

## Research and Applications

# UMLS to DBpedia link discovery through circular resolution

John Cuzzola,<sup>1</sup> Ebrahim Bagheri,<sup>1</sup> and Jelena Jovanovic<sup>2</sup>

<sup>1</sup>Laboratory for Systems, Software and Semantics (LS3), Ryerson University, Ontario, Canada and <sup>2</sup>Faculty of Organizational Sciences (FOS), University of Belgrade, Belgrade, Serbia

Correspondence Author: John Cuzzola, Laboratory for Systems, Software and Semantics (LS3), Ryerson University, Ontario, Canada. E-mail: jcuzzola@ryerson.ca

Received 28 May 2017; Revised 16 January 2018; Editorial Decision 23 February 2018; Accepted 26 February 2018

## ABSTRACT

**Objective:** The goal of this work is to map Unified Medical Language System (UMLS) concepts to DBpedia resources using widely accepted ontology relations from the Simple Knowledge Organization System (skos:exactMatch, skos:closeMatch) and from the Resource Description Framework Schema (rdfs:seeAlso), as a result of which a complete mapping from UMLS (UMLS 2016AA) to DBpedia (DBpedia 2015-10) is made publicly available that includes 221 690 skos:exactMatch, 26 276 skos:closeMatch, and 6 784 322 rdfs:seeAlso mappings.

**Methods:** We propose a method called circular resolution that utilizes a combination of semantic annotators to map UMLS concepts to DBpedia resources. A set of annotators annotate definitions of UMLS concepts returning DBpedia resources while another set performs annotation on DBpedia resource abstracts returning UMLS concepts. Our pipeline aligns these 2 sets of annotations to determine appropriate mappings from UMLS to DBpedia.

**Results:** We evaluate our proposed method using structured data from the Wikidata knowledge base as the ground truth, which consists of 4899 already existing UMLS to DBpedia mappings. Our results show an 83% recall with 77% precision-at-one (P@1) in mapping UMLS concepts to DBpedia resources on this testing set.

**Conclusions:** The proposed circular resolution method is a simple yet effective technique for linking UMLS concepts to DBpedia resources. Experiments using Wikidata-based ground truth reveal a high mapping accuracy. In addition to the complete UMLS mapping downloadable in n-triple format, we provide an online browser and a RESTful service to explore the mappings.

**Key words:** Link Discovery, UMLS, DBpedia, Automated Vocabulary Mapping, Instance Matching

## INTRODUCTION

DBpedia is a crowd-sourced community project for extracting structured, multilingual information from Wikipedia to be made freely available on the Web in machine intelligible format based on Semantic Web standards.<sup>1</sup> It is the central component and the main inter-linking hub in the Linked Open Data (LOD) (<https://www.w3.org/wiki/SweoIG/TaskForces/CommunityProjects/LinkingOpenData>) cloud, a network of open structured datasets published on the Web according to the Linked Data principles.<sup>2</sup> LOD consists of several billion interlinked data points and covers a wide variety of domains such as geography, government, life sciences, media, social network-

ing, and scientific publications, to name a few. Whereas biomedical datasets constitute a large portion of the LOD cloud (As obvious from the LOD cloud diagram: <http://lod-cloud.net/>), and several of these datasets are connected to DBpedia, the complete integration of the UMLS Metathesaurus is still missing. If available, a mapping between DBpedia resources and UMLS concepts could provide several benefits to the biomedical community.

The work presented in this paper aims at providing a bridge connecting UMLS to DBpedia, in a manner that is both efficient, i.e., fully automated, and effective, i.e., highly accurate. In particular, the contribution of the presented work is 2-fold:

1. we introduce a method of automated link discovery between equivalent, near-equivalent, and related concepts originating from 2 large-scale knowledge bases (KBs), namely, UMLS Metathesaurus and DBpedia and
2. we release a publicly available complete mapping set between UMLS and DBpedia that can facilitate the integration of many biomedical and medical KBs, through UMLS, to the Linked Open Data cloud.

## BACKGROUND

### Significance

The significance of the UMLS to DBpedia mapping presented in this paper is multifold:

- Sophisticated text/data mining tasks depend on the availability of KBs built from diverse sources.<sup>3</sup> Wikipedia contains large amounts of scientific and medical data, and thus has been recognized as highly useful for setting up initial KB for biomedical projects.<sup>4</sup> It has also proven useful for estimating semantic similarity of gene pairs.<sup>5</sup> In particular, Dessi and Atzori demonstrated that Wikipedia's 10K+ articles about human genes allow for highly accurate assessment of gene similarity and detection of functional groups of genes. The machine-readable version of Wikipedia, DBpedia, is also a highly rich knowledge source with the additional advantage of enabling automated and machine-intelligible access to the knowledge it contains. For instance, Yamamoto et al.<sup>6</sup> used DBpedia to automatically extend a life science database of abbreviations and their long forms with additional descriptions of the long forms, thus enabling users to more easily select the correct long form for a particular abbreviation.
- As the central hub in the LOD cloud, DBpedia offers connection to numerous biomedical and other related datasets and KBs. Based on the latest statistics, DBpedia is connected to other LOD datasets through an estimated 50 million links. This indicates that DBpedia can serve as a hub for accessing diverse types of data for building rich KBs.
- Based on search engine ranking and page view statistics, the English Wikipedia is a prominent source of online health information.<sup>7</sup> DBpedia has the potential to be even more useful, as it provides grounds for building advanced applications that not only facilitate information search and retrieval, but also act proactively, e.g., applications that recommend resources a user has not explicitly asked for but might benefit from (see, e.g.,<sup>8</sup>). In addition, it can be used to further advance the current approaches for assessing the trustworthiness of online health information. For example, Park et al.<sup>9</sup> demonstrated that online health-related content annotated with Wikipedia concepts can be effectively used for building page-level and site-level classifiers aimed at differentiating between trustworthy and suspicious sites. It is reasonable to expect that the performance of such classifiers could be further improved if the Wikipedia concepts, identified in Web pages, are mapped to the corresponding UMLS concepts, thus allowing for a more precise semantic representation of health-related content of Web pages.
- Finally, a UMLS to DBpedia mapping can be relevant for bridging the gap between health-related jargon used by professionals and that used by the general public.<sup>10</sup> For instance, having examined 10 large online question corpora, Roberts and Demner-Fushman found that consumers, i.e., the general public, used significantly less medical terms than medical

