that are *systematically* inconsistent, and thus most likely to be wrong.

For manual review, suspected clusters were evaluated by three experts. Available ratings were: *definitely good*, *probably good*, *not sure*, *probably wrong*, and *definitely wrong*. The order of clusters was picked at random, independently for each expert. The experts were allowed to use external sources to aid their work, and were instructed to simply use their best judgment when making the decision.

The intra-cluster similarity method identified 115 clusters of potentially invalid synonyms. Using the simple majority vote rule, 84 of these clusters were deemed *wrong* by the reviewers. Moreover, 95 clusters were marked as such by at least one reviewer, and all three reviewers agreed that 58 clusters are wrong. To investigate inter-reviewer agreement, scores were converted to an ordinal scale from -2 to 2, with *definitely wrong* corresponding to -2, and *definitely good* corresponding to 2. Mean scores and standard deviations for each reviewer are presented in Table 4. It can be observed that reviewers had significantly different score distributions (Fleiss’ $\kappa = 0.501$, Krippendorff’s $\alpha = 0.474$). This indicates moderate agreement, which is desirable when seeking varied opinions on the synonymy of terms.

⁷https://github.com/Ostrzyciel/ontology-quality-2022/blob/main/cso/5_intra_ref_consistency/results/suggestion.ttl

⁸<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

Table 4. Term conflation: reviewer’s scores distributions

Reviewer	Mean score	St. dev.
Reviewer 1	-1.16	1.60
Reviewer 2	-0.85	1.67
Reviewer 3	-0.13	1.79

It was found that, in CSO, frequently conflated were terms surrounding a single subject. Here, as an example one can consider *classification system*, *classification tasks*, *classification results*. Other synonyms seemed of little use, introducing only singular/plural variations, or different variants of hyphen placement. Abbreviations were also occasionally mixed into the labels – for example *neural network (nn)* exists as a single term.

Additionally, the reviewers made plentiful remarks about scrutinized clusters. In particular, many out-of-scope clusters were found, frequently from the fields of pedagogy and genetics. This is consistent with the observations of Han et al. [36], who point out that CSO fails to distinguish computer science topics from those of related fields.

4.7. Summary

Overall, in above experiments, 3 137 synonym clusters and 4 287 topics were found to be affected by at least one issue. For cases where the fix is unambiguous, suggested patches in the form of Turtle files were prepared and made available.

5. Discussion and Concluding Remarks

Over the course of a use case-driven investigation, 31 heterogeneous ontologies were scrutinized with regard to their quality. Not only use case-specific issues were found, but also those related to basic principles of ontology design, clarity and consistency.

Overall, OBO Foundry ontologies present an extremely varied landscape of sizes, quality assurance methods, and stakeholders involved. Thus, it is interesting to observe that out of 30 studied ontologies, 28 were found to involve at least one class of problems. Moreover, in majority of cases, found problems concerned a very large number of triples (these were not “localized problems” but clear cases of systemic ones). One method to tackle them would be, for OBO Foundry, to further expand upon their set of guidelines, for example by specifying, which properties to use for describing metadata. Defining formally the ranges of properties could also prove valuable, as it would enable automatic tools to check for potential problems. Cross-ontology references also require improvements, to better conform with Linked Data standards and have clearer meaning. Finally, the OBO to OWL translation should be further improved, to better preserve semantic information.

In the case of CSO, it was found that some issues stem from the decisions made when designing its overall structure, while others can be attributed to the non-exhaustive quality assurance of the automatically generated ontology. It can be hoped that the reported results and suggested fixes will be valuable to the further development of CSO.

In the experiments, the usefulness of natural language processing methods and network analysis for ontology quality assurance were clearly shown. Most importantly,

these techniques can greatly reduce the amount of human labor required to perform the evaluation. Thus, ontological engineers should not be afraid to broaden their toolboxes with less conventional techniques. This is particularly the case since the *joint* use of approaches described above can substantially improve quality of ontologies, while considerably reducing the amount of work that has to be completed by human experts.

Although the automatic ontology inspection methods, such as the OBO Dashboard, do help with the problem of ontology quality, they cannot solve it on their own. There will always be limits to the scope of checks performed by automated, general tools. Thus, it is vital to pay special attention to the issues observed during reuse, and *continuously* integrate resulting observations to improve the ontologies. Reusability is, arguably, the *essence* of the Linked Open Data vision. Therefore, it should be the core driving factor for ontology development. It was shown that using multiple ontologies together is particularly hard, with numerous, still largely unresolved, issues regarding cross-ontology references. Thus, this topic deserves special attention in future research.

