

# Using logical constraints to validate information in collaborative knowledge graphs: a study of COVID-19 on Wikidata

Houcemeddine Turki<sup>a</sup> (ORCID: 0000-0003-3492-2014), Dariusz Jemielniak<sup>b</sup> (ORCID: 0000-0002-3745-7931), Mohamed Ali Hadj Taieb<sup>c</sup> (ORCID: 0000-0002-2786-8913), Jose Emilio Labra Gayo<sup>d</sup> (ORCID: 0000-0001-8907-5348), Mohamed Ben Aouicha<sup>c</sup> (ORCID: 0000-0002-2277-5814), Mus'ab Banat<sup>e</sup> (ORCID: 0000-0001-9132-3849), Thomas Shafee<sup>f</sup> (ORCID: 0000-0002-2298-7593), Eric Prud'Hommeaux<sup>g</sup> (ORCID: 0000-0003-1775-9921), Tiago Lubiana<sup>h</sup> (ORCID: 0000-0003-2473-2313), Diptanshu Das<sup>i</sup> (ORCID: 0000-0002-7221-5022), Daniel Mietchen<sup>j,1\*</sup> (ORCID: 0000-0001-9488-1870), on behalf of WikiProject COVID-19<sup>2</sup>

<sup>a</sup> *Faculty of Medicine of Sfax, University of Sfax, Sfax, Tunisia*

<sup>b</sup> *Department of Management in Networked and Digital Societies, Kozminski University, Warsaw, Poland*

<sup>c</sup> *Faculty of Sciences of Sfax, University of Sfax, Sfax, Tunisia*

<sup>d</sup> *Web Semantics Oviedo (WESO) Research Group, University of Oviedo, Spain*

<sup>e</sup> *Faculty of Medicine, Hashemite University, Zarqa, Jordan*

<sup>f</sup> *La Trobe University, Melbourne, Victoria, Australia*

<sup>h</sup> *Computational Systems Biology Laboratory, University of São Paulo, São Paulo, Brazil*

<sup>i</sup> *Institute of Child Health (ICH), Kolkata, India*

<sup>i</sup> *Medica Superspecialty Hospital, Kolkata, India*

<sup>j</sup> *School of Data Science, University of Virginia, Charlottesville, Virginia, United States of America*

**Abstract.** Urgent global research demands real-time dissemination of precise data. Wikidata, a collaborative and openly licensed knowledge graph available in RDF format, provides an ideal forum for exchanging structured data that can be verified and consolidated using validation schemas and bot edits. In this research paper, we catalog an automatable task set necessary to assess and validate the portion of Wikidata relating to the COVID-19 disease, its causative virus, and key aspects of the resulting pandemic. These tasks assess relational and statistical data and are implemented in SPARQL, a query language for semantic databases. We demonstrate the efficiency of our methods for evaluating structured information on COVID-19 in Wikidata, and its applicability in collaborative ontologies and knowledge graphs more broadly. We show the advantages and limitations of our proposed approach by comparing it to other methods for validation of linked web data.

**Keywords:** SPARQL, Public health surveillance, Wikidata, Knowledge graph refinement, COVID-19, Validation constraints

---

<sup>1\*</sup> Corresponding author. E-mail: dm7gn@virginia.edu.

<sup>2</sup> Project Member: Project members: Jan Ainali, Susanna Ånäs, Erica Azzellini, Mus'ab Banat, Mohamed Ben Aouicha, Alessandra Boccone, Jane Darnell, Diptanshu Das, Lena Denis, Rich Farmbrough, Daniel Fernández-Álvarez, Konrad Foerstner, Jose Emilio Labra Gayo, Mauricio V. Genta, Mohamed Ali Hadj Taieb, James Hare, Alejandro González Hevia, David Hicks, Toby Hudson, Netha Hussain, Jinoy Tom Jacob, Dariusz Jemielniak, Krupal Kasyap, Will Kent, Samuel Klein, Jasper J. Koehorst, Martina Kutmon, Antoine Logean, Tiago Lubiana, Andy Mabbett, Kimberli Mäkräinen, Tania Maio, Bodhisattwa Mandal, Nandhini Meenakshi, Daniel Mietchen, Nandana Mihindukulasooriya, Mahir Morshed, Peter Murray-Rust, Minh Nguyễn, Finn Årup Nielsen, Mike Nolan, Shay Nowick, Julian Leonardo Paez, João Alexandre Peschanski, Alexander Pico, Lane Rasberry, Mairélys Lemus-Rojas, Diego Saez-Trumper, Magnus Säljö, John Samuel, Peter J. Schaap, Jodi Schneider, Thomas Shafee, Nick Sheppard, Adam Shorland, Ranjith Siji, Michal Josef Špaček, Ralf Stephan, Andrew I. Su, Hilary Thorsen, Houcemeddine Turki, Lisa M. Verhagen, Denny Vrandečić, Andra Waagmeester, and Egon Willighagen.

## 1. Introduction

Since December 2019, the COVID-19 disease has spread to become a global pandemic. This disease is caused by a zoonotic coronavirus called SARS-CoV-2 (Severe Acute Respiratory Syndrome CoronaVirus 2) and is characterized by the onset of acute pneumonia and respiratory distress. The global impact, with more than 77 million infections and almost 1.7 million deaths globally (as of December 21, 2020<sup>3</sup>), is frequently compared to the 1918 Spanish Flu [1]. Emerging mRNA vaccines entail serious distribution and storage challenges and no therapies are especially effective against late-stages of the disease. As with all zoonotic diseases, its abrupt introduction to humans demands an outsized effort for data acquisition, curation and integration to drive evidence-based medicine, predictive modeling and public health policy [2, 3].

Agile data sharing and computer-supported reasoning about the COVID-19 pandemic and SARS-CoV-2 virus allow us to quickly understand more about the disease's epidemiology, pathogenesis, and physiopathology. This understanding can then inform the required clinical, scholarly and public health measures to fight the condition and handle its non-medical ramifications [4-6]. Consequently, initiatives have rapidly emerged to create datasets, web services and tools to analyse and visualise COVID-19 data. Examples include Johns Hopkins University's COVID-19 dashboard [2] and the Open COVID-19 Data Curation Group's epidemiological data [3]. Some of these resources are interactive and return their results based on combined clinical and epidemiological information, scholarly information and social network analysis [7-9]. However, a significant shortfall in interoperability is common: although these dashboards facilitate examination of their own slice of the data, most lack general integration with other sites or datasets. The lack of technical support for interoperability is exacerbated by legal restrictions: despite being free to access, most are issued under *All Rights Reserved* terms or licenses. Similarly, >80% of the 96608 COVID-19-related projects on the GitHub repository for

computing projects are under *All Rights Reserved*<sup>4</sup> terms (as of 21 December 2020). Restrictive licensing of data sets and applications severely impedes their dissemination and integration, ultimately undermining their value. For complex and multifaceted phenomena such as the COVID-19 pandemic, there is a particular need for a collaborative, free, machine-readable, interoperable and open knowledge graph to integrate the varied data.

Wikidata<sup>5</sup> just fits the need as a CC0<sup>6</sup> licensed, large-scale, multilingual knowledge graph used to represent human knowledge in a structured format (Resource Description Framework or RDF) [10, 11]. It therefore has the advantage of being inherently findable, accessible, interoperable, and reusable, i.e. FAIR [12]. It was initially developed in 2012 as an adjunct to Wikipedia but has grown significantly beyond its initial parameters. As of now, it is a centralized, cross-disciplinary meta-database and knowledge base for storing structured information in a format optimized to be easily read and edited by both machines and humans [13]. Thanks to its flexible representation of facts, Wikidata can be automatically enriched using information retrieved from multiple public domain sources or inferred from synthesised data [11]. This database includes a wide variety of pandemic-related information, including clinical knowledge, epidemiology, biomedical research, software development, geographic, demographic and genetics data. It can consequently be a vital large-scale reference database to support research and medicine during the COVID-19 pandemic [11, 12].

The key hurdle to overcome for projects such as Wikidata is that several of their features can make them at-risk of inconsistent structure or coverage: 1) collaborative projects use decentralised contribution rather than central oversight, 2) large-scale projects

<sup>3</sup> "COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU)". ArcGIS. Johns Hopkins University. Retrieved 25 August 2020.

<sup>4</sup> 80002 of 96608 as of 2020-12-21: <https://github.com/search?q=covid-19+OR+covid19+OR+coronavirus+OR+cord19+OR+cord-19>

<sup>5</sup> <https://www.wikidata.org/>

<sup>6</sup> CC0 is a rights waiver similar to Creative Commons licenses, used to publish material into the public domain. It waives as much copyright as possible within a given jurisdiction. Further information can be found at <https://creativecommons.org/publicdomain/zero/1.0/>

operate at a scale where manual checking is not possible, and 3) interdisciplinary projects script the acquisition of data to integrate a wide variety of data sources. To maximise usability of the data, it is therefore important to minimise inconsistencies in its structure and coverage. As a result, methods of evaluating the existing knowledge graphs and ontologies, integral to knowledge graph maintenance and development, are of crucial importance. Such an evaluation is particularly relevant in the case of collaborative semantic databases, such as Wikidata.

Knowledge graph evaluation is therefore necessary to assess the quality, correctness, or completeness of a given knowledge graph against a set of predetermined criteria [14]. There are a number of possible approaches to evaluating a knowledge graph based on external information (so-called extrinsic evaluation), including: comparing its structure to a paragon ontology, comparing its coverage to source data, applying it to a test problem and judging the outcomes, and manual expert review of its ontology [15]. Different systematic approaches have been proposed for the comparison of ontologies and knowledge graphs, including NLP techniques, machine learning, association rule mining, and other methods [16-18]. The criteria for evaluating ontologies typically include: *Accuracy*, which determines if definitions, classes, properties and individual entries in the evaluated ontology are correct; *Completeness*, referring to the scope of coverage of a given knowledge domain in the evaluated ontology; *Adaptability*, determining the range of different anticipated uses of the evaluated ontology (versatility); and *Clarity*, determining the effectiveness of communication of intended meanings of defined terms by the evaluated ontology [14, 19-21]. However, extrinsic methods are not the only ones that are used for evaluating such a set of criteria. Knowledge graphs can be also assessed through an intrinsic evaluation that assesses the structure of the analyzed knowledge graph thanks to the inference of internal description logics and consistency rules [14].

In this research paper, we emphasize the usefulness of intrinsic methods to evaluate knowledge graphs by presenting our solution to the quality assurance checks and corrections of COVID-19 semantic data in Wikidata. This consists of a catalogue of automatable tasks based on logical constraints expected of the knowledge graph. Most of

these constraints were not explicitly available in the RDF validation resources of Wikidata before the pandemic and are designed in this work to support new types of COVID-19 information in the assessed knowledge graph, including epidemiological and social data. We implement these constraints with SPARQL and test them on Wikidata using the SPARQL endpoint of this knowledge graph, available at <https://query.wikidata.org>. We introduce the value of Wikidata as a multi-purpose collaborative knowledge graph for the flexible and reliable representation (Section 2) and validation (Section 3) of COVID-19 knowledge. Furthermore, we cover the use of SPARQL to query this knowledge graph (Section 4). Then, we demonstrate how logical constraints can be captured in structural schemas and consequently used to validate and encourage the consistent usage of relation types to represent COVID-19 knowledge (Section 5) and we show how statistical constraints can be applied to verify epidemiological data related to the pandemic (Section 6). Finally, we compare our constraint-based approaches with other methods through the analysis of the outcomes of previous research papers related to knowledge graph validation (Section 7), and draw conclusions for future directions (Section 8).

## 2. Wikidata as a collaborative knowledge graph

Wikidata currently serves as a semantic framework for a variety of scientific initiatives, such as GeneWiki [22], allowing different teams of scholars to upload valuable academic data into a collective and standardized pool. Its versatility and interconnectedness are making it a standard for interdisciplinary data integration and dissemination across fields as diverse as linguistics, information technology, film studies, and medicine [11, 23-28], although its popularity and recognition across fields still vary significantly [29].

It contains concepts, linked by their taxonomic relations, allowing embedding and creating instances of subclasses of classified data and links between them. Its multilingual nature enables fast-updating dynamic data reuse across different language versions of a resource such as Wikipedia [30], with fewer inconsistencies from local culture [31] or language biases [32, 33].

The data structure employed by Wikidata is intended to be highly standardized, whilst

maintaining the flexibility to be applied across highly diverse use-cases. There are mainly two essential components: Items, which represent objects, concepts or topics; and properties, which describe how one item relates to another. A statement, therefore, consists of a subject item (*S*), a property that describes their nature of the statement (*P*), and an object (*O*) that can be an item, a value, an external ID, or a string, etc. While items can be freely created, new properties require community discussion and vote, with 7851 properties<sup>7</sup> currently available. Statements can be further modified by any number of qualifiers to make them more specific and be supported by references to indicate the source of the information.

Thus, Wikidata forms a continuously growing, single, unified network graph, with 88M items forming the nodes, and 1127M statements<sup>8</sup> forming the edges. A live SPARQL endpoint and query service, regular RDF dumps, as well as linked data APIs and visualization tools, form a backbone of Wikidata uses [34, 35].

Importantly, Wikidata is based on free and open-source philosophy and software and is a database that anyone can edit, similarly to the very popular online encyclopedia, Wikipedia [36]. As a result, the emerging ontologies are created entirely collaboratively, without centralized coordination [37], and developed in a community-driven fashion [38]. This approach allows for the dynamic development of areas of interest for the user community but poses challenges, e.g., to systematize and proportionate class completeness across topics [39]. Also, since the edit history is available to anyone, tracing human and non-human contributions, as well as detecting and reverting vandalism is available by design and relies on community management [40] as well as on software tools like ORES [41].