professionals.<sup>11</sup> Likewise, consumers' questions were found to be closer to an open-domain language model, built on newswire and Wikipedia, than to a medical model, built on a sample from PubMed Central. This was further confirmed by Mrabet et al.<sup>12</sup> who demonstrated that combining an open-domain KB (i.e., DBpedia) with a biomedical KB (i.e., UMLS) could lead to a substantial improvement in identifying the main topics of consumer health questions. These findings suggest that DBpedia could be more suitable for semantic annotation, i.e., entity linking, of consumer questions, whereas UMLS would be more suitable for questions/answers coming from medical professionals; therefore, a mapping between UMLS and DBpedia can facilitate automated matching between (annotated) customers' questions and medical professionals' answers. In addition, it can be used to further improve the discovery and retrieval performance of systems for search and exploration of online content related to health and life sciences, such as DeepLife.<sup>13</sup> DeepLife's knowledge base covers a wide spectrum of biomedical entities, originating from UMLS and KnowLife,<sup>14</sup> thus covering the needs and terminology of health and life science professionals. If extended with DBpedia/Wikipedia entities, through the proposed mapping, it would be better able to match search requests by the general public.

There has already been work within the biomedical and health-care domains that employ open instance mapping platforms, such as Silk<sup>15</sup> and Link discovery framework for Metric Spaces (LIMES)<sup>16</sup> to map across medical terminologies. For instance, Tilahun et al.<sup>17</sup> used Silk to automatically link HIV-related data elements with data elements from Bio2RD, and LinkedCT. Bing et al.<sup>18</sup> used Silk to map concepts between biomedical entities to help discover the side effects of using thiazolidinedione classed drugs such as Rosiglitazone. Luciano et al.<sup>19</sup> used Silk to link proteomic, disease, and treatment data, to health records to find candidate patients for active clinical trials. Similarly, The Cancer Genome Atlas (<https://cancergenome.nih.gov/>) used LIMES to build a massive, publicly available, 30 billion triple datastore of genetic genome mutations to advance discoveries against this disease.<sup>20</sup> There has also been work that has performed terminology mapping without using open mapping platforms. For example, Lee et al.<sup>21</sup> have used heuristics for mapping laboratory terminology to Logical Observation Identifiers Names and Codes. Likewise, Kahn<sup>22</sup> has used semi-automated string matching to map Orphanet Rare Disease Ontology terms to the terms in the Radiology Gamuts Ontology. However, to the best of our knowledge, there has been no prior work that attempted to systematically map UMLS concepts to concepts from the widely used DBpedia knowledge base, thus facilitating the integration of UMLS with the Linked Open Data cloud.

### Ontological Representation of Equality Relations

When formally expressing links between 2 knowledge bases, the most common relation is "equal-to,"<sup>23</sup> often asserted using the predicate *sameAs* in the *Web Ontology Language* (<https://www.w3.org/OWL/>), or by *exactMatch* in the *Simple Knowledge Organization System* (SKOS) (<https://www.w3.org/2004/02/skos/>). The primary difference is *owl:sameAs* represents true equivalence in that every property of concept *x* is in the ontology of *y* and vice versa, whereas *skos:exactMatch* asserts that resource *x* is an exact match to resource *y* when both *x* and *y* can be used interchangeably for a wide range of information retrieval tasks. The predicate *skos:closeMatch* is similar to *skos:exactMatch* but does not necessarily preserve *transitivity*. We intentionally avoid making the assertion of *owl:sameAs*

because of strict equivalence requirements opting for `skos:exactMatch/closeMatch` as better choices given the published W3C standards. Furthermore, our method also considers the “seeAlso” property of the *Resource Description Framework Schema* (RDFS) that asserts that information about *x* might be available through resource *y*.

METHODS

Algorithm

We pair 4 semantic annotation tools to perform link discovery between UMLS and DBpedia. Two pairings of annotators link UMLS concepts to DBpedia resources while the remaining pair links from DBpedia to UMLS concept-unique-identifiers (CUI). We label the DBpedia annotators and the UMLS annotators as D1 and D2, and U1 and U2, respectively.

Figure 1 outlines our link discovery method. The method starts with a UMLS concept of *Stem Cell Factor* (C0143630). The first step is to obtain the concept definition from UMLS (“expressed during embryogenesis and provides key signal in multiple aspects of mast cell differentiation and function; hematopoietic growth factor and ligand of c-kit receptor CD117”). Next, we construct a query string with all known labels and aliases for this UMLS concept and concatenate it with the concept definition, as shown in Table 1 (left).