Overall, it was demonstrated how reusing ontologies, especially in demanding applications, serves as the real test of their quality. The purpose of the performed experiments is not to present them as general approaches to ontology quality assurance. Rather, they exemplify the process of continuous ontology reuse, evaluation, and improvement.

Acknowledgments

We would like to thank Anastasiya Danilenka and Jan Sawicki for their invaluable help in reviewing CSO's entities.

References

- [1] Noy N, McGuinness D. Ontology Development 101: A Guide to Creating Your First Ontology. Knowledge Systems Laboratory. 2001 01;32.
- [2] Bizer C, Heath T, Berners-Lee T. Linked data: The story so far. In: Semantic services, interoperability and web applications: emerging concepts. IGI global; 2011. p. 205-27.
- [3] Fernández-López M, Poveda-Villalón M, Suárez-Figueroa MC, Gómez-Pérez A. Why are ontologies not reused across the same domain? Journal of Web Semantics. 2019;57:100492.
- [4] Cota G, et al. The landscape of ontology reuse approaches. Appl Practices Ontol Des, Extraction, Reason. 2020;49:21.
- [5] Brank J, Grobelnik M, Mladenic D. A survey of ontology evaluation techniques. In: Proceedings of the conference on data mining and data warehouses (SiKDD 2005). Citeseer Ljubljana, Slovenia; 2005. p. 166-70.
- [6] Raad J, Cruz C. A survey on ontology evaluation methods. In: Proceedings of the International Conference on Knowledge Engineering and Ontology Development, part of the 7th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management; 2015. p. 179-86.
- [7] Amith M, He Z, Bian J, Lossio-Ventura JA, Tao C. Assessing the practice of biomedical ontology evaluation: Gaps and opportunities. Journal of biomedical informatics. 2018;80:1-13.

- [8] Al-Sayed MM, Hassan HA, Omara FA. Towards evaluation of cloud ontologies. *Journal of Parallel and Distributed Computing*. 2019;126:82-106.
- [9] Casellas N. Ontology Evaluation through Usability Measures. In: Meersman R, Herrero P, Dillon T, editors. *On the Move to Meaningful Internet Systems: OTM 2009 Workshops*. Berlin, Heidelberg: Springer Berlin Heidelberg; 2009. p. 594-603.
- [10] Laadhar A, Abrahão E, Jonquet C. Investigating one million XRefs in thirty ontologies from the OBO world. In: *11th International Conference on Biomedical Ontologies (ICBO)*; 2020. p. G.1-12.
- [11] Szmeja P, Wasielewska K, Ganzha M, Drozdowicz M, Paprzycki M, Fidanova S, et al. Reengineering and extending the Agents in Grid Ontology. In: *International Conference on Large-Scale Scientific Computing*. Springer; 2013. p. 565-73.
- [12] Ganzha M, Paprzycki M, Pawlowski W, Szmeja P, Wasielewska K. Towards Common Vocabulary for IoT Ecosystems - preliminary Considerations. In: Nguyen NT, Tojo S, Nguyen LM, Trawinski B, editors. *Intelligent Information and Database Systems - 9th Asian Conference, ACIIDS 2017, Kanazawa, Japan, April 3-5, 2017, Proceedings, Part I*. vol. 10191 of *Lecture Notes in Computer Science*; 2017. p. 35-45. Available from: https://doi.org/10.1007/978-3-319-54472-4_4.
- [13] Wasielewska-Michniewska K, Ganzha M, Paprzycki M, Denisiuk A. Application of ontologies in the enterprise – overview and critical analysis. In: *Proceedings of the Third International Conference on Information Management and Machine Intelligence: ICIMMI 2021, Algorithms for Intelligent Systems (AIS), TO APPEAR*. Springer; 2022. .
- [14] Sowiński P, Wasielewska-Michniewska K, Ganzha M, Paprzycki M. Topical Classification of Food Safety Publications with a Knowledge Base. *arXiv preprint arXiv:220100374*. 2022.
- [15] Smith B, Ashburner M, Rosse C, Bard J, Bug W, Ceusters W, et al. The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration. *Nature biotechnology*. 2007;25(11):1251-5.
- [16] Donnelly K, et al. SNOMED-CT: The advanced terminology and coding system for eHealth. *Studies in health technology and informatics*. 2006;121:279.
- [17] Jackson RC, Balhoff JP, Douglass E, Harris NL, Mungall CJ, Overton JA. ROBOT: a tool for automating ontology workflows. *BMC bioinformatics*. 2019;20(1):1-10.
- [18] Jackson RC, Matentzoglou N, Overton JA, Vita R, Balhoff JP, Buttigieg PL, et al. OBO Foundry in 2021: operationalizing open data principles to evaluate ontologies. *Database: The Journal of Biological Databases and Curation*. 2021;2021.
- [19] Moreira D, Mungall C, Shah N, Aitken S, Richter JD, Redmond T, et al. The NCBO OBOF to OWL Mapping. *Nature Precedings*. 2009:1-1.
- [20] Miles A, Bechhofer S. SKOS Simple Knowledge Organization System Reference. *W3C*; 2009.
- [21] Hoyt CT. An integrative registry of biological databases, ontologies, and nomenclatures; v0.0.6. Zenodo; 2020. Available from: <https://doi.org/10.5281/zenodo.4404608>.
- [22] Courtot M, Gibson F, Lister AL, Malone J, Schober D, Brinkman RR, et al. MIREOT: The minimum information to reference an external ontology term. *Applied Ontology*. 2011;6(1):23-33.