Other ontological databases and knowledge graphs exist [42, 43]. However, much like the factors that led Wikipedia to rise to be a dominant encyclopedia [33, 44], Wikidata's close connection to Wikimedia volunteer communities and wide readership provided by Wikipedia have quickly given it a competitive edge. The system, therefore, aims to combine the wisdom of the crowds with advanced algorithms. For instance, Wikidata editors are assisted by a property

suggesting system, proposing additional properties to be added to entries [45]. Wikidata has subsequently exhibited the highest growth rate of any Wikimedia project and was the first amongst them to pass one billion contributions [46].

As a collaborative venture, its governance model is similar to Wikipedia [47], but with some important differences. Wide permissions to edit Wikidata are manually granted to approved bots and to Wikimedia accounts that are at least 4 days old and have made at least 50 edits using manual modifications or semi-automated tools for editing Wikidata<sup>9</sup>. These accounts are supervised by a limited number of experienced administrators to prevent misleading editing behaviors (such as vandalism, harassment, and abuse) and to ensure a sustainable consistency of the information provided by Wikidata<sup>10</sup>. As such, Wikidata is highly relevant to the computer-supported collaborative work (CSCW) field, yet the number of studies of Wikidata from this perspective is still very limited [48]. To understand the value of using SPARQL to validate the usage of relation types in collaborative ontologies and knowledge graphs, it is important to understand the main distinctive features of Wikidata as a collaborative project.

Much as Wikidata is developed collaboratively by international editors, it is also designed to be language-neutral. As a result, it is quite possible to contribute to Wikidata with only a limited command of English and to effectively collaborate whilst sharing no common human language - an aspect unique even in the already rich ecosystem of collaborative projects [49]. It may well be an early sign of other language-independent cooperative knowledge creation initiatives, such as Wikilambda, which is an abstract Wikipedia currently developed on the basis of Wikidata [50].

It is also possible to build Wikipedia articles, especially in underrepresented languages, based on Wikidata data only, and create article placeholders to stimulate encyclopedia articles' growth [51]. This stems from combining concepts that are relatively easily intertranslatable between languages (e.g. professions, causes of death, capitals) with language-agnostic data (e.g. numbers, geographical coordinates, dates). As a result, Wikidata is a paragon example of not only cross-cultural cooperation but

<sup>7</sup> For an updated list of available Wikidata property, please see <https://tools.wmflabs.org/hay/propbrowse/>.

<sup>8</sup> To track the evolution of the number of Wikidata statements, please see <https://grafana.wikimedia.org/d/000000182/wikidata-datamodel-references?orgId=1&refresh=30m>.

<sup>9</sup> For an overview of the semi-automated editing tools for Wikidata, please see <https://www.wikidata.org/wiki/Wikidata:Tools>

<sup>10</sup> Further information about the rights and governance of users in Wikidata is shown at [https://www.wikidata.org/wiki/Wikidata:User\\_access\\_levels](https://www.wikidata.org/wiki/Wikidata:User_access_levels)

also human-bot collaborative efforts [37, 52]. Given the large-scale crowdsourcing efforts in Wikidata and the use of bots and semi-automated tools to mass edit Wikidata, its current volume is higher than what can be reviewed and curated by administrators manually. It is quite intuitive: as the general number of edits created by bots grows, so grows the number of administrative tasks to be automated. Automation may include simplifying alerts, fully and semi-automated reverts, better user tracking, or automated corrections. However, the creation of automated methods for the verification and validation of the ontological relations it contains is required most.

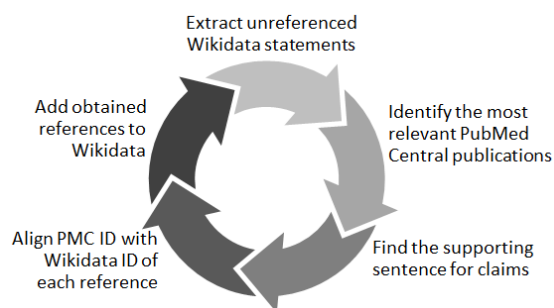
### 3. Knowledge graph validation of Wikidata

As Wikidata properties are assigned labels, descriptions and aliases in multiple languages (Red in Fig. 6), multilingual information of these properties can be used alongside the labels, descriptions, and aliases of Wikidata items to verify and find sentences supporting biomedical statements in scholarly outputs [53]. Such a process can be based on various natural language processing techniques, including word embeddings [53, 54] and semantic similarity [55]. These techniques are robust enough to achieve an interesting level of accuracy, and some of them can achieve better accuracy when the Wikidata classes of the subject and object of semantic relations are given as inputs [56, 57].

The subjects and objects of Wikidata relations can likewise be aligned to other biomedical semantic resources such as MeSH and UMLS Metathesaurus [11]. Thus, benchmarks for relation extraction based on one of the major biomedical ontologies can be converted into a Wikidata friendly format and used to automatically enrich Wikidata with novel biomedical relations or to automatically find statements supporting existing biomedical Wikidata relations [58]. Furthermore, MeSH keywords of scholarly publications can be converted into their Wikidata equivalents, refined using citation and co-citation analysis [59], and used to verify and add biomedical Wikidata relations, e.g. by applying deep learning-based bibliometric-enhanced information retrieval techniques [60, 61].

Another option of validating biomedical statements based on the labels of their subjects, predicates, and objects in Wikidata can be the use of these labels for the reformulation of a query to search bibliographic databases and consequently to find appropriate references for the assessed Wikidata

statements (Example in Fig. 5). Several bots and bot frameworks have been successfully built using this principle such as Wikidata Integrator<sup>11</sup> that extracts the Wikidata statements of a given gene or protein using SPARQL, compare them with their equivalents in other structured databases like NCBI's Gene resources and Uniprot and adjust them if needed, and RefB<sup>12</sup> (Fig. 1) that extracts biomedical Wikidata statements not supported by references using SPARQL and identifies the sentences supporting them in scholarly publications using PubMed Central search engine and a variety of techniques such as concept proximity analysis.



**Fig. 1.** Process of RefB, a bot that adds scholarly references to biomedical Wikidata statements based on PubMed Central [Source: [https://w.wiki/an\\$](https://w.wiki/an$), License: CC BY 4.0]. The source code of RefB is available at <https://github.com/Data-Engineering-and-Semantics/refb/>

In addition to their multilingual set of labels and descriptions, Wikidata properties are assigned object types using `wikibase:propertyType` relations (Blue in Fig. 2). These relations allow the assignment of appropriate objects to statements, so that non-relational statements cannot have a Wikidata item as an object, while objects of relational statements are not allowed to have data types like a value or a URL [10].

<sup>11</sup> Wikidata Integrator is a bot framework for automatically curating genetic information provided by Wikidata (<https://github.com/SuLab/WikidataIntegrator>). For Wikidata bots using this framework, refer to [https://www.wikidata.org/wiki/Wikidata:WikiProject\\_Gene\\_Wiki\\_Bot\\_accounts](https://www.wikidata.org/wiki/Wikidata:WikiProject_Gene_Wiki_Bot_accounts).

<sup>12</sup> RefB: *Description* at [https://www.wikidata.org/wiki/Wikidata:Requests\\_for\\_permissions/Bot/RefB\\_\(WikiCred\)](https://www.wikidata.org/wiki/Wikidata:Requests_for_permissions/Bot/RefB_(WikiCred)), *Source code* at <https://github.com/Data-Engineering-and-Semantics/refb/>, *Wikidata edits* at [https://www.wikidata.org/wiki/Special:Contributions/RefB\\_\(WikiCred\)](https://www.wikidata.org/wiki/Special:Contributions/RefB_(WikiCred)).



**symptoms** (P780)

possible symptoms of a medical condition edit

• In more languages

Language	Label	Description	Also known as
English	symptoms	possible symptoms of a medical condition	
French	symptômes	manifestations ressenties par le patient atteint d'une maladie, plaintes exprimées par celui-ci	signes fonctionnels
Central Atlas Tamazight	No label defined	No description defined	
Arabic	الأعراض	No description defined	

All entered languages

Data type

Item

Statements

instance of Wikidata property related to medicine edit

• 0 references + add reference + add value

subject item of this property symptom edit

• 0 references + add reference + add value

Wikidata property example meningitis edit

symptoms headache + add reference + add value

equivalent property <https://schema.org/signOrSymptom> edit

• 0 references + add reference + add value

Constraints

property constraint

value type constraint edit

class clinical sign + add reference

symptom + add reference

relation instance or subclass of + add value

type constraint edit

class physiological condition + add reference

fictional medical condition + add reference

relation instance or subclass of + add value

citation needed constraint edit

• 0 references + add reference + add value

**Fig. 2.** Wikidata page of a clinical property [Source: <https://w.wiki/aeF>, Derived from: <https://w.wiki/aeG>, License: CC BY-SA 4.0]. It includes the labels, descriptions and aliases of the property in multiple languages (Red), the object data type (Blue), statements where the property is the subject (Green) as well as property constraints (Brown).

Just like a Wikidata item, a property can be described by statements (Green in Fig. 2). The predicates of these statements link a property to its class (*instance of* [P31]), to its corresponding Wikidata item (*subject item of this property* [P1629]), to example usages (*Wikidata property*

*example* [P1855]), to equivalents in other IRIs<sup>13</sup> (*equivalent property* [P1628]), to Wikimedia categories that track its usage on a given wiki (*property usage tracking category* [P2875]), to its inverse property (*inverse property* [P1696]), or to its proposal discussion (*property proposal discussion* [P3254]), etc.

These statements can be interesting for various knowledge graph validation purposes. In fact, the class, the usage examples and the proposal discussion of a Wikidata property can be useful through the use of several natural language processing techniques, particularly semantic similarity, to provide several features of the use of the property such as its domain of application (e.g. the subject or object of a statement using a Wikidata property related to medicine should be a medical item) and consequently to eliminate some of erroneous use by screening the property usage tracking category. The class of the Wikidata item corresponding to the property can be used to identify the field of work of the property and thus flag some inappropriate applications. In addition, the external identifiers of such an item can be used for the verification of biomedical relations by their identification within the semantic annotations of scholarly publications built using the SAT+R (Subject, Action, Target, and Relations) model [62]. The inverse property relations can identify missing Wikidata statements ( $C_1$ ,  $P$ ,  $C_2$ ), which are implied by the presence of inverse statements ( $C_2$ ,  $P^{-1}$ ,  $C_1$ ) in other Wikidata resources. Here,  $P^{-1}$  is the Wikidata property that is the inverse of  $P$ ,  $C_S$  is a common class of the subjects of  $P$ , and  $C_O$  is a common class of the objects of  $P$ .

Despite the importance of these statements defining properties, *property constraint* [P2302] relations (Brown in Fig. 2) are the semantic relations that are primarily used for the validation of the usage of a property. In essence, they define a set of conditions for the use of a property, including several heuristics for the type and format of the subject or the object, information about the characteristics of the property, and several description logics for the usage of the property as shown in Table 1. Property constraints are either manually added by Wikidata users or inferred with an excellent accuracy from the knowledge graph of Wikidata or the history of human changes to Wikidata statements [63, 64].

<sup>13</sup> Internationalized Resource Identifier (IRI) is a standardized character string (e.g. a URL) that recognizes a given item in a semantic resource

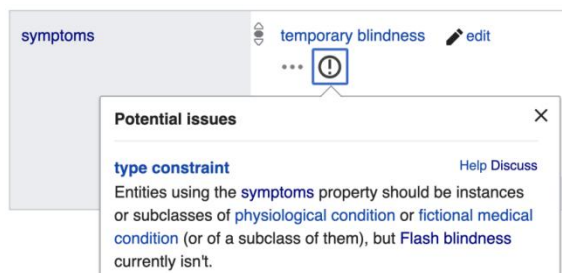
Table 1  
Constraint types for the usage of Wikidata properties

Wikidata ID	Constraint type	Description
Q19474404	single value constraint	Constraint used to specify that this property generally contains a single value per item
Q21502404	format constraint	Constraint used to specify that the value for this property has to correspond to a given pattern
Q21502408	mandatory constraint	status of a Wikidata property constraint: indicates that the specified constraint applies to the subject property without exception and must not be violated
Q21502410	distinct values constraint	Constraint used to specify that the value for this property is likely to be different from all other items
Q21510852	Commons link constraint	Constraint used to specify that the value must link to an existing Wikimedia Commons page
Q21510854	difference within range constraint	Constraint used to specify that the value of a given statement should only differ in the given way. Use with qualifiers minimum quantity/maximum quantity
Q21510856	mandatory qualifier constraint	Constraint used to specify that the listed qualifier has to be used
Q21510862	symmetric constraint	Constraint used to specify that the referenced entity should also link back to this entity
Q21510863	used as qualifier constraint	Constraint used to specify that a property must only be used as a qualifier
Q21510864	value requires statement constraint	Constraint used to specify that the referenced item should have a statement with a given property
Q21510495	relation of type constraint	relation establishing dependency between types/metalevels of its members
Q21510851	allowed qualifiers constraint	Constraint used to specify that only the listed qualifiers should be used. No value disallows any qualifier
Q21510865	value type constraint	Constraint used to specify that the referenced item should be a subclass or instance of a given type
Q21514353	allowed units constraint	Constraint used to specify that only listed units may be used
Q21510857	multi-value constraint	Constraint used to specify that a property generally contains more than one value per item
Q21510859	one-of constraint	Constraint used to specify that the value for this property has to be one of a given set of items
Q21510860	range constraint	Constraint used to specify that the value must be between two given values
Q21528958	used for values only constraint	Constraint used to specify that a property can only be used as a property for values, not as a qualifier or reference
Q21528959	used as reference constraint	Constraint used to specify that a property must only be used in references or instances of citation (Q1713)
Q25796498	contemporary	Constraint used to specify that the subject

	constraint	and the object have to coincide or coexist at some point of history
Q21502838	conflicts-with constraint	Constraint used to specify that an item must not have a given statement
Q21503247	item requires statement constraint	Constraint used to specify that an item with this statement should also have another given property
Q21503250	type constraint	Constraint used to specify that the item described by such properties should be a subclass or instance of a given type
Q54554025	citation needed constraint	Constraint specifies that a property must have at least one reference
Q62026391	suggestion constraint	status of a Wikidata property constraint: indicates that the specified constraint merely suggests additional improvements, and violations are not as severe as for regular or mandatory constraints
Q64006792	lexeme value requires lexical category constraint	Constraint used to specify that the referenced lexeme should have a given lexical category
Q42750658	value constraint	class of constraints on the value of a statement with a given property. For constraint: use specific items (e.g. "value type constraint", "value requires statement constraint", "format constraint", etc.)
Q51723761	no bounds constraint	Constraint specifies that a property must only have values that do not have bounds
Q52004125	allowed entity types constraint	Constraint used to specify that only listed entity types are valid for this property
Q52060874	single best value constraint	Constraint used to specify that this property generally contains a single "best" value per item, though other values may be included as long as the "best" value is marked with preferred rank
Q52558054	none of constraint	Constraint specifying values that should not be used for the given property
Q52712340	one-of qualifier value property constraint	Constraint used to specify which values can be used for a given qualifier when used on a specific property
Q52848401	integer constraint	Constraint used when values have to be integer only
Q53869507	property scope constraint	Constraint to define the scope of the property (main value, qualifier, references, or combination); only supported by KrBot currently