The query string is partitioned by a placeholder `DIV_DESCR`. This placeholder is used to divide the query string into 2 parts: labels with aliases (left-side) and UMLS definition (right-side). The right side is used by the semantic annotators to disambiguate the aliases on the left side of the placeholder. Similarly, the labels and aliases are kept separated from each other using a placeholder `DIV_NAME` to discourage semantic annotators from seeing incorrect multi-word n-grams by chance because of aliases situated next to each other. The generated query string is passed through 2 DBpedia semantic annotators (D1 and D2), each of which returns entity links to

DBpedia resources (Step 3). The DBpedia resources found to the left of the `DIV_DESCR` placeholder are collected as *link candidates*. For each of these link candidates, a new query is constructed, also shown in Table 1, but using the labels, aliases, and the abstract from DBpedia (Step 4). Each of these newly generated queries (from D1 and D2 link candidates) are then passed onto 2 UMLS semantic annotators (U1 and U2) in order to produce 4 UMLS annotated result sets: D1U1, D1U2, D2U1, and D2U2 (Step 5). Given these 4 result sets, we examine the UMLS annotations that appear to the left of the `DIV_DESCR` placeholder looking for an annotation with the CUI that we began with in Step 1 (i.e., C0143630). If such an annotation exists, then the candidate DBpedia resource is set aside to be later identified as either `skos:exactMatch` or `skos:closeMatch` (Step 6). Those candidates that do not produce the same CUI as the one used in Step 1 are delegated as having the weaker `rdfs:seeAlso` relationship.

In order to reduce disambiguation errors on the `rdfs:seeAlso` candidates, we discard those DBpedia resources that do not circularly resolve to *any* UMLS concepts in all 4 pairings of the annotators. In other words, all 4 pairings (D1U1, D1U2, D2U1, D2U2) must agree that the DBpedia resource resolves to some UMLS concept in order for the resource to remain as an `rdfs:seeAlso` relation.

Lastly, the `skos:exactMatch/closeMatch` set is separated into `skos:exactMatch` and `skos:closeMatch` relations by computing a *Jaccard coefficient* on all concept labels and aliases then testing for a minimum threshold. Formally, suppose UMLS concept CUI and DBpedia resource RES are related using *exact/close-match* as determined by our method (Figure 1). Let *C* and *T* be the set of all aliases/labels for CUI and RES, respectively. Let function *A*(*s*) return a set of individual characters from string *s*. Then, CUI is a `skos:exactMatch` to RES, if some label/alias of *C* and *T* meets the minimum threshold:

max\_{c \in C, t \in T} \frac{A(c) \cap A(t)}{A(c) \cup A(t)} \geq \text{Threshold} \tag{1}

We will show later in the paper that our method is not sensitive to specific threshold values.

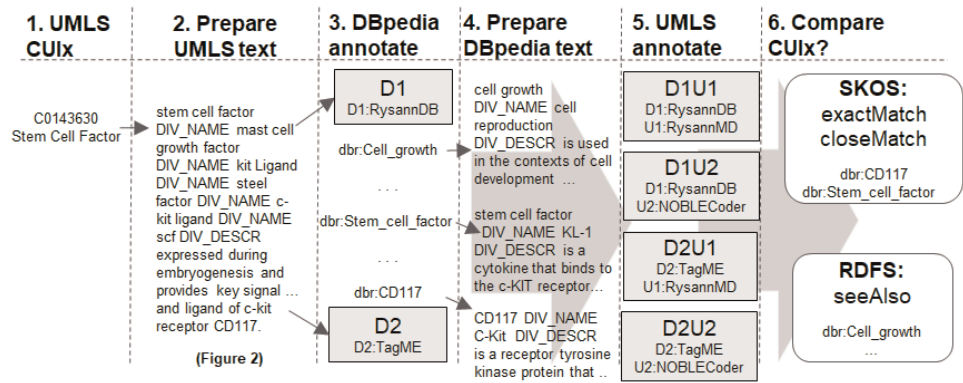


Figure 1. Pipeline for linking UMLS concepts to DBpedia using circular resolution method.

Table 1. Query String Constructed for the UMLS Concept

Stem Cell Factor C0143630	DBpedia resource Stem Cell Factor
Stem cell factor <b>DIV_NAME</b> mast cell growth factor <b>DIV_NAME</b> kit Ligand <b>DIV_NAME</b> steel factor <b>DIV_NAME</b> c-kit ligand <b>DIV_NAME</b> scf <b>DIV_DESCR</b> expressed during embryogenesis and provides key signal ... and ligand of c-kit receptor CD117.	Stem cell factor <b>DIV_NAME</b> steel factor <b>DIV_NAME</b> KITLG <b>DIV_NAME</b> KIT ligand <b>DIV_DESCR</b> Stem cell factor (also known as SCF, KIT-ligand, KL, or steel factor) is a cytokine that binds to the c-KIT receptor (CD117).

The bolded terms are added by our algorithm to partition them (force boundaries).

**rdfs:seeAlso —**  
 dbpedia:Cell\_growth  
 dbpedia:Coagulation  
 dbpedia:Chemical\_reaction  
 dbpedia:Stem\_cell  
 dbpedia:Ligand  
 dbpedia:Embryogenesis  
 dbpedia:Cell\_(biology)  
 dbpedia:Gene\_expression  
 dbpedia:Mast\_cell

**skos:closeMatch —**  
 dbpedia:SCF\_complex  
 dbpedia:CD117

**skos:exactMatch —**  
 dbpedia:Stem\_cell\_factor

**Figure 2.** The result of link discovery for *Stem Cell Factor* (C0143630) using the circular resolution method and Jaccard coefficient based [close|exact]-match classification.