- [23] Salatino AA, Thanapalasingam T, Mannocci A, Osborne F, Motta E. The Computer Science Ontology: a large-scale taxonomy of research areas. In: International Semantic Web Conference. Springer; 2018. p. 187-205.
- [24] Auer S, Bizer C, Kobilarov G, Lehmann J, Cyganiak R, Ives Z. DBpedia: A nucleus for a web of open data. In: The semantic web. Springer; 2007. p. 722-35.
- [25] Vrandečić D, Krötzsch M. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*. 2014;57(10):78-85.
- [26] Tanon TP, Weikum G, Suchanek F. YAGO 4: A reason-able knowledge base. In: European Semantic Web Conference. Springer; 2020. p. 583-96.
- [27] Osborne F, Motta E. Klink-2: integrating multiple web sources to generate semantic topic networks. In: International Semantic Web Conference. Springer; 2015. p. 408-24.
- [28] Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology: tool for the unification of biology. *Nature genetics*. 2000;25(1):25-9.
- [29] Jaradeh MY, Oelen A, Farfar KE, Prinz M, D'Souza J, Kismihók G, et al. Open Research Knowledge Graph: next generation infrastructure for semantic scholarly knowledge. In: Proceedings of the 10th International Conference on Knowledge Capture; 2019. p. 243-6.
- [30] Hagberg A, Swart P, S Chult D. Exploring network structure, dynamics, and function using NetworkX. Los Alamos National Lab.(LANL), Los Alamos, NM (United States); 2008.
- [31] Bollobas B, Gehrung FW, Halmos PR. *Modern Graph Theory*. vol. 184 of Graduate Texts in Mathematics. New York, NY: Springer; 2013.
- [32] Cordasco G, Gargano L. Community detection via semi-synchronous label propagation algorithms. In: 2010 IEEE international workshop on: business applications of social network analysis (BASNA). IEEE; 2010. p. 1-8.
- [33] Bechhofer S, Van Harmelen F, Hendler J, Horrocks I, McGuinness DL, Patel-Schneider PF, et al. OWL web ontology language reference. *W3C recommendation*. 2004;10(2):1-53.
- [34] Ghazvinian A, Noy NF, Musen MA. Creating mappings for ontologies in biomedicine: simple methods work. In: AMIA Annual Symposium Proceedings. vol. 2009. American Medical Informatics Association; 2009. p. 198.
- [35] Song K, Tan X, Qin T, Lu J, Liu TY. MPNet: Masked and permuted pre-training for language understanding. *arXiv preprint arXiv:200409297*. 2020.
- [36] Han K, Yang P, Mishra S, Diesner J. WikiCSSH: extracting computer science subject headings from Wikipedia. In: ADBIS, TPDF and EDA 2020 Common Workshops and Doctoral Consortium. Springer; 2020. p. 207-18.