As shown in Fig. 2, a property constraint is defined as a relation where the property type is featured as an object and the detailed conditions of the constraint to be applied on Wikidata statements are integrated as qualifiers to the relation. When a property constraint is violated, the corresponding statement is automatically included in a report of

property constraint violations<sup>14</sup> and is marked by an exclamation mark on the page of the subject item (Fig. 3) so that it can be quickly processed and adjusted by the community or by Wikidata bots if applicable.



**Fig. 3.** Example of a property constraint violation marked in the page of a Wikidata item, Q3603152 (flash blindness) [Available on Wikimedia Commons: <https://w.wiki/ZuJ>, license: CC0]

Although these methods are important to verify and validate Wikidata, they are not the only ones that are used for these purposes. In 2019, Wikidata announced the adoption of Shape Expressions language (ShEx) as part of the Mediawiki entity schemas extension<sup>15</sup>. ShEx was proposed following an RDF validation workshop that was organized by W3C<sup>16</sup> in 2014 as a concise, high-level language to describe and validate RDF data [65]. This Mediawiki extension uses ShEx to store structure definitions (EntitySchemas or Shapes) for sets of Wikidata entities which are selected by some query pattern (frequently the involvement of said entities in a Wikidata class). This provides collaborative quality control where the community can iteratively develop a schema and refine the data to conform to that schema. For those familiar with XML, ShEx is analogous to XML Schema or RelaxNG. *SHACL* (Shapes Constraint Language), another language used to constraint RDF data models, uses a flat list of constraints, analogous to XML's Schematron. It was adapted from SPIN (SPARQL Inference Notation) by the W3C Data Shapes working group in 2014 and became a W3C recommendation in 2017 [66]. However, ShEx was chosen to represent EntitySchemas in Wikidata, as it has a compact syntax which makes it more human-friendly, supports recursion, and is designed to support

distributed networks of reusable schemas [67]. Besides the possibility to infer ShEx expressions from the screening of a large set of concerned items, they can be easily written by humans in an intuitive way.

In Wikidata, ShEx-based EntitySchemas are assigned an identifier (a number beginning with an E) as well as labels, descriptions, and aliases in multiple languages, so that they can be easily identified by users. Entity schemas are defined using the ShEx-compact syntax<sup>17</sup>, which is a concise, human-readable syntax. A schema usually begins by some prefix declarations similar to SPARQL. An optional start definition declares the shape which will be used by default. In the example (Fig. 4), the shape <app> will be used, and its declaration contains a list of properties, possible values, and cardinalities. By default, shapes are open, which means that other properties apart from the ones declared are allowed. In this example, the values of property `wdt:P31` are declared to be either a COVID-19 dashboard (`wd:Q90790055`), a search engine (`wd:Q91136116`) or a dataset (`wd:Q91137337`). The *EXTRA* directive indicates that there can be additional values for property `wdt:P31` that differ from the specified ones. The value for property `wdt:P1476` is declared to be zero or more literals. The cardinality indicators come from regular expressions, where '?' means zero or one, '\*'; means zero or more, and '+' means one or more. While the values for the other properties are declared to be anything (the dot indicates no constraint) zero or more times, except for the properties `wdt:P577` and `wdt:P7103` that are marked as optional using the question mark. Further documentation about ShEx can be found at <http://shex.io/> and in Labra Gayo et al. (2017) [67].

```
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX wdt: <http://www.wikidata.org/prop/direct/>
PREFIX wd: <http://www.wikidata.org/entity/>

start = @<app>

<app> EXTRA wdt:P31 {
  wdt:P31 [ wd:Q90790055 # instance of
COVID-19 dashboard or
          wd:Q91136116 # search engine or
          wd:Q91137337 # dataset
        ] ;
  wdt:P1476 LITERAL * ; #title
```

<sup>14</sup>

[https://www.wikidata.org/wiki/Wikidata:Database\\_reports/Constraint\\_violations](https://www.wikidata.org/wiki/Wikidata:Database_reports/Constraint_violations)

<sup>15</sup> <https://www.mediawiki.org/wiki/Extension:EntitySchema>

<sup>16</sup> <https://www.w3.org/2012/12/rdf-val/report>

<sup>17</sup> ShEx schemas can also be defined in RDF-based representations like Turtle or JSON-LD.



```

wdt:P366 . * ; #use
wdt:P123 . * ; #publisher
wdt:P178 . * ; #developers
wdt:P495 . * ; #country of origin
wdt:P306 . * ; #operating system
wdt:P856 . * ; #official website
wdt:P921 . * ; #main subject
wdt:P144 . * ; #based on
wdt:P577 . ? ; #publication date
wdt:P7103 . ? ; #start of covered
period
wdt:P275 . * ; #copyright license
wdt:P5008 . * ; #on focus list of
Wikimedia project
}

```

**Fig. 4.** EntitySchema for COVID-19 dashboards, search engines and datasets [Source: <https://www.wikidata.org/wiki/EntitySchema:E205>]

Due to the ease of using ShEx to define EntitySchemas, it has been used successfully to validate Danish lexemes in Wikidata [68] and biomedical Wikidata statements [69]. During the COVID-19 pandemic, Wikidata's data model of every COVID-19-related class as well as of all major biomedical classes has been converted to an EntitySchema, so that it can be used to validate the representation of COVID-19 Wikidata statements [12]. These EntitySchemas were successfully used to enhance the development and the robustness of the semantic structure of the data model underlying the COVID-19 knowledge graph in Wikidata and are accordingly made available at a subpage of Wikidata's WikiProject COVID-19, accessible via [https://www.wikidata.org/wiki/Wikidata:WikiProject\\_COVID-19/Data\\_models](https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Data_models).

#### 4. SPARQL as a semantic query language

SPARQL was officially created in 2008 as a query language and protocol to search, add, modify or delete RDF data available over the Internet. Its name is a recursive acronym which stands for "SPARQL Protocol and RDF Query Language". SPARQL allows a query to be composed of triple patterns, conjunctions, disjunctions, and optional patterns and can consequently be used to retrieve contextualized information from knowledge graphs. As it has been designed to extract a searched pattern from a semantic graph [70], SPARQL queries have also been used to query the competency questions<sup>18</sup>, so as to evaluate ontologies and knowledge graphs in a

<sup>18</sup> Competency questions: A set of requirements ensuring consistency of a knowledge graph, constraints determining what knowledge to be involved in a knowledge graph [71].

context-sensitive way [72-74]. Indeed, a sister project presents how SPARQL can be used to generate data visualisations<sup>19</sup> [35, 75]. Validating RDF data portals using SPARQL queries has been regularly proposed as an approach that gives great flexibility and expressiveness [76]. However, academic literature is still far from revealing a consensus on methods and approaches to evaluate ontologies using this query language [77], and other approaches have been proposed for validation [69, 78].

SPARQL<sup>20</sup> is a human-friendly language based on defining triples as conditions [70] and defines prefixes to abbreviate IRIs similarly to like ShEx (Blue in Fig. 4). It also uses the skeleton of SQL to define queries to knowledge graphs in RDF format [79]. For example, SPARQL shares most of SQL's clauses used to retrieve variables and aggregate functions used to compute new variables, as shown in Table 2 [79-81]. SPARQL also defines new aggregate function-based variables in the SELECT clause using the (function(variable) AS new\_variable) format, and constant values and strings are put between quotation marks. It defines logical conditions in the HAVING clause for variables based on aggregate functions or in the WHERE clause for variables to be retrieved from the source database as FILTER (condition) [80, 81]. The declaration of the logical conditions also uses the same operators (AND [also &&], || [OR], and NOT [also !]), values (True, False, and Null), logical functions (EXISTS [verifies the existence of a condition or a statement], NOT EXISTS [the opposite of EXISTS], and MINUS [eliminate the set of values having a given characteristic from the results]), and mathematical operators (> [superior to], < [inferior to], = [equal to], >= [superior or equal to], <= [inferior or equal to], != [different from], + [plus], - [minus], \* [times], and / [divide]) [81].

Table 2

List of clauses and aggregate functions available in SQL and SPARQL [80, 81]

Clauses	Description
SELECT	define variables to show
SELECT DISTINCT	define variables and omit repeated results

<sup>19</sup> For SPARQL-based visualizations of COVID-19 information in Wikidata, see <https://speed.ieee.tn/https://egonw.github.io/SARS-CoV-2-Queries/>, [https://www.wikidata.org/wiki/Wikidata:WikiProject\\_COVID-19/Queries](https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Queries), and <https://scholia.toolforge.org/topic/Q84263196>.

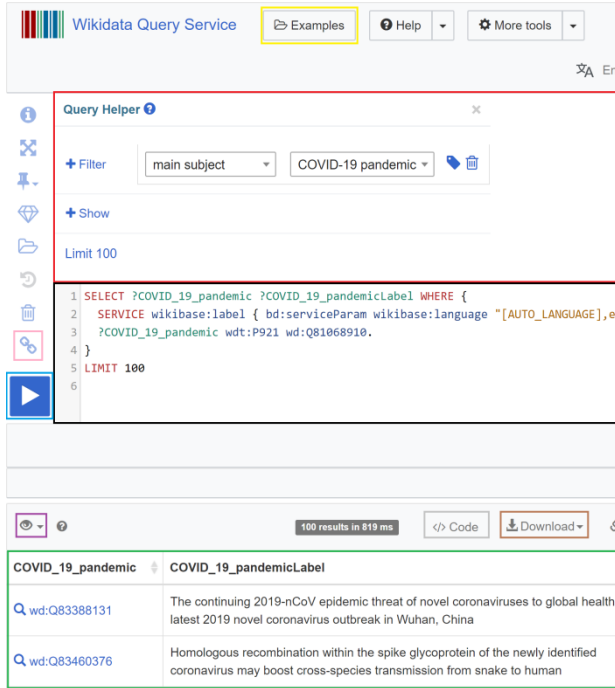
<sup>20</sup> An open license SPARQL textbook available in multiple languages can be found at <https://en.wikibooks.org/wiki/SPARQL>.

FROM	define the source database of the defined query
WITH	define a subquery
WHERE	restrict returned results by conditions in the form of triples
LIMIT	restrict returned results by to a total number
OFFSET	skip a number of first returned results
HAVING	restrict returned results by based on aggregate functions
GROUP BY	group entries to compute an aggregate function
ORDER BY	sort results according to a given variable
Functions	Description
AVG	Average of non-NULL values in a set.
COUNT	Number of results in a group, including the ones with NULL values
MAX	Maximum in a set of non-NULL values
MIN	Minimum in a set of non-NULL values
STDEV	Standard deviation of all values provided in the expression based on a limit set of results
STDEVP	Standard deviation for all values in the provided expression based on all the returned results
SUM	Sum of all non-NULL values in a set
VAR	Statistical variance of values in an expression based on a limit set of results
VARP	Statistical variance of values in an expression based on all the returned results

In contrast to SQL, the variables in SPARQL are preceded by an interrogation mark and are not separated by a comma in the SELECT clause [81, 82], and even the declaration of statements in the WHERE clause using SPARQL is different from the one using SQL. In the latter, the declaration of the statements in a WHERE clause can only be done in a single line [79]. When multiple statement conditions should be fulfilled, they have to be linked using the INTERSECT operator [83]. When a unique condition from a list of statements should be respected, the list's statements should be linked using the UNION operator [83]. Where results fulfilling a given condition should be eliminated, the condition must be preceded by the MINUS operator [83]. In SPARQL, the WHERE clause can include multiple lines between curly brackets, where each line is in the form of a subject-predicate-object triple [79]. When the statements between brackets are in the form of a triple, they should end with a period. When two successive statements have the same subject, the first statement can end with a semicolon. In this particular situation, the subject of the second statement can be omitted [81, 82]. An exception to this is the FILTER() function allowing the definition of a logical condition to be considered or the BIND() function allowing the creation of a new variable based on the retrieved characteristic of a single result

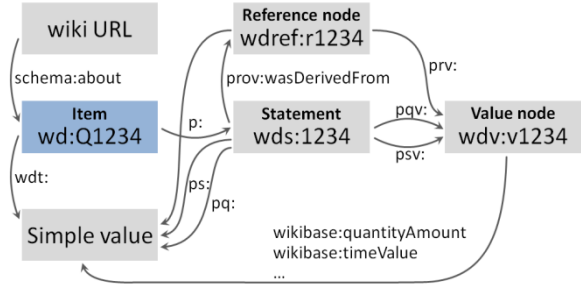
row [81, 82]. Although the MINUS and UNION operators can be used as in SQL, the INTERSECT operator is useless and is forsaken in SPARQL and the MINUS and UNION operators should be preceded and followed by statements between curly brackets like the WHERE clause [81, 82]. SPARQL has also the advantage to allow including entries where a set of statements in the WHERE clause is not respected by putting these statements after the OPTIONAL operator between curly brackets [81, 82].