We name the above method *circular resolution* given the fact that we begin with a UMLS concept (C0143630); annotate a query string (composed of label + aliases + definition) with DBpedia resources; construct a similar query string for each of the returned DBpedia resources; then annotate these DBpedia query strings using UMLS semantic annotators hoping to loop back to the original UMLS concept (C0143630). We complete the method with Equation 1 to produce the results as shown in Figure 2.

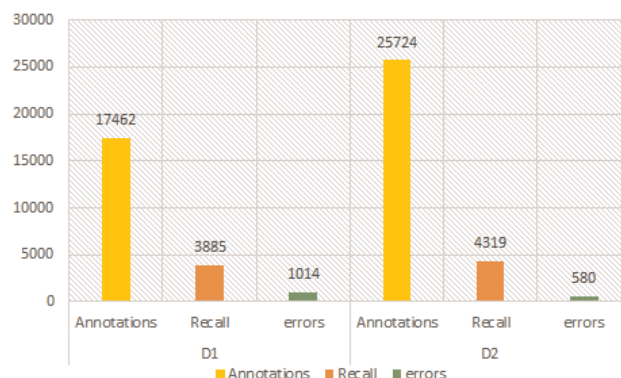
## Evaluation

To evaluate the effectiveness of our method, we queried Wikidata (<https://www.wikidata.org>) for all entries that have a UMLS mapping to DBpedia. The query returned 5006 entries. We disregarded mappings whose UMLS CUIs did not appear in our installation of UMLS because of licensing restrictions on the Metathesaurus. The final size of our testing set (ground truth) was 4899 entries.

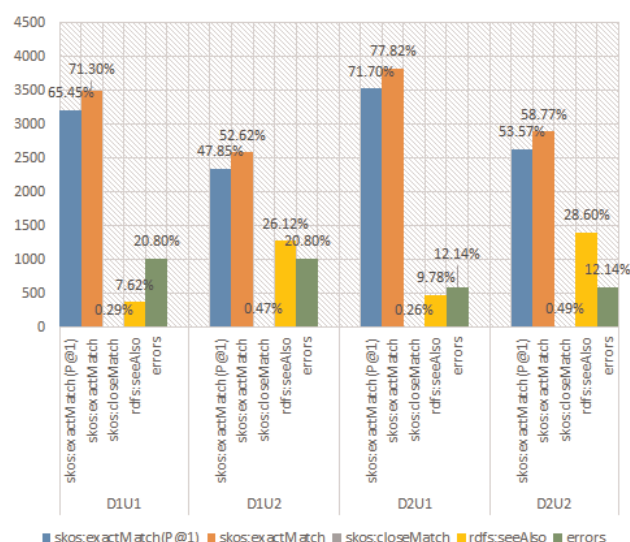
We performed an extensive search of the literature for reports on existing mappings of UMLS concepts to DBpedia that could serve as a benchmark for our algorithm and mapping. However, we found no such mapping, which suggests that our mapping is the first publicly available one. We also extensively searched for existing software tools that we could use to evaluate our algorithm and mapping. This proved quite difficult as we encountered numerous issues ranging from the lack of documentation to the unavailability of the systems themselves. Nonetheless, despite these issues, we were able to set up an additional baseline for comparison by using 2 annotators, i.e., RysannDB<sup>24</sup> and TagME,<sup>25</sup> which have been used for biomedical named entity recognition. We evaluate our cooperative circular resolution algorithm by comparing it against RysannDB and TagME annotators, using the Wikidata ground truth.

## RESULTS

We first focus our experiments on the output produced when only DBpedia annotators are used. In particular, we used RysannDB (<http://denote.rnet.ryerson.ca/RysannDB>)<sup>24</sup> as D1, and TagME (<https://tagme.d4science.org/tagme>)<sup>25</sup> as D2. Figure 3 shows why



**Figure 3.** Counts of links produced by RysannDB (D1) and TagME (D2) when annotating the Wikidata ground truth. Includes counts of matching (Recall) and nonmatching (errors) links.



**Figure 4.** Count and percentage of ground truth mappings resolved as *skos:exactMatch*, *skos:closeMatch*, *rdfs:seeAlso*, or neither (error) against Wikidata including precision-at-1 for each annotator pairing.

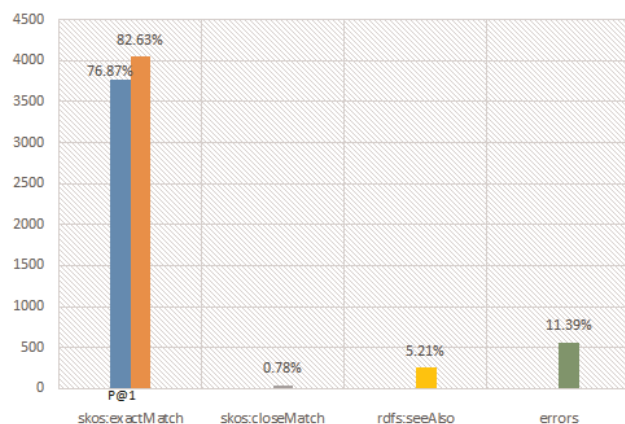
annotating UMLS definitions with DBpedia annotators alone would be ineffective.

RysannDB (D1) offered 17 462 entity links to DBpedia of which 3885 matched the Wikidata ground truth. TagME (D2) produced 25 724 links with 4319 matching links. Note that the Wikidata ground truth only contains 4899 entries. The matching counts of 3885 and 4319 measure the recall, whereas precision is negatively affected by the additional 13 577 and 21 405 links provided by the 2 annotators.

Next, we pair D1/D2 with UMLS annotators RysannMD (<http://denote.rnet.ryerson.ca/RysannMD>)<sup>23</sup> (U1) and Noble Coder (<http://noble-tools.dbmi.pitt.edu>)<sup>26</sup> (U2) to produce pairings of D1U1, D1U2, D2U1, and D2U2. Figure 4 shows how each pairing separately placed the ground truth into *skos:exactMatch*, *skos:closeMatch*, *rdfs:seeAlso*, or neither (disambiguation or recall error) using circular resolution.

From among the 4 pairings, the pairing of TagME and RysannMD (D2U1) was the most effective at linking UMLS to DBpedia with a 77.82% recall in identifying ground truth mappings as the expected *skos:exactMatch* relationship type. This pairing also achieved the smallest number of errors at 12.14%. The next best per-





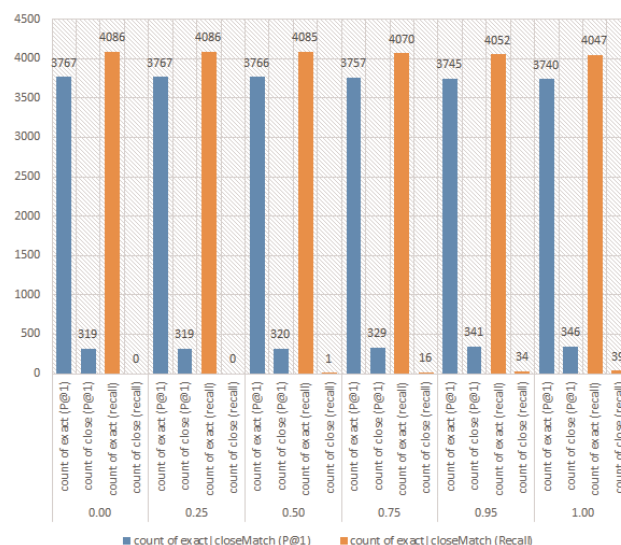
**Figure 5.** Count of ground truth mappings resolved as *skos:exactMatch*, *skos:closeMatch*, *rdfs:seeAlso*, or neither (error) including precision-at-1.

forming pair based on RysannDB and RysannMD (D1U1) achieved a recall of 71.30% with an error of 20.80%. Although the pairings of D1U2 and D2U2 performed weaker with a 52.62% and 58.77% recall agreement, it will be shown in Figure 5 that collectively they contribute to producing a better result. This is because each individual pairing provides some unique mappings that the others do not.

We include in our analysis a precision-at-one (P@1) metric on the *skos:exactMatch* type to better judge the effectiveness of our circular resolution method. Specifically, the Wikidata ground truth assumes a 1-to-1 *skos:exactMatch* mapping between a UMLS Concept and a DBpedia resource. However, our technique may return multiple *skos:exactMatch* links for a single UMLS concept. Consequently, we report on the method's performance when a 1-to-1 mapping is strictly required by selecting the resource with the highest Jaccard coefficient that also meets the minimum threshold (Equation 1). We found our aforementioned top pairings of D1U1 and D2U1 still bested D1U2 and D2U2 with a precision-at-one of 65.45% and 71.70%.

Next, as per our pipeline (Figure 1), we pool together the mappings of each of the 4 pairings as a single solution, then report on *skos:[exact|close]Match*, *rdfs:seeAlso*, errors, and P@1 in Figure 5. Our findings show this combined mapping performs best in exactMatch (recall), errors, and exactMatch (P@1) than any of the individual pairings.

We conclude our tests by examining the sensitivity of our approach to the threshold for the Jaccard coefficient introduced in Equation 1. We show how the threshold affects precision and recall on *skos:exactMatch* classifications when a 1-to-1 UMLS concept to DBpedia resource mapping is required (i.e., high precision P@1), and also when multiple DBpedia resources are allowed to link to a single UMLS concept, i.e., high recall. As shown in Figure 6, when the threshold value is set to zero, we observed 3767 correctly mapped concepts at P@1 vs 4086 correctly linked when a 1-to-many mapping is allowed. There was no change when the threshold was set to 0.25 and a negligible change of 1 exact-match to a close-match reclassification at a threshold of 0.5. Changes occurred when the threshold was set to 0.95 when a difference of 22 and 34 exact matches were observed. Furthermore, when the threshold was set to one, 27 and 39 exact match relationship changes were observed. The impact of varying the threshold from 0 to 1 results in an overall performance change of around 0.5%; hence showing insensitivity to the threshold. From these results, we can conclude that the 4 annotators (D1/D2/U1/U2) are effectively leveraging their semantic capabilities to provide high quality candidates for close/exact-match



**Figure 6.** Counts on the number of exact/close matches with a Jaccard threshold of 0 (no threshold), 0.25, 0.50, 0.75, 0.95, and 1.00.

determination, and thus our method is relatively stable with respect to any chosen Jaccard threshold. Consequently, the best configuration would be utilizing Equation 1 to solely rank the candidates (using value zero as the threshold) then selecting the highest computed Jaccard for a 1-to-1 exact match (i.e., P@1).