In Wikidata, the Wikidata Query Service (<https://www.wikidata.org>) allows to query the knowledge graph using SPARQL [11, 34]. The required Wikidata prefixes are already supported in the backend of the service and do not need to be defined [34]. What the user needs to do is to formulate their SPARQL query (Black in Fig. 5) and click on the Run button (Blue in Fig. 5). After a compilation period, the results will appear (Green in Fig. 5) and can be downloaded in different formats (Brown in Fig. 5), including JSON, TSV, CSV, HTML, and SVG. Different modes for the visualization of the query results can be chosen (Purple in Fig. 5), particularly table, charts (line, scatter, area, bubble), image grid, map, tree, timeline, and graph. The query service also allows users to use a query helper (Red in Fig. 5) that can generate basic SPARQL queries and get inspired by sample queries (Yellow in Fig. 5), especially when they lack experience. It also allows us to generate a short link for the query (Pink in Fig. 5) and codes to embed the query results in web pages and computer programs (Brown in Fig. 5) [34].



**Fig. 5.** Web interface of Wikidata Query Service [Source: <https://w.wiki/aeH>, Derived from: <https://query.wikidata.org>, License: CC BY-SA 4.0]. It involves a query field (Black), a query builder (Red), a short link button (Pink), a Run button (Blue), a visualization mode button (Purple), a download button (Brown), an embedding code generation button (Grey), a results field (green), and a sample query button (Yellow).

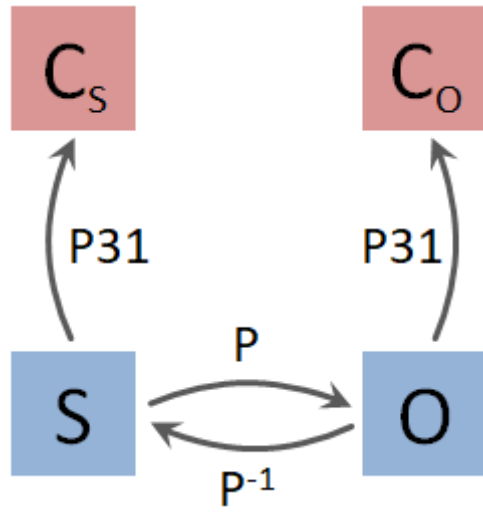
The statements in the WHERE clause should be defined such that known subjects and objects are preceded by *wd* prefix whatever they are Wikidata items or properties and that the predicate should be a Wikidata property, and it is preceded by *wdt* prefix as clearly shown in Fig. 6. Other Wikidata prefixes can be used to parse Wikidata qualifiers (*pq* and *pqv*) and references (*pr* and *prv*) or to link between a Wikidata statement to one of its components (*p*, *prov*, *ps*, and *psv*). The *wikibase* prefix can be used to return the characteristics of an item, a property or a statement. For example, *wikibase:directClaim* and *wikibase:Claim* can shift a property from a Wikidata prefix to another one (e.g. shifting Wikidata properties from *wdt* to *wd*), and *wikibase:rank* can be useful to return the level of importance assigned by the community to a statement.



**Fig. 6.** RDF data structure of Wikidata knowledge graph [Available at: <https://w.wiki/any>, adapted from source: <https://w.wiki/ZUA>, Michael F. Schöntzner, CC-BY 4.0]

## 5. Constraint-driven inference of biomedical property constraints

As described above, Wikidata properties are assigned property constraints and statements as logical conditions for the use of the types of triples to represent knowledge in Wikidata (Fig. 2). Screening Wikidata items in a class to identify common features of the assessed entities based on a set of formal rules has been previously proposed [64, 84]. These features involve common characteristics of the data model of the concerned class as well as patterns of used Wikidata properties such as symmetry and are later used to verify the completeness of the class and validate the statements related to the evaluated class using SPARQL queries. In this work, we propose a similar protocol fully based on logical constraints fully implementable using SPARQL queries to infer constraints for the assessment of the usage of relation types (*P*) on Wikidata based on the most frequently used corresponding inverse statements ( $C_O$ ,  $P^{-1}$ ,  $C_S$ ). These constraints can be later used to define COVID-19-related Wikidata statements and to generate ShEx schemas for COVID-19-related Wikidata classes. Fig. 7 represents the scheme of the given relation type that will be used to assess and validate the use of Wikidata properties. If we consider COVID-19 <drug used for treatment> tocilizumab as an accurate relational statement in Wikidata, *COVID-19* is the subject (*S*), *drug used for treatment* is the relation type (*P*), tocilizumab is the object (*O*), *medical condition treated* is the inverse relation type ( $P^{-1}$ ), *disease* is the subject class ( $C_S$ ) and *medication* is the object class ( $C_O$ ).



**Fig. 7.** Scheme of a given Wikidata property [Source: <https://w.wiki/anw>, License: CC BY 4.0]: S and O are respectively the subject and the object of the statement, P is the predicate of the statement, P-1 is the inverse property of P, CS is the class of the subject, and CO is the class of the object.

Once retrieved, the common inverse property statements ( $C_O$ ,  $P^{-1}$ ,  $C_S$ ) of the given Wikidata property  $P$  will be used to identify relations that use  $P$  in an uncommon and probably wrong way, to identify missing inverse relations of  $P(S,O)$  corresponding to the most used ( $C_O$ ,  $P^{-1}$ ,  $C_S$ ) scheme, and to identify the Wikidata items missing statements using  $P$  as shown in Table 3. The assessment of the usage of the given Wikidata property will not be restricted to these tasks, as it also involves the identification of relations using  $P$  not supported by references and the identification of Wikidata properties used to define references for relations using  $P$ .

Table 3

Tasks for quality assessment of the usage of Wikidata relation types using the Wikidata SPARQL endpoint

Task	Description
<b>Defining the scheme of a Wikidata property</b>	
T1	Identify common use cases <sup>21</sup> of $P$ : ( $C_S, C_O$ ) pairs
T2	Identify inverse properties of $P$ corresponding to each common use case: ( $C_S$ , $P^{-1}$ , $C_O$ ) statements
<b>Identifying the deficiencies of the scheme</b>	
T3	For each returned $P^{-1}$ , identify $P(S,O)$ relations supported by references and corresponding to the most common ( $C_S$ , $P^{-1}$ , $C_O$ ) statement but not available in Wikidata
T4	Identify $P(S,O)$ relations not corresponding to the most common scheme of $P$
<b>Assessing the reference support of relations using the studied</b>	

<sup>21</sup> Use case: A set of conditions for the use of a relation type  $P$ .

Wikidata property	
T5	Identify Wikidata properties used to define references for relations using $P$

This task set is useful to assess and adjust the reference support, the language support, the quantity, and the quality of the relations using  $P$  and  $P^{-1}$  at a given point in time and can be easily completed using the Wikidata Query Service. The SPARQL query of each task is given in Appendix A, where  $\langle \text{PropertyID} \rangle$  is the Wikidata ID of the studied property  $P$ ,  $\langle \text{SubjectID} \rangle$  is the Wikidata ID of the subject class  $C_S$  that is most used with this property, and  $\langle \text{ObjectID} \rangle$  is the Wikidata ID of the most used object class  $C_O$ .

For Tasks T1 and T2, we eliminated property use cases where classes  $C_S$  and  $C_O$  are first-order metaclasses (Q24017414), so that we do not get nonspecific use cases. Additionally, we only considered use cases applied to more than a defined usage threshold (here set as 100 but can change according to context) in order to omit statements that are not widely used in Wikidata. For Task T4, we used logical constraints to find statements where the subject is not an instance of the most used subject class  $C_S$  (G1), then to find statements where the object is not an instance of the most used object class  $C_O$  (G2). After that, we identified the statements that exist in both G1 and G2 as the most likely wrong statements ( $G1 \cap G2$ ) as they correspond neither to the most used subject class nor to the most used object class of the studied property. Such a task can either identify an accurate relation where the subject and object are not assigned to the corresponding Wikidata class due to the lack of completeness of Wikidata taxonomy or recognize a wrong Wikidata statement. For Task T5, Wikidata properties used to define fewer than a threshold number of references using  $P$  were not considered (again, here set to 100). Our analysis was performed on September 20, 2019, following the Zika outbreak as a proactive action to build the data model infrastructure to support clinical information about future infectious epidemics in Wikidata (the date is relevant due to the rapidly expanding nature of the database).

To assess the effectiveness of the use of logical constraints to generate conditions for the verification and validation of the use of relation types to enrich the Wikidata ontology, we applied our method to the main six Wikidata properties that can be used to represent COVID-19-related knowledge (Table 4).

Table 4  
Wikidata properties assessed in this study

Property	Description	Statements
Drug used for treatment (P2176)	drug, procedure, or therapy that can be used to treat a medical condition	6344
Significant drug interaction (P769)	clinically significant interaction between two pharmacologically active substances (i.e., drugs and/or active metabolites) where concomitant intake can lead to altered effectiveness or adverse drug events.	1850
Medical condition treated (P2175)	disease that this pharmaceutical drug, procedure, or therapy is used to treat	6499
Symptoms (P780)	possible symptoms of a medical condition	8068
Route of administration (P636)	path by which a drug, fluid, poison, or other substance is taken into the body	2900
Therapeutic area (P4044)	disease area in which a medical intervention is applied	1320

Task T1 was effective at sorting the common use cases of the studied Wikidata properties as shown in Table 5. All the retrieved use cases were proven to be logically accurate when compared to the descriptions

Table 5  
Common use cases of the studied Wikidata properties

Wikidata ID	Property	Subject Class	Object Class	Number of Statements
P2176	Drug used for treatment	Disease (Q12136)	medication (Q12140)	4777
		Disease (Q12136)	essential medicine (Q35456)	1558
		Infectious disease (Q18123741)	medication (Q12140)	558
		Disease (Q12136)	Heterocyclic compound (Q193430)	484
		Disease (Q12136)	Biopharmaceutical (Q679692)	471
P636	Route of administration	medication (Q12140)	route of administration (Q621636)	179
P4044	Therapeutic area	Pharmaceutical product (Q28885102)	disease (Q12136)	1147
		mixture (Q169336)	disease (Q12136)	1142
		Pharmaceutical product (Q28885102)	rare disease (Q929833)	142
		mixture (Q169336)	rare disease (Q929833)	141
		Pharmaceutical product (Q28885102)	Designated intractable/rare diseases (Q42303753)	115
P769	Significant drug interaction	medication (Q12140)	medication (Q12140)	1729
		medication (Q12140)	essential medicine (Q35456)	524
		essential medicine (Q35456)	medication (Q12140)	507
		medication (Q12140)	Heterocyclic compound	342

of Wikidata properties available in Table 4. The most common use cases for *drugs used for treatment* [P2176], *therapeutic area* [P4044], *significant drug interactions* [P769], or *medical condition treated* [P780] corresponded to 72 percent or more of the supported statements. However, there was a significant lack of availability of common use cases for *route of administration* [P636] and *symptoms* [P780]. This data deficiency may be due to human limitations (inexperience with wikidata or the medical logic being entered) from inconsistencies between languages (which often derive from slight differences in the naming and framing of articles in different language Wikipedias). These shortfalls could be alleviated by clearer taxonomy in attributing Wikidata items to corresponding classes.



			(Q193430)	
		Heterocyclic compound (Q193430)	medication (Q12140)	338
P2175	Medical condition treated	medication (Q12140)	Disease (Q12136)	4729
		essential medicine (Q35456)	Disease (Q12136)	1520
		medication (Q12140)	Infectious disease (Q18123741)	557
		Heterocyclic compound (Q193430)	Disease (Q12136)	487
		Biopharmaceutical (Q679692)	Disease (Q12136)	449
P780	Symptoms	disease (Q12136)	symptom (Q169872)	338
		disease (Q12136)	disease (Q12136)	264

Task T2 successfully sorted the inverse properties of Wikidata relation types for each corresponding use case as shown in Table 6. Here, we found that three relations had clear inverse properties: *medical condition treated* [P2175], *significant drug interaction* [P769] and *drug used for treatment* [P2176] are the inverse properties, respectively, for *drug used for treatment* [P2176], *significant drug interaction* [P769] and *medical condition treated* [P2175]. Hence, P2175 and P2176 are inverse to each other, and P769 is inverse to itself. However, we did not find any common inverse properties for *route of administration* [P636], *therapeutic area* [P4044] or *symptoms* [P780]. Consequently, the Task T2 can be used not only to find inverse properties of Wikidata relation types but also to identify Wikidata relation types where inverse properties do not exist or are not used as intended. In such a situation, the user should manually search for any inverse property to verify whether it exists or propose to the Wikidata community to create it as a new property<sup>22</sup> if it does not exist [11].

<sup>22</sup>e.g. Risk factor property proposal: [https://www.wikidata.org/wiki/Wikidata:Property\\_proposal/risk\\_factor](https://www.wikidata.org/wiki/Wikidata:Property_proposal/risk_factor)

Table 6

Inverse properties corresponding to each common use case of the studied Wikidata relation types

Wikidata ID	Property	Inverse property	Use case		Number of Statements
			Subject Class	Object Class	
P2176	Drug used for treatment	medical condition treated (P2175)	Disease (Q12136)	Medication (Q12140)	4576
		medical condition treated (P2175)	Disease (Q12136)	essential medicine (Q35456)	1482
		medical condition treated (P2175)	Infectious disease (Q18123741)	Medication (Q12140)	549
		medical condition treated (P2175)	Disease (Q12136)	Heterocyclic compound (Q193430)	477
		medical condition treated (P2175)	Disease (Q12136)	Biopharmaceutical (Q679692)	442
P636	Route of administration	NA			
P4044	Therapeutic area	NA			
P769	Significant drug interaction	Significant drug interaction (P769)	Medication (Q12140)	Medication (Q12140)	1330
		Significant drug interaction (P769)	Medication (Q12140)	essential medicine (Q35456)	359
		Significant drug interaction (P769)	essential medicine (Q35456)	Medication (Q12140)	359
		Significant drug interaction (P769)	Heterocyclic compound (Q193430)	Medication (Q12140)	288
		Significant drug interaction (P769)	Medication (Q12140)	Heterocyclic compound (Q193430)	288
P2175	Medical condition treated	Drug used for treatment (P2176)	Medication (Q12140)	Disease (Q12136)	4576
		Drug used for treatment (P2176)	essential medicine (Q35456)	Disease (Q12136)	1482
		Drug used for treatment (P2176)	Medication (Q12140)	Infectious disease (Q18123741)	549
		Drug used for treatment (P2176)	Heterocyclic compound (Q193430)	Disease (Q12136)	477
		Drug used for treatment (P2176)	Biopharmaceutical (Q679692)	Disease (Q12136)	442
P780	Symptoms	NA			

Table 7

Number of missing inverse statements of Wikidata relations supported by references and corresponding to the most used scheme of each Wikidata property

Wikidata ID	Property	Most used scheme			Missing inverse statements
		Inverse property	Subject Class	Object Class	
P2176	Drug used for treatment	medical condition treated (P2175)	Disease (Q12136)	Medication (Q12140)	160
P636	Route of administration	NA			
P4044	Therapeutic area	NA			
P769	Significant drug interaction	Significant drug interaction (P769)	Medication (Q12140)	Medication (Q12140)	385
P2175	Medical condition treated	Drug used for treatment (P2176)	Medication (Q12140)	Disease (Q12136)	143
P780	Symptoms	NA			

Task T3 effectively extracted those statements that use *drug used for treatment* [P2176], *significant drug interaction* [P769] and *medical condition treated* [P2175] as a Wikidata relation type where related

inverse relations do not exist in Wikidata as clearly stated in Table 7. Only relations corresponding to the most common use case of the related Wikidata property and supported by references are considered.

For the studied Wikidata relation types, 688 missing inverse statements were identified. An example of these statements is *COVID-19* [Q84263196] <*drug used for treatment* [P2176]> *dexamethasone* [Q422252] that does not exist in Wikidata despite the availability of its inverse relation (*dexamethasone* [Q422252] <*medical condition treated* [P2176]> *COVID-19* [Q84263196]) in the same knowledge graph.

These statements can be directly added to Wikidata using tools for the automatic enrichment of Wikidata, particularly QuickStatements [11], as they are supported by external references and are already stated in a Wikidata-friendly format.

Task T4 efficiently identified the statements not corresponding to the most common use case of the related Wikidata property as shown in Table 8. In fact, 11236 statements not corresponding to the most used subject class of the studied Wikidata properties and 7354 statements not corresponding to the most used object class of the studied Wikidata properties were identified.

Table 8

Number of statements not corresponding to the most common use case of each Wikidata property: Statements where the subject class is not the most used one (G1), statements where the object class is not the most used one (G2)

Wikidata ID	Property	G1	G2	G1∩G2
P2176	Drug used for treatment	858	390	72
P636	Route of administration	2656	1255	1255
P4044	Therapeutic area	5	171	3
P769	Significant drug interaction	82	42	3
P2175	Medical condition treated	620	1036	135
P780	Symptoms	7015	4460	3749

Among these statements, 5217 relations corresponded neither to the most common subject class nor to the most common object class of the considered properties. When applying expert validation to 800 randomly selected relations among the 5217 studied ones, we found that only 6.6% of these relations (53) were truly inaccurate and that the remaining 93.4% (747) were accurate but identified due to the lack of assignment of their subjects and objects to their hypernyms (i.e. a significant lack in defining relations between Wikidata items and

corresponding classes). An example of such accurate relations is (*alcohol withdrawal syndrome* [Q2914873], *Drug used for treatment* [P2176], (*RS*)-*baclofen* [Q413717]) where *alcohol withdrawal syndrome* is not an instance of disease [Q12136] and (*RS*)-*baclofen* is not declared as an instance of *medication* [Q12140]. The precision rate of the identification of deficient relations using this method seems to vary considering the studied property but does not exceed 10% (Fig. 8).

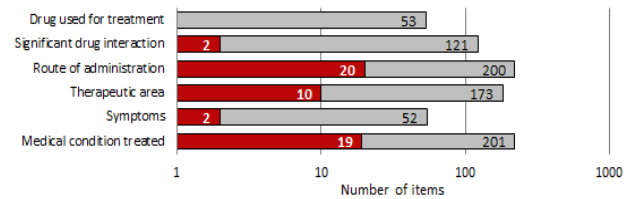


Fig. 8. Relations returned by Task T4 for the studied Wikidata properties [Available at: <https://w.wiki/ao2>, License: CC BY 4.0]. Extracted relations verified by expert validation as deficient are represented in red. Note: log x-axis.

Accordingly, the results sorted by Task T4 should be manually verified and validated by experts, so that users can use true identified relations (False positive) to enrich their respective subject and object Wikidata items with corresponding missing classes and find the reasons behind the deficiency of wrong identified relations (True positive) to develop automatic methods to solve them. The insufficiencies of wrong relations can either be due to ontological reasons (64%) or medicine-related reasons (36%) as shown in Fig. 9 and cannot consequently be handled only by computer scientists. Ontological reasons include the substitution of the accurate subject or object of a given relation with a semantically related item (e.g. The object of (*Renvela* [Q29006419], *Therapeutic area* [P4044], *hemodialysis* [Q391744]) should be *renal dialysis* [Q202301]) as well as the Subject-Object Inversion (e.g. the subject and object in (*botulism* [Q154845], *Medical condition treated* [P2175], *Heptavalent botulism antitoxin* [Q17148719]) have to be permuted) and the use of wrong property (e.g. The Wikidata property in (*MK-608* [Q23309937], *Significant drug interaction* [P769], *Zika virus RNA-dependent RNA polymerase NS5* [Q22954521]) should be *target of action*). Efforts in crowdsourcing ontology verification of other biomedical ontologies such as SNOMED-CT confirmed the existence of both types of errors and stipulated that not adjusting these lexical resources and using them in clinical decision support can generate harmful recommendations [85].

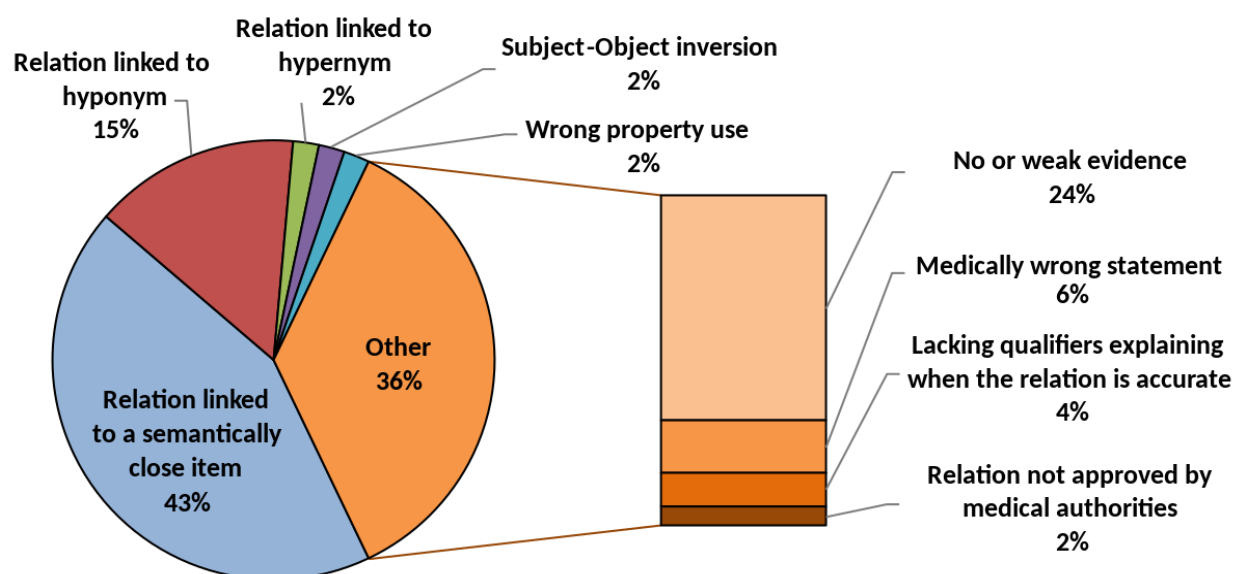


Fig. 9. Reasons of the inaccuracy of the truly deficient identified relations (True positive) [Source: <https://w.wiki/ao3>, License: CC BY 4.0]

Task T5 was efficient in finding the Wikidata properties used to define the references of the statements for each studied relation type (Table 9). For the studied Wikidata relation types, we found that references are mainly defined using three properties: *stated in* [P248], *retrieved* [P813], and *reference URL* [P854] (Example in Fig. 10). One of the highest priority tasks on Wikidata is for experts to find and add appropriate references using these three

properties to currently unsupported Wikidata relations. Once the references are in the system, further refinement is possible, e.g. a *reference URL* [P854] containing (or pointing to a page that contains) an external identifier for which Wikidata has a suitable property - e.g. *Digital Object Identifier* [P356] - then that property could be added to an item about the cited references, and the P854 statement replaced by a P248 statement pointing to that item.

Table 9

Number of statements not corresponding to the most common use case of each Wikidata property: Statements where the subject class is not the most used one (G1), statements where the object class is not the most used one (G2) Wikidata properties used to define references for studied Wikidata relation types: *stated in* (P248), *retrieved* (P813), *language of work or name* (P407), *National Drug File Reference Terminology ID* (P2115), *reference URL* (P854), and *European Medicines Agency product number* (P3637)

Wikidata ID	Property	P248	P813	P407	P2115	P854	P3637
P2176	Drug used for treatment	6654	6636	3617	3522	1626	
P636	Route of administration					2647	
P4044	Therapeutic area	1310	1310			1313	1310
P769	Significant drug interaction	1757					
P2175	Medical condition treated	6683	6672	3533	3516	1719	
P780	Symptoms	257	114			7094	

<b>drug used</b>	for treatment
	tocilizumab
nature of statement	clinical trial
alternate names	Actemra
manufacturer	Hoffmann-La Roche
quantity	400 milligram
route of administration	intravenous drip
▼ 3 references	
stated in	Effective Treatment of Severe COVID-19 Patients with Tocilizumab
reference URL	<a href="http://chinaxiv.org/abs/202003.00026">http://chinaxiv.org/abs/202003.00026</a>
publication date	5 March 2020



2020 COVID-19 pandemic in Tunisia (Q87343682)

viral outbreak in Tunisia  
2020 coronavirus outbreak in Tunisia [edit](#)

Statements

number of deaths	51	point in time	8 August 2020	<a href="#">edit</a>
			1 reference	
	53	point in time	14 August 2020	<a href="#">edit</a>
			1 reference	
case fatality rate	0.039	point in time	8 April 2020	<a href="#">edit</a>
			1 reference	
	0.038	point in time	4 April 2020 7 April 2020	<a href="#">edit</a>
			1 reference	
number of cases	879	point in time	18 April 2020	<a href="#">edit</a>
			1 reference	
	909	point in time	21 April 2020	<a href="#">edit</a>
			2 references	
number of hospitalized cases	93	point in time	22 April 2020	<a href="#">edit</a>
			1 reference	
	85	point in time	20 April 2020	<a href="#">edit</a>
			1 reference	
number of recoveries	190	point in time	21 April 2020	<a href="#">edit</a>
			2 references	
	170	point in time	20 April 2020	<a href="#">edit</a>
			1 reference	
number of clinical tests	12,531	point in time	13 April 2020	<a href="#">edit</a>
			1 reference	
	11,941	point in time	12 April 2020	<a href="#">edit</a>
			1 reference	

**Fig. 11.** Sample statistical data about the COVID-19 pandemic in Tunisia [Adapted from: <https://www.wikidata.org/wiki/Q87343682>, Source: <https://w.wiki/uUr>, License: CC BY 4.0].