## DISCUSSION

### Ground Truth Error Analysis

Our method achieved noticeable recall (83%) and precision scores (77%) during experimentation using the Wikidata ground truth benchmark. However, it did make mistakes depending on how an expected exact match concept was classified at various stages of the pipeline. We classify these errors as follows:

1. *Candidate Selection Omission.* The DBpedia annotators (D1/D2) did not select the correct resource as a candidate. The outcome is that the correct resource does not appear as an exact match, close match, or see also.
2. *Failed to Promote Error.* The UMLS annotators (U1/U2) did not produce any links from the candidate resource back to the target UMLS concept. In this case, the correct resource remains as an *rdfs:seeAlso*.
3. *Failed to Meet Threshold.* Annotators U1/U2 correctly promoted a resource as a *skos:closeMatch* or *skos:exactMatch* but the correct resource either failed to meet the threshold in Equation 1 or a higher calculated Jaccard for the wrong concept was selected for P@1. This results in the correct resource being classified as *skos:closeMatch*.
4. *Wrongly Promoted and Failed Jaccard Filtering.* A wrong concept was incorrectly promoted by U1/U2 and satisfied the Jaccard threshold or best P@1. This leads to linking the UMLS concept to an incorrect DBpedia resource as a *skos:exactMatch* (disambiguation error).

Table 2 summarizes the counts of the errors encountered during ground truth testing with 10 examples for each error type. For example, CUI concept C0001815 “Primary Myelofibrosis” failed as

**Table 2.** Summary of Circular Resolution Error Counts (types 1–4) with Showcase Examples of Expected Ground Truth (G.T) and Circular Resolution (C.R) Answer

(1) Candidate Selection Omission	(2) Failed to Promote Error	(3) Failed to meet Threshold	(4) Wrongly Promoted + Jaccard Filter
388 (No Link) Sample CUIs C0496758 C0302182 C2937300 C0153620 C0022441 C0023234 C0025534 C0477373 C0795950 C1841679 Showcase example CUI: C0477373 “Other forms of migraine” G.T: Familial_hemiplegic_migraine C.R.: no entity link	255 (rdfs:seeAlso) C2931205 C0006111 C0008684 C0155937 C1854540 C2607929 C1412004 C1335473 C1514284 C0279607 CUI: C1514284 “Potassium Deficiency Disorder” G.T: Hypokalemia C.R: Linked as rdfs:seeAlso	38 (skos:closeMatch) C0031946 C0153241 C0041341 C1274184 C0019284 C0018553 C0266611 C0001815 C0007134 C0343065 CUI: C0001815 “Primary Myelofibrosis” G.T: Myelofibrosis C.R: CIMF-FM (exact) Myelofibrosis (close)	170 (error) C0751782 C0039753 C0020433 C0795690 C0741160 C0026697 C0917990 C0032290 C0072826 C1337224

**Table 3.** Two Examples of *rdfs:seeAlso* Mappings Where No Exact Match is Available

C3175196 “Other people frequently tell me that what I’ve said is impolite even though I think it is polite: d:Pt: ^Patient: Ord: PhenX” <i>rdfs:seeAlso</i> — dbpedia:Taboo dbpedia:Time dbpedia:Patient dbpedia:Thought dbpedia:Level_of_measurement	C0370538 “Punch graft for hair transplant; more than 15 punch grafts” <i>rdfs:seeAlso</i> — dbpedia:Bone_grafting dbpedia:Hair dbpedia:Organ_transplantation dbpedia:Hair_transplantation dbpedia:Graft_(surgery)
---	---

a type (3) error resulting in a close match classification. This same CUI also suffered a type (4) error as it was wrongly linked to DBpedia resource *CIMF-FM*. The reader is encouraged to use our online browser (<http://denote.rnet.ryerson.ca/umlsMap/browser>) to further investigate each of these errors.

## Mapping the UMLS

We applied our method to the UMLS Metathesaurus to produce 221 690 *skos:exactMatch*, 26 276 *skos:closeMatch*, and 6 784 322 *rdfs:seeAlso* relations. The total number of concepts in our license-free version of the UMLS was 2 397 167. This gives a percentage of mapping from the UMLS to DBpedia for *skos:[close|exact]Match* of 10.34% and an average of 2.83 *rdfs:seeAlso* relationships per concept. Although this may seem a low percentage, consider that our ground truth from all of Wikidata contained only 5006 mapped UMLS concepts compared to our 221 690 mappings (a factor of 50x increase). The difficulty in mapping a large portion of the UMLS as an exact match occurs largely because many concepts are so specific as to not have a corresponding entry in DBpedia, as illustrated in Table 3. This is not very surprising to those familiar with UMLS. In order to gain further insight, we performed a simple experiment in which we surmised that the one-word concepts in UMLS were more likely to have a corresponding exact match DBpedia entry than those comprising 2 or more words. To further challenge our method, we excluded those one-word concepts that appeared directly within the DBpedia URL itself thus making it more difficult for the annotators to perform the alignment (e.g., *C0018081:Gonor-*

*rhea* mapped to *dbpedia:Gonorrhea* was excluded from this experiment). A cursory inspection of a random sampling of the 241 791-word mappings revealed good results with success and error rates equivalent to those observed in Figure 5 and Table 2. For example, our method correctly mapped *C0001429:Adenolymphoma* with *dbpedia:Warthin’s\_tumor*, but mistakenly matched *C1174791:Basen* to *dbpedia:Basen,\_Armenia*. We have provided this one-word mapping as a supplementary document for further inspection.