From simple count statistics ( $c$ ,  $t$ ,  $d$ ,  $h$ , and  $r$  statements), it is possible to compare regional epidemiological variables and their variance for a

given date ( $Z$ ) or date range, and relate these to the general disease outbreak (each component defined as a *part of* [P361] of the general outbreak) as shown in Table 10. Tasks V1 and V2 have been generated from the evidence that COVID-19 started in late 2019 and that its clinical discovery can only be done through medical diagnosis techniques [87]. Tasks V3 and V4 have been derived from the fact that  $c$ ,  $d$ ,  $r$ , and  $t$  are cumulative counts. Consequently, these variables are only subjects to remain constant or increase over days. Task V5 is motivated by the fact that a simple epidemiological count cannot return negative values. Tasks V6, V7, V8, and V9 are due to the evidence that  $d$ ,  $r$ , and  $h$  cannot be superior to  $c$  as a patient needs to be affected by SARS-CoV-2 to die or be hospitalized due to the contraction of COVID-19 [86] and that a patient needs to undergo COVID-19 testing to be confirmed as a case of the disease [87]. V10 is built upon the assumption that  $c$ ,  $d$ ,  $r$ ,  $h$ , and  $t$  values can be geographically aggregated [86].

This task set has easily been applied using ten simple SPARQL queries that can be found in Appendix B where  $\langle \text{PropertyID} \rangle$  is the Wikidata property to be analyzed and has returned 5496 deficiencies in the COVID-19 epidemiological information as shown in Table 11. Among these mistaken statements, 2856 were *number of cases* statements, 2467 were *number of deaths* statements, 189 were *number of recoveries* statements, 9 were *number of clinical tests* statements, and 10 were *number of hospitalized cases* statements. This distribution of the deficiencies among epidemiological properties is explained by the dominance of *number of cases* and *number of deaths* statements on the COVID-19 epidemiological information. Most of these mistakes are linked to a violation of the cumulative pattern of major variables. These deficiencies can be removed using tools for the automatic enrichment of Wikidata like QuickStatements (cf. Turki, et al., 2019 [11]) or adjusted one by one by active members of WikiProject COVID-19.

Table 10

Tasks for the heuristics-based evaluation of epidemiological data using the Wikidata SPARQL endpoint

Task	Description	Sample filtered deficient statement
Validating qualifiers of COVID-19 epidemiological statements		
V1	Verify $Z$ as a date > November 01, 2019	<i>COVID-19 pandemic in X</i> <number of cases> 5 <point in

		time> <b>March 25, 20</b>
V2	Verify $Q$ as any subclass of (P279*) of medical diagnosis (Q177719)	<i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> March 25, 2020 <determination method> <b>COVID-19 Dashboard</b>
Ensuring the cumulative pattern of $c$ , $d$ , $r$ , and $t$		
V3	Identify $c$ , $d$ , $r$ and $t$ statements having a value in date $Z+1$ not superior or equal to the one in date $Z$ (Verify if $d_Z \leq d_{Z+1}$ , $r_Z \leq r_{Z+1}$ , $t_Z \leq t_{Z+1}$ , and $c_Z \leq c_{Z+1}$ )	( <i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> March 25, 2020) AND ( <i>COVID-19 pandemic in X</i> <number of cases> 6 <point in time> March 24, 2020)
V4	Find missing values of $c$ , $d$ , $r$ and $t$ in date $Z+1$ where corresponding values in dates $Z$ and $Z+2$ are equal	( <i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> March 24, 2020) AND ( <i>COVID-19 pandemic in X</i> <number of cases> 6 <point in time> March 26, 2020) AND ( <i>COVID-19 pandemic in X</i> <number of cases> <b>no value</b> <point in time> March 25, 2020)
Validating values of epidemiological data for a given date		
V5	Identifying $c$ , $d$ , $r$ , $h$ , and $t$ statements with negative values	<i>COVID-19 pandemic in X</i> <number of cases> -5 <point in time> March 25, 2020
V6	Identify $h$ statements having a value superior to the number of cases for a date $Z$	( <i>COVID-19 pandemic in X</i> <number of hospitalized cases> 15 <point in time> March 25, 2020) AND ( <i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> March 25, 2020)
V7	Identify $c$ statements having a value superior or equal to the number of clinical tests for a date $Z$	( <i>COVID-19 pandemic in X</i> <number of clinical tests> 4 <point in time> March 25, 2020) AND ( <i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> March 25, 2020)
V8	Identify $c$ statements having a value inferior to the number of deaths for a date $Z$	( <i>COVID-19 pandemic in X</i> <number of deaths> 10 <point in time> March 25, 2020) AND ( <i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> March 25, 2020)
V9	Identify $c$ statements having a value inferior to the number of recoveries for a date $Z$	( <i>COVID-19 pandemic in X</i> <number of recoveries> 10 <point in time> March 25, 2020) AND ( <i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> March 25, 2020)
V10	Comparing the epidemiological variables of a general outbreak with the ones of its components	( <i>COVID-19 pandemic in X</i> <number of cases> 10 <point in time> March 25, 2020) AND ( <i>COVID-19 pandemic in Y</i> <number of cases> 5 <point in time> March 25, 2020) WHERE $X$ is a district of $Y$

Table 11

Number of deficient statements for every type of epidemiological Wikidata property identified by each task (As of August 8, 2020)

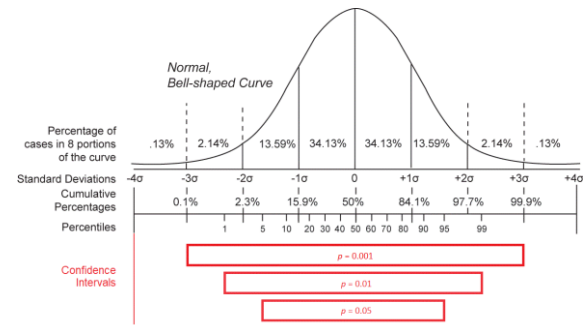
	$c$	$d$	$r$	$t$	$h$	Overall
V1	18	9	10	2	1	40
V2	2	91	6	0	0	99
V3	660	92	6	5		763
V4	2081	2247	149	1		4478
V5	0	0	0	0	0	0
V6	8				8	8
V7	1			1		1
V8	9	9				9
V9	17		17			17
V10	60	19	1	0	1	81

Overall	2856	2467	189	9	10	5496
---------	------	------	-----	---	----	------

Concerning the variables issued from the integration of basic epidemiological counts ( $m$ ,  $R_0$ ,  $mn$  and  $mx$  statements), they give a summary overview of the statistical behavior of the studied infectious pandemic and that is why they can be useful to identify if the stated evolution of the morbidity and mortality caused by the outbreak is reasonable [88]. However, the validation of these variables is more complicated due to the complexity of their definition [88-90]. The basic reproduction number ( $R_0$ ) is meant to be a constant that characterizes the dissemination power of an infectious disease. It is defined as the expected number of people (within a community with no prior exposure to the disease) that can contract a disease via the same infected individual. This variable should exceed the threshold of 1 to define a contagious disease [88]. Although  $R_0$  can give an idea about the general behavior of an outbreak of a given disease, any calculated value depends on the model used for its computation (e.g. *SIR* Model) as well as the

underlying data and is consequently a bit imprecise and variable from one study to another [88]. That is why it is not reliable to use this variable to evaluate the accuracy of simple epidemiological counts for a given pandemic. The only heuristic that can be applied to this variable is to verify if its value exceeds 1 for diseases causing large outbreaks. The incubation period of a disease gives an overview of the silent time required by an infectious agent to become active in the host organism and cause notable symptoms [89, 90]. This variable is very important as it reveals how many days an inactive case can spread the disease in the host's environment before the host is being symptomatically identified. As a result, it can give an idea about the contagiousness of the infectious disease and its basic reproduction number ( $R_0$ ). However, the determination of the incubation period - especially for a novel pathogen - is challenging, as a patient often cannot identify with precision the day when they had been exposed to the disease, at least if they did not travel to an endemic region or had not been in contact with a person they knew to be infected. This factor was behind the measurement of falsely small incubation periods for COVID-19 at the beginning of COVID-19 epidemic in China [89]. Furthermore, the use of minimal ( $mn$ ) and maximal ( $mx$ ) incubation periods in Wikidata to epidemiologically describe a disease instead of the median incubation period is a source of a lack of accuracy of the extracted values [89, 90]. In fact, minimal and maximal incubation periods for a given disease are obtained in the function of the mean ( $\bar{X}$ ) and standard deviation ( $\sigma$ ) of the measures of the confidence interval of observed incubation periods in patients. Effectively,  $mn$  is equal to  $\bar{X} - z * \frac{\sigma}{\sqrt{n}}$  and  $mx$  is equal to  $\bar{X} + z * \frac{\sigma}{\sqrt{n}}$  where  $n$  is the number of analyzed observations and  $z$  is a characteristic of the hypothetical statistical distribution and of the statistical confidence level adopted for the estimation [91]. As a consequence,  $mn$  and  $mx$  variables are modified according to the number of observations ( $n$ ) with a smaller difference between the two variables for higher values of  $n$ . As well, the two measures also vary according to the used statistical distribution and that is why different values of  $mn$  and  $mx$  were reported for COVID-19 when applying different distributions (Weibull, gamma and log-normal distribution) using a confidence level of 0.95 on the same set of observed cases [89]. Similarly, the two variables can change according to the adopted confidence level ( $p - 1$ ) when using the same statistical distribution where a higher confidence

level is correlated with a higher difference between the calculated  $mn$  and  $mx$  values, as shown in Fig. 12 [91, 92]. Given these reasons and despite the significant importance of the two measures, these two statistical variables cannot be used to evaluate statistical epidemiological counts for COVID-19 due to their lack of precision and difficulty of determination.



**Fig. 12.** Confidence intervals for different p-values ( $p$ ) when using a normal distribution [Source: <https://w.wiki/aKT>, License: Public Domain] (after Ward & Murray-Ward, 1999 [92]).

As for the reported case fatality rate ( $m$ ), its definition is less intricate than the ones of the basic reproduction number and of the incubation period, as  $m$  is only the quotient of the cumulative number of deaths ( $d$ ) by the cumulative number of cases ( $c$ ) as stated in official reports. It is consequently easy to validate for a given disease by comparing its values with simple reported counts of cases and deaths [86]. Here, two simple heuristics can be applied using SPARQL queries as shown in Appendix C. As the number of deaths is less than or equal to the number of cases of a given disease,  $m$  values should be set between 0 and 1. That is why Task M1 is defined to extract  $m$  statements where  $m > 1$  or  $m < 0$ . Also, as  $m = d / c$  for a date  $Z$ ,  $m$  values that are not close to the corresponding quotients of deaths by disease cases should be identified as deficient and  $m$  values should be stated for a given date  $Z$  if mortality and morbidity counts exist. Thus, Task M2 is created to extract  $m$  values where the absolute value of  $(m - d/c)$  is superior to 0.001, and Task M3 is developed to identify (item, date) pairs where  $m$  statements are missing and  $c$  and  $d$  statements are available in Wikidata. Absolute values for Task M2 are obtained using SPARQL's ABS function, and deficient (item, date) pairs are eliminated in Task M3 where  $m > 1$  and  $c < d$ .

As a result of these three tasks, we interestingly identified 143 deficient  $m$  statements and 7116

missing  $m$  statements. 133 of the mistaken statements are identified thanks to Task M2 and concern 25 Wikidata items and 31 distinct dates and only 10 deficient statements related to 3 Wikidata items and 8 distinct dates are found using Task M1. These statements should be verified against reference datasets to verify their values and to determine the reason behind their deficiency. Such a reason can be the integration of the wrong case and death counts in Wikidata or a bug or inaccuracy within the source code of the bot making or updating such statements. The verification process can be automatically done using an algorithm that compares Wikidata values ( $c$ ,  $d$  and  $m$  statements) with their corresponding ones in other databases (using file or API reading libraries) and subsequently adjusts statements using the Wikidata API directly or via tools like QuickStatements [11]. As for the missing  $m$  statements returned by M3, they are linked to 395 disease outbreak items and to 205 distinct dates and concern 70% (7116/10168) of the (case count, death count) pairs available in Wikidata. The outcome of M3 proves the efficiency of comparative constraints to enrich and assess the completeness of epidemiological data available in a knowledge graph, particularly Wikidata, based on existing information. Consequently, derivatives of Task M3 can build to infer  $d$  values based on  $c$  and  $m$  statements or to find  $c$  values based on  $d$  and  $m$  statements. The missing statements found by such tasks can be integrated in Wikidata using a bot based on Wikidata API and Wikidata Query Service to ameliorate the completeness and integrity of available mortality data for epidemics, mainly the COVID-19 pandemic [11].

## 7. Discussion

The results presented here demonstrate the value of our relational and statistical constraints-based validation approach for knowledge graphs like Wikidata across a range of features. In particular: identifying use cases of key relation types (Tables 5 and 6), verifying the completeness of inverse statements (Table 7), and aiding experts in finding deficiencies within the taxonomy and the non-taxonomic relations to manually address (Table 8 and Figures 8 and 9). These tasks successfully address most of the competency questions, particularly conceptual orientation (*clarity*), coherence (*consistency*), strength (*precision*) and full coverage (*completeness*). Combined with previous findings in the context of bioinformatics [84, 93-94], this proves

that the efficiency of rule-based approaches to evaluate semantic information from scratch displays a similar accuracy as other available ontology evaluation algorithms [95, 96].

The efficiency of these constraint-based assessment methods can be further enhanced by using machine learning techniques to perform imputations and adjustments on deficient data [97]. The scope of rule-based methods can be similarly expanded to cover other competency questions such as non-redundancy (*conciseness*) through the proposal of other logical constraints to tackle them such as a condition to find taxonomic relations to trim in a knowledge graph (Examples can be found at [https://www.wikidata.org/wiki/Wikidata:Database\\_evaluation](https://www.wikidata.org/wiki/Wikidata:Database_evaluation)). The main limitation of applying the logical constraints using SPARQL in the context of Wikidata is that the runtime of a query that infers or verifies a complex condition or that analyzes a huge amount of class items or property use cases can exceed the timeout limit of the used endpoint [34].