## Maintaining the UMLS Mapping

From the perspective of the choice of the semantic annotators, RysannDB (D1) and TagME (D2) were selected as the DBpedia linkers because of their accuracy and speed of processing natural language text. Speed is a particular concern since our goal was to map the entire UMLS to DBpedia. Some other well-known annotators, although of comparable accuracy, are too slow to be practical for this task. The same consideration was given to the choice of RysannMD (U1) and Noble Coder (U2) based on the findings in.<sup>24</sup>

The time to map UMLS to DBpedia required ~60 h of processing for each pairing (D1U1, D1U2, D2U1, D2U2) on an Intel 3.00 GHz Xeon CPU-based server with 128GB of RAM. Although this may seem time intensive, one should consider the following:

1. Our implementation of circular resolution was focused on link discovery challenges, not on processing time optimization. Efficiency-oriented implementations would execute the processing of pairs D1/D2 and U1/U2 concurrently, thus reducing the

mapping time by a factor of 4. Further improvements can be gained by dividing the UMLS database into smaller datastores and processing in parallel.

2. Updates of the mapping require the processing of only new UMLS entries allowing for incremental updates.
3. Like other open datasets, the burden of (1) and (2) falls to the authors of this work as the dataset maintainers. We intend to maintain this dataset and make it available through our website and officially through the LOD cloud.

## Alternative Approaches

Link discovery and instance matching is an active area of research, with many open challenges. A comprehensive survey by Nentwig et al.<sup>23</sup> gives a good summary of the current state-of-the-art. In this survey, 9 out of 11 examined frameworks could only determine owl:sameAs relationships. The remaining frameworks (Silk<sup>15</sup> and LIMES<sup>16</sup>), do support additional link types through heuristic rules. However, the user is responsible for manually constructing the necessary heuristic patterns for detecting a particular relationship type, e.g., rdfs:seeAlso. In contrast, our method operates at a higher level of abstraction relying on underlying semantic annotation engines. This allows our method to easily take advantage of a wide combination of techniques that have already been incorporated into existing semantic annotators by choosing different annotators to fill in the role of D1, D2, U1, and U2. Furthermore, the heuristic rules approach taken by Silk and LIMES may not be interchangeable between different pairs of KBs. That is, rules designed to map from KB1 to KB2 may not be the same rules needed to map from KB1 to KB3 even for the same link type. Comparatively, our method performs the alignment by only considering textual information from readily available concept labels/definitions and through the use of the natural language processing capabilities of existing semantic annotators. It should be noted, however, that the heuristic rules approach undertaken by Silk and LIMES does allow for flexibility in the relationship type sought after, whereas our method is limited to skos:[exact|close]Match and rdfs:seeAlso; an important point in the conceptual distinction between Silk and LIMES and our work. Both Silk and LIMES are customizable and extensible frameworks on top of which specific link discovery processes are implemented to interconnect different datasets. Both of these frameworks are primarily developed to allow experts to design mapping pipelines from existing components that are shipped with the 2 frameworks or can be added to the frameworks as third party add-ons. However, our work focuses on one specific mapping process and, hence, would not be considered as an extensible framework. In this light, circular resolution could be integrated into the LIMES or Silk pipeline that could prove valuable for a wider range of mapping tasks.

Lastly, we considered numerous designs for circular resolution before settling on the method proposed here. One such consideration involved the treatment of the primary label and alternative names of a concept as separate annotation problems, which would then be merged. This approach would have eliminated the use of the separation tokens, i.e., DIV\_NAME and DIV\_DESCR. Details of this alternative method, and the reason for its dismissal, are given in a supplementary document (Supplementary Appendix A).

## CONCLUSION

In this paper, we have presented a method, called circular resolution, to map UMLS concepts to DBpedia resources using rdfs:seeAlso, sko-

s:closeMatch, and skos:exactMatch relations. Our technique reports a recall of 83% with 77% precision-at-one when benchmarked against Wikidata. A full UMLS to DBpedia mapping is also made publicly available. In addition, we provide an online browser to easily explore the mappings and a RESTful interface for querying the mappings (<http://denote.rnet.ryerson.ca/umlsMap>). We hope that this mapping can become an integral part of the *Linked Open Data cloud* and facilitate the effective interchange and integration of different knowledge bases with medical and biomedical knowledge bases. To this end, our future work includes creating UMLS mappings for the various ontologies openly available through “*The Open Biological and Biomedical Ontology (OBO) Foundry*” (<http://www.obofoundry.org/>) which provides open access to medical and biological vocabularies.

## FUNDING

This work was supported by the *Natural Sciences and Engineering Research Council of Canada* under grant number RGPIN-2015-06118

## COMPETING INTERESTS

None.

## CONTRIBUTORS

The authors declare that this manuscript is a product of original work and each author contributed to the design and interpretation of the results. Furthermore, the authors have critically evaluated the content before final approval for publication. The authors are accountable for all aspects of this work and believe in the accuracy of their results, interpretation thereof, and content of this manuscript.

## SUPPLEMENTARY MATERIAL

Supplementary material is available at *Journal of the American Medical Informatics Association* online.

## REFERENCES

1. Lehmann J, Isele R, Jakob M, et al. DBpedia - a large-scale, multilingual knowledge base extracted from Wikipedia. *Semantic Web J* 2012; 6 (2): 167–195.
2. Heath T, Bizer C. Linked data. *Evolving the Web into a Global Data Space (1st edition)*. *Synthesis Lectures on the Semantic Web: Theory and Technology*. San Rafael, CA, USA: Morgan & Claypool; 2011: 1–136.
3. Pai VM, Rodgers M, Conroy R, et al. Workshop on using natural language processing applications for enhancing clinical decision making: an executive summary. *J Am Med Inform Assoc* 2014; 21 (e1): e2–e5.
4. Friedlin J, McDonald, CJ. An evaluation of medical knowledge contained in Wikipedia and its use in the LOINC database. *J Am Med Inform Assoc* 2010; 17 (3): 283–287.
5. Dessi N, Atzori M. Is Wikipedia a Latent Gene Ontology? In *2017 IEEE 26th International Conference on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE)*. 2017: 164–169.
6. Yamamoto Y, Yamaguchi A, Yonezawa A. Building linked open data towards integration of biomedical scientific literature with DBpedia. *J Biomed Semant* 2013; 4: 8.
7. Laurent MR, Vickers TJ. Seeking health information online: does Wikipedia matter? *J Am Med Inform Assoc* 2009; 16 (4): 471–479.
8. Wiesner M, Pfeifer D. Health recommender systems: concepts, requirements, technical basics and challenges. *Int J Environ Res Public Health* 2014; 11 (3): 2580–2607.

9. Park M, Sampathkumar H, Luo B, Chen XW. Content-based assessment of the credibility of online healthcare information. In *2013 IEEE International Conference on Big Data*. 2013: 51–58.
10. Keselman A, Smith CA, Divita G, *et al*. Consumer health concepts that do not map to the UMLS: where do they fit? *J Am Med Inform Assoc* 2008; 15 (4): 496–505.
11. Roberts K, Demner-Fushman D. Interactive use of online health resources: a comparison of consumer and professional questions. *J Am Med Inform Assoc* 2016; 23 (4): 802–811.
12. Mrabet Y, Kilicoglu H, Roberts K, Demner-Fushman D. Combining open-domain and biomedical knowledge for topic recognition in consumer health questions. *AMIA Ann Symp Proc* 2016; 914–923.
13. Ernst P, Siu A, Milchevski D, Hoffart J, Weikum G. DeepLife: An entity-aware search, analytics and exploration platform for health and life sciences. In Pradhan S, Apidianaki M, eds. *Proceedings of ACL-2016 System Demonstrations* Stoudsborg, PA: ACL; 2016: 19–24.
14. Ernst P, Siu A, Weikum G. Knowlife: a versatile approach for constructing a large knowledge graph for biomedical sciences. *BMC Bioinformatics* 2015; 16 (1): 1–13.
15. Volz J, Bizer C, Gaedke M, Kobilarov G. Silk – a link discovery framework for the web of data. *Proceedings of the 2nd Linked Data on the Web Workshop*. 2009: 559–572.
16. Ngomo A, Auer S. LIMEs: a time-efficient approach for large-scale link discovery on the web of data. In *proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence (IJCAI'11)*. 2011; 2312–2317.
17. Tilahun B, Kauppinen T, Keßler C, Fritz F. Design and development of a linked open data-based health information representation and visualization system: potentials and preliminary evaluation. *JMIR Med Inform* 2014; 2 (2): e31.
18. Bing H, Tang J, Ding Y, *et al*. Mining relational paths in integrated biomedical data. *PLoS One* 2011; 6 (12): e27506.
19. Luciano J, Andersson B, Batchelor C, *et al*. The translational medicine ontology and knowledge base: driving personalized medicine by bridging the gap between bench and bedside. *J Biomed Semantics* 2011; 2 (Suppl 2): S1.
20. Saleem M, Padmanabhuni S, Ngomo A, Almeida J, Decker S, Deus H. Linked cancer genome atlas database. In *Proceedings of the 9th International Conference on Semantic Systems. I-SEMANTICS '13*. 2013; 129–134.
21. Lee LH, Groß A, Hartung M, Liou DM, Rahm E. A multi-part matching strategy for mapping LOINC with laboratory terminologies. *J Am Med Inform Assoc* 2014; 21 (5): 792–800.
22. Kahn CE. Integrating ontologies of rare diseases and radiological diagnosis. *J Am Med Inform Assoc* 2015; 22 (6): 1164–1168.
23. Nentwig M, Hartung M, Ngonga Ngomo AC, Rahm E. A survey of current Link Discovery frameworks. *Semantic Web* 2017; 8 (3): 419–436.
24. Cuzzola J, Jovanovic J, Bagheri E. RysannMD: a biomedical semantic annotator balancing speed and accuracy. *J Biomed Inform* 2017; 71: 91–109.
25. Ferragina P, Scaiella U. TAGME: on-the-fly annotation of short text fragments (by Wikipedia entities). In *Proceedings of the 19th ACM International Conference on Information and Knowledge Management (CIKM '10)*. New York, NY, USA: ACM; 2010: 1625–1628.
26. Tseytlin E, Mitchell K, Legowski E, Corrigan J, Chavan G, Jacobson RS. NOBLE - Flexible concept recognition for large-scale biomedical natural language processing. *BMC Bioinformatics* 2016; 17: 32.