These evaluation assignments covered by our approach can be done by other rule-based (*structure-based* and *semantic-based*) ontology evaluation methods. Structure-based methods verify if a knowledge graph is defined according to a set of formatting constraints and semantic-based methods check if concepts and statements of a knowledge graph meet logical conditions [14]. Some of these methods are software tools, particularly Protégé extensions such as OWLET [98] and OntoCheck [99]. OWLET infers the JSON schema logics of a given knowledge graph, converts them into OWL-DL axioms, and uses the semantic rules to validate the assessed ontological data [98]. OntoCheck screens an ontology to identify structural conventions and constraints for the definition of the analyzed relational information and consequently to homogenize the data structure and quality of the ontology by eliminating typos and pattern violations [99]. Here, the advantage of applying constraints using SPARQL is that its runtime is faster, as it does not require the download of the full dumps of the evaluated knowledge graph [34]. The benefit of our method and other structure-based and semantic-based web-based tools for knowledge graph validation like OntoKeeper [95] and adviseEditor [100] when compared to software tools is that the maximal size of the knowledge graphs that can be assessed by web services is larger than the one that can be evaluated by software tools because the latter depends on the requirements and capacities of the host computer [98, 99]. It is true that these drawbacks of other structure-

based tools can be solved through the simplification of the knowledge graph by reducing redundancies using techniques like ontology trimming [101] or through the construction of an abstraction network to decrease the complexity of the analyzed knowledge graph [14, 102]. However, knowledge graph simplification processes are time-consuming and resulting time gain can consequently be insignificant [14, 101-102].

Such tasks can be also solved using data-driven ontology evaluation methods. These techniques process texts in natural languages to validate the concepts and statements of a knowledge graph and currently include intrinsic (*lexical-based*) and extrinsic (*cross-validation*, *big data-based* and *corpus-based*) methods [14].

Lexical-based methods use rules implemented in SQL or SPARQL to retrieve items and glosses corresponding to a concept and their semantic relations (mostly *subclass of* statements) [103, 104]. These items are then compared against a second set of rules to identify inconsistencies in their labels, descriptions or semantic relations [14]. The output can then be analyzed using natural language processing techniques such as hamming distance measures [104], semantic annotation tools [103] and semantic similarity measures [14] to comparatively identify deficiencies in the semantic representation, labelling and symmetry of the assessed knowledge graph.

Conversely, extrinsic data-based methods extract the usage and linguistic patterns from raw text corporuses such as bibliographic databases and clinical records (*Corpus-based methods*) or from gold standard semantic resources like large ontologies and knowledge graphs (*Cross-validation methods*) or from social media posts and interactions, Internet of Things data or web service statistics (*Big data-based methods*) [14, 105-107] using structure-based and semantic-based ontology evaluation methods as explained above [106] as well as a range of techniques including machine learning [58, 108], topic modelling using latent dirichlet analysis [109], word embeddings [53], statistical correlations [110] and semantic annotation methods [111]. The returned features of the analyzed resources are compared to the ones of the analyzed knowledge graph to assess the accuracy and completeness of the definition and use of concepts and properties [14].

When compared to our proposed approach, lexical-based methods have the advantage to identify and adjust characteristics of a knowledge graph item based on its natural language information of a

knowledge graph item, particularly terms and glosses [103, 104]. The drawbacks of using semantic similarity, word embeddings and topic modelling approaches in such approaches is that these techniques are sensitive to the used parameters, to input characteristics and to the chosen models of computation and can consequently give different results according to the context of determination [56, 57]. The current role of constraints in the extraction of lexical information and respective semantic relations [103, 104] proves that the scope of constraint-based validation should not only be restricted to rule-based evaluation but also to lexical-based evaluation. Yet, the function of logical conditions should be expanded to refine the list of pairs (lexical information, semantic relation) to more accurately identify deficient and missing semantic relations and defective lexical data and to support multilingual lexical-based methods. This would build on the many SPARQL functions that analyze strings in knowledge graphs<sup>24</sup> such as STRLEN (length of a string), STRSTARTS (verification of a substring beginning a given string), STREND (verification of a substring finishing a given string), and CONTAINS (verification of a substring included in a given string) [81, 82].

As for the extrinsic data-driven methods, they are mainly based on large-scale resources that are regularly curated and enriched. Raw-text corporuses are mainly composed of scholarly publications [21] and blog posts [112]. Information in scholarly publications is ever-changing according to the dynamic advances in scholarly knowledge, particularly medical data [113]. This expansion of scientific information in scholarly publications is highly recognized in the context of COVID-19 where detailed information about COVID-19 disease and the SARS-CoV-2 virus is published within less than six months [9]. Big data is the set of real-time statistical and textual information that is generated by web services including search engines and social media and by Internet of Things objects including sensors [105]. This data is characterized by its value, variety, variability, velocity, veracity and volume [105] and can be consequently used to track the changes of the community knowledge and consciousness over time [109, 114]. Large semantic resources are ontologies and knowledge graphs that are built and curated by a community of specialists and that are regularly verified, updated and enriched

<sup>24</sup> Detailed information about string functions in SPARQL can be found at <https://www.w3.org/TR/sparql11-query/#func-strings>.

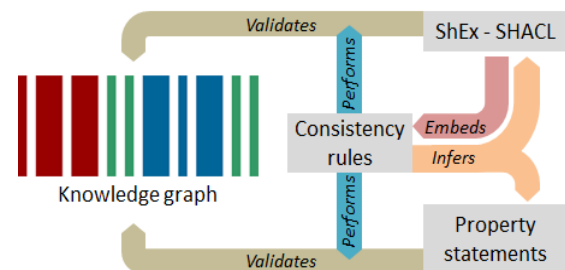


using human efforts and computer programs [115]. These resources represent broad and reliable information about a given specialty through machine learning techniques [58] and the crowdsourcing of scientific efforts [85] and can be consequently compared to other semantic databases for validation purposes. Examples of these resources are COVID-19 Disease Map [8] and SNOMED-CT<sup>25</sup> [115].

Large-scale knowledge graphs are dynamic corpuses. Changes in the logical and semantic conditions for the definition of knowledge in a particular domain need to be identified to adjust the assessed knowledge graph accordingly. Rule-based and lexical-based approaches (especially constraints-based methods) are therefore less simple to apply than extrinsic data-driven methods [14]. Nonetheless, the growing and changing nature of gold standard resources require continuous human efforts and an advanced software architecture to maintain (e.g. structure-based and semantic-based methods), process (e.g. word embeddings and latent dirichlet analysis) and store (e.g. Hadoop and MapReduce) these reference resources [85, 105, 115]. This architecture has advanced hardware requirements and its results are subject to change according to used parameters [105].

These tasks are in line with the usage of Shape Expressions as well as property constraints and relations for the validation of data quality and completeness of the semantic information of class items in knowledge graphs as shown in the “Knowledge graph validation of Wikidata” section. A ShEx ShapeMap is a pair of a triple pattern for selecting entities to validate and a shape against which to validate them. This allows for the definition of the properties to be used for the items of a given class [12, 65] and property constraints and relations based on the meta-ontology (i.e. data skeleton) of Wikidata. Expressions written in shape-based property usage validation languages for RDF (e.g. SHACL) can be used to state conditions and formatting restrictions for the usage of relational and non-relational properties [13, 69, 107]. SPARQL can be more efficient in inferring such information than the currently existing techniques that screen all the items and statements of a knowledge graph one by one to identify the conditions for the usage of properties (e.g. SQID) mainly because SPARQL is meant to directly extract information according to a pattern without having to evaluate all the conditions against all items of a knowledge graph [64, 70, 84].

The separate execution of value-based constraints is common in the quality control of XML data. Typically, structural constraints are managed by RelaxNG or XML Schemas, while value-based constraints are captured as Schematron. Much as Schematron rules are typically embedded in RelaxNG, the consistency constraints presented above can be embedded in Shape Expressions Semantic Actions or in SHACL-SPARQL as shown in Fig. 11 [116]. These supplement structural schema languages with mechanisms to capture value-based constraints and in doing so, provide context for the enforcement of those constraints. The implementation of value-based constraints shown in the “Constraint-driven heuristics-based validation of epidemiological data” section can likewise be implemented in a shapes language [78]. Parsing the rules in Table 3 and 10 would allow the mechanical generation or augmentation of shapes, providing flexibility for how the rules are expressed while still exploiting the power of shapes languages for validation. More generally, ontology-based and knowledge graph-based software tools have the potential to provide wide data and platform interoperability, and thus their semantic interoperability is relevant for a range of downstream applications such as IoT and WoT technologies [117].



**Fig. 11.** Interactions between consistency rules, property statements and RDF validation languages [Source: <https://w.wiki/ao5>, License: CC BY 4.0]

## 8. Conclusion

In this paper, we investigate how to best assess COVID-19 knowledge in collaborative ontologies and knowledge graphs (particularly Wikidata) using relational and statistical constraints. Collaborative databases produced through the cumulative edits of thousands of users are able to generate huge amounts of structured information [11] but as a result of their entirely uncoordinated development, they often result

<sup>25</sup> Systematized Nomenclature Of Medicine - Clinical Terms

in uneven coverage of crucial information and inconsistent expression of that information. The resulting gaps are a significant problem (false negatives, false positives, reasoning deficiencies, and missing references). Avoiding, identifying, and closing these gaps is therefore of top importance. We presented a standardized methodology for auditing key aspects of data quality and completeness for these resources<sup>26</sup>.

This approach complements and informs shape-based methods for data conformance to community-decided schemas. The SPARQL execution does not require any pre-processing, and is not only restricted to the validation of the representation of a given item according to a reference data model but also to the comparison of the assessed relational and statistical statements. Our method is demonstrated as useful for measuring the overall accuracy and data quality on a subset of Wikidata and is consequently a necessary first step in any pipeline for detecting and fixing issues in collaborative ontologies and knowledge graphs.

This work has shown the state of the knowledge graph as a snapshot in time. Future work will extend this to investigate how the knowledge base evolves as biomedical knowledge is integrated into it over time. This will require incorporating the edit history in the SPARQL endpoint APIs of knowledge graphs [40, 118] to dynamically visualize time-resolved SPARQL queries. We will also couple the information inferred using this method<sup>27</sup> with Shape Expressions and the explicit constraints of relation types to provide a more effective enrichment, refinement, and adjustment of collaborative ontologies and knowledge graphs.

## 9. Author statements

**Data availability:** All the SPARQL queries used in this research work are provided in the appendices. The Internet Archive links of the URLs cited by this paper are made available at [https://web.archive.org/web/20200913162813/https://www.wikidata.org/w/index.php?title=User:Daniel\\_Mietchen/sandbox&oldid=1276710909](https://web.archive.org/web/20200913162813/https://www.wikidata.org/w/index.php?title=User:Daniel_Mietchen/sandbox&oldid=1276710909).

<sup>26</sup> This method can be adapted to meet the needs of the user. For instance, the SPARQL queries can be slightly adjusted to assess other patterns in collaborative ontologies such as the usage of classes.

<sup>27</sup> This information can be represented in the form of RDF triples where the subject is the studied relation type and integrated to Wikidata.

**Conflict of interest:** All the co-authors of this paper except EP are active members of WikiProject Medicine, the community curating clinical knowledge in Wikidata, and of WikiProject COVID-19, the community developing multidisciplinary COVID-19 information in Wikidata. DJ is a non-paid voluntary member of the Board of Trustees of Wikimedia Foundation, the non-profit publisher of Wikipedia and Wikidata. EP is a co-creator of SPARQL. EP and JELG are co-creators of ShEx.

## 10. Acknowledgements

The work done by Houcemeddine Turki, Mohamed Ali Hadj Taieb and Mohamed Ben Aouicha was supported by the Ministry of Higher Education and Scientific Research in Tunisia (MoHESR) in the framework of Federated Research Project PRFCOV19-D1-P1, by Wikimedia Foundation through a rapid grant, and by WikiCred Grants Initiative of Craig Newmark Philanthropies, Facebook, and Microsoft. The work done by Jose Emilio Labra Gayo was partially funded by the Spanish Ministry of Economy and Competitiveness (Society challenges: TIN2017-88877-R). The work done by Daniel Mietchen was supported in part by the Alfred P. Sloan Foundation under grant number G-2019-11458. The work done by Dariusz Jemielniak was funded by Polish National Science Center grant no 2019/35/B/HS6/01056. We thank the Wikidata community, Olivier Corby (Université Côte d’Azur, France), Odile Papini (Aix-Marseille Université, France), Egon Willighagen (Maastricht University, Netherlands), and Mahir Morshed (University of Illinois at Urbana-Champaign, United States of America) for useful comments and discussions about the topic of this research paper.

## References

- [1] L. Krishnan, S. M. Ogunwale, and L. A. Cooper. Historical Insights on Coronavirus Disease 2019 (COVID-19), the 1918 Influenza Pandemic, and Racial Disparities: Illuminating a Path Forward. *Annals of Internal Medicine* (2020). doi:10.7326/M20-2223.
- [2] E. Dong, H. Du, and L. Gardner. An interactive web-based dashboard to track COVID-19 in real time. *The Lancet infectious diseases*, 20.5 (2020), pp. 533-534. doi: 10.1016/S1473-3099(20)30120-1.
- [3] B. Xu, M. U. Kraemer, and Data Curation Group. Open access epidemiological data from the COVID-19 outbreak. *The Lancet Infectious Diseases*, 20.5 (2020), pp. 534. doi:10.1016/S1473-3099(20)30119-5.

- [4] D. L. Heymann. Data sharing and outbreaks: best practice exemplified. *The Lancet*, 395.10223 (2020), pp. 469-470. doi: 10.1016/S0140-6736(20)30184-7.
- [5] D. Mitchen and J. Li. Quantifying the Impact of Data Sharing on Outbreak Dynamics (QIDSOD). *Research Ideas and Outcomes*, 6 (2020), p. e54770. doi: 10.3897/rio.6.e54770.
- [6] RDA COVID-19 Working Group. *RDA COVID-19: recommendations and guidelines, 5th release 28 May 2020*, Research Data Alliance, 2020. doi:10.15497/RDA00046.
- [7] J. Y. Cuan-Baltazar, M. J. Muñoz-Perez, C. Robledo-Vega, M. F. Pérez-Zepeda, and E. Soto-Vega. Misinformation of COVID-19 on the internet: infodemiology study. *JMIR public health and surveillance*, 6.2 (2020), p. e18444. doi:10.2196/18444.
- [8] M. Ostaszewski, A. Mazein, M. E. Gillespie, I. Kuperstein, A. Niarakis, H. Hermjakob, et al. COVID-19 Disease Map, building a computational repository of SARS-CoV-2 virus-host interaction mechanisms. *Scientific data*, 7.1 (2020), p. 136. doi:10.1038/s41597-020-0477-8.
- [9] D. Kagan, J. Moran-Gilad, and M. Fire. Scientometric trends for coronaviruses and other emerging viral infections. *GigaScience*, 9.8 (2020), p. g1aa085. doi:10.1093/gigascience/giaa085.
- [10] D. Vrandečić and M. Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57.10 (2014), pp. 78-85. doi:10.1145/2629489.
- [11] H. Turki, T. Shafee, M. A. Hadj Taieb, M. Ben Aouicha, D. Vrandečić, D. Das, and H. Hamdi. Wikidata: A large-scale collaborative ontological medical database. *Journal of Biomedical Informatics*, 99 (2019), p. 103292. doi:10.1016/j.jbi.2019.103292.
- [12] A. Waagmeester, E. L. Willighagen, A. I. Su, M. Kutmon, J. E. Labra Gayo, D. Fernández-Álvarez, et al. A protocol for adding knowledge to Wikidata, a case report. *BioRxiv* (2020). doi:10.1101/2020.04.05.026336.
- [13] F. Erxleben, M. Günther, M. Krötzsch, J. Mendez, and D. Vrandečić (2014). Introducing Wikidata to the linked data web, in: *International semantic web conference*, Springer, Cham., 2014, pp. 50-65. doi:10.1007/978-3-319-11964-9\_4.
- [14] M. Amith, Z. He, J. Bian, J. A. Lossio-Ventura, and C. Tao. Assessing the practice of biomedical ontology evaluation: Gaps and opportunities. *Journal of Biomedical Informatics*, 80 (2018), pp. 1-13. doi:10.1016/j.jbi.2018.02.010.
- [15] J. Brank, M. Grobelnik, and D. Mladenic. A survey of ontology evaluation techniques, in: *Proceedings of the Conference on Data Mining and Data Warehouses (SiKDD 2005)*, Ljubljana, Slovenia: Citeseer, 2005, pp. 166-170.
- [16] A. Lozano-Tello and A. Gomez-Perez. Ontometric: A Method to Choose the Appropriate Ontology. *Journal of Database Management (JDM)*, 15.2 (2004), pp. 1-18.
- [17] A. Degbelo. A Snapshot of Ontology Evaluation Criteria and Strategies, in: *Proceedings of the 13th International Conference on Semantic Systems*, New York: ACM, 2017, pp. 1-8. doi:10.1145/3132218.3132219.
- [18] H. Paulheim. Knowledge graph refinement: A survey of approaches and evaluation methods. *Semantic Web*, 8.3 (2017), pp. 489-508. doi:10.3233/SW-160218.
- [19] D. Vrandečić. Ontology Evaluation, in: R. S. S. Staab, *Handbook on Ontologies*, Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 293-313. doi:10.1007/978-3-540-92673-3\_13.
- [20] L. Obrst, W. Ceusters, I. Mani, S. Ray, and B. Smith. The Evaluation of Ontologies, in: *Semantic Web*, Boston, MA: Springer, 2007, pp. 139-158. doi:10.1007/978-0-387-48438-9\_8.
- [21] J. Raad and C. Cruz, C. A survey on ontology evaluation methods, in: *Proceedings of the International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management*, ACM, 2015, pp. 179-186. doi:10.5220/0005591001790186.
- [22] S. Burgstaller-Muehlbacher, A. Waagmeester, E. Mitraka, J. Turner, T. Putman, J. Leong, et al. Wikidata as a semantic framework for the Gene Wiki initiative. *Database*, 2016 (2016), p. baw015. doi:10.1093/database/baw015.
- [23] E. Mitraka, A. Waagmeester, S. Burgstaller-Muehlbacher, L. M. Schriml, A. I. Su, and B. M. Good. Wikidata: A platform for data integration and dissemination for the life sciences and beyond. *bioRxiv* (2015), p. 031971. doi:10.1101/031971.
- [24] D. Mitchen, G. Hagedorn, E. Willighagen, M. Rico, A. Gómez-Pérez, E. Aibar, K. Rafes, C. Germain, A. Dunning, L. Pintscher, and D. Kinzler. Enabling open science: Wikidata for research (Wiki4R). *Research Ideas and Outcomes*, 1 (2015), p. e7573. doi:10.3897/rio.1.e7573.
- [25] A. Waagmeester, L. Schriml, and A. I. Su. Wikidata as a linked-data hub for Biodiversity data. *Biodiversity Information Science and Standards*, 3 (2019), p. e35206. doi:10.3897/biss.3.35206.
- [26] H. Turki, D. Vrandečić, H. Hamdi, and I. Adel. Using WikiData as a Multi-lingual Multi-dialectal Dictionary for Arabic Dialects, in: *2017 IEEE/ACS 14th International Conference on Computer Systems and Applications (AICCSA)*, IEEE, 2017, pp. 437-442. doi:10.1109/AICCSA.2017.115.
- [27] S. Wasi, M. Sachan, and M. Darbari. Document Classification Using Wikidata Properties, in: *Information and Communication Technology for Sustainable Development*, Singapore: Springer, 2020, pp. 729-737. doi:10.1007/978-981-13-7166-0\_73.
- [28] A. Heftberger, J. Höper, C. Müller-Birn, and N.-O. Walkowski. Opening up Research Data in Film Studies by Using the Structured Knowledge Base Wikidata, in: H. Kremers, *Digital Cultural Heritage*, Springer International Publishing, 2020, pp. 401-410. doi:10.1007/978-3-030-15200-0\_27.
- [29] M. Mora-Cantalops, S. Sánchez-Alonso, and E. García-Barriocanal. A systematic literature review on Wikidata. *Data Technologies and Applications*, 53 (2019), pp. 250-268. doi:10.1108/DTA-12-2018-0110.
- [30] C. Müller-Birn, B. Karran, J. Lehmann, and M. Luczak-Rösch, M. Peer-production System or Collaborative Ontology Engineering Effort: What is Wikidata?, in: *Proceedings of the 11th International Symposium on Open Collaboration*, New York: ACM, 2015, pp. 20:1-20:10. doi:10.1145/2788993.2789836.
- [31] M. Miquel-Ribé and D. Laniado. Wikipedia Culture Gap: Quantifying Content Imbalances Across 40 Language Editions. *Frontiers in Physics*, 6 (2018), p. 54. doi:10.3389/fphy.2018.00054.
- [32] Kaffee, L. A., Piscopo, A., Vougiouklis, P., Simperl, E., Carr, L., & Pintscher, L. (2017). A glimpse into babel: An analysis of multilinguality in Wikidata, in: *Proceedings of the 13th International Symposium on Open Collaboration*, ACM, 2017, p. 14. doi:10.1145/3125433.3125465.
- [33] D. Jemielniak and M. Wilamowski. Cultural diversity of quality of information on Wikipedias. *Journal of the Association for Information Science and Technology*, 68.10 (2017), pp. 2460-2470. doi:10.1002/asi.23901.
- [34] S. Malyshev, M. Krötzsch, L. González, J. Gonsior, and A. Bielefeldt. Getting the most out of Wikidata: Semantic technology usage in wikipedia's knowledge graph, in: *International Semantic Web Conference*, Springer, Cham., 2018, pp. 376-394. doi:10.1007/978-3-030-00668-6\_23.

- [35] F. Å. Nielsen, D. Mietchen, and E. Willighagen. Scholia, scientometrics and Wikidata, in: *European Semantic Web Conference*, Springer, Cham., 2017, pp. 237-259. doi:10.1007/978-3-319-70407-4\_36.
- [36] D. Jemielniak. *Common knowledge?: An ethnography of Wikipedia*, Stanford: Stanford University Press, 2014. ISBN:978-0804789448.
- [37] A. Piscopo and E. Simperl. Who Models the World?: Collaborative Ontology Creation and User Roles in Wikidata. *Proceedings of the ACM on Human-Computer Interaction*, 2. CSCW (2018), pp. 141:1–141:18. doi:10.1145/3274410.
- [38] J. Samuel. Collaborative Approach to Developing a Multilingual Ontology: A Case Study of Wikidata, in: *Research Conference on Metadata and Semantics Research*, Springer, 2017, pp. 167–172. doi:10.1007/978-3-319-70863-8\_16.
- [39] M. Luggen, D. Difallah, C. Sarasua, G. Demartini, and P. Cudré-Mauroux. Non-parametric Class Completeness Estimators for Collaborative Knowledge Graphs—The Case of Wikidata, in: *The Semantic Web – ISWC 2019*, Springer International Publishing, 2019, pp. 453–469. doi:10.1007/978-3-030-30793-6\_26.
- [40] T. Pellissier Tanon and F. Suchanek. Querying the Edit History of Wikidata, in: *The Semantic Web: ESWC 2019 Satellite Events*, Springer International Publishing, 2019, pp. 161–166. doi:978-3-030-32327-1\_32.
- [41] A. Sarabadani, A. Halfaker, and D. Taraborelli. Building automated vandalism detection tools for Wikidata, in: *Proceedings of the 26th International Conference on World Wide Web Companion*, ACM, 2017, pp. 1647-1654.
- [42] M. Färber, F. Bartscherer, C. Menne, and A. Rettinger. Linked data quality of DBpedia, Freebase, OpenCyc, Wikidata, and YAGO. *Semantic Web*, 9.1 (2018), pp. 77–129. doi:10.3233/SW-170275.
- [43] S. Pillai, L.-K. Soon, and S.-C. Haw. Comparing DBpedia, Wikidata, and YAGO for Web Information Retrieval, in: *Intelligent and Interactive Computing*, Springer Singapore, 2019, pp. 525–535. doi:10.1007/978-981-13-6031-2\_40.
- [44] T. Shafee, G. Masukume, L. Kipersztok, D. Das, M. Häggström, and J. Heilman. Evolution of Wikipedia’s medical content: past, present and future. *J Epidemiol Community Health*, 71.11 (2017), pp. 1122-1129. doi:10.1136/jech-2016-208601.
- [45] E. Zangerle, W. Gassler, M. Pichl, S. Steinhauser, and G. Specht. An Empirical Evaluation of Property Recommender Systems for Wikidata and Collaborative Knowledge Bases, in: *Proceedings of the 12th International Symposium on Open Collaboration*, New York: ACM, 2016, pp. 18:1–18:8. doi:10.1145/2957792.2957804.
- [46] A. Waagmeester, G. Stupp, S. Burgstaller-Muehlbacher, B. M. Good, G. Malachi, O. L. Griffith, et al. Wikidata as a knowledge graph for the life sciences. *eLife*, 9 (2020), p. e52614. doi:10.7554/eLife.52614.
- [47] A. Lanamäki and J. Lindman. Latent Groups in Online Communities: a Longitudinal Study in Wikipedia. *Computer Supported Cooperative Work (CSCW)*, 27.1 (2018), pp. 77-106. doi:10.1007/s10606-017-9295-8.
- [48] C. Sarasua, A. Checco, G. Demartini, D. Difallah, M. Feldman, and L. Pintscher. The evolution of power and standard Wikidata editors: comparing editing behavior over time to predict lifespan and volume of edits. *Computer Supported Cooperative Work (CSCW)*, 28.5 (2019), pp. 843-882. doi:10.1007/s10606-018-9344-y.
- [49] D. Jemielniak and A. Przegalińska. *Collaborative Society*, Cambridge, MA: MIT Press, 2020. ISBN:978-0262537919.
- [50] D. Vrandečić. Architecture for a multilingual Wikipedia. *arXiv preprint arXiv:2004.04733* (2020).
- [51] L.-A. Kaffee and E. Simperl. Analysis of Editors’ Languages in Wikidata, in: *Proceedings of the 14th International Symposium on Open Collaboration*, ACM, 2018, p. 21. doi:10.1145/3233391.3233965.
- [52] M. Farda-Sarbas, H. Zhu, M. F. Nest, and C. Müller-Birn. Approving automation: analyzing requests for permissions of bots in Wikidata, in: *Proceedings of the 15th International Symposium on Open Collaboration*, 2019, pp. 1-10. doi:10.1145/3306446.3340833.
- [53] Y. Zhang, Q. Chen, Z. Yang, H. Lin, and Z. Lu. BioWordVec, improving biomedical word embeddings with subword information and MeSH. *Scientific data*, 6.1 (2019), pp. 52:1-52:9. doi:10.1038/s41597-019-0055-0.
- [54] Q. Chen, K. Lee, S. Yan, S. Kim, C. H. Wei, and Z. Lu. BioConceptVec: Creating and evaluating literature-based biomedical concept embeddings on a large scale. *PLoS computational biology*, 16.4 (2020), p. e1007617. doi:10.1371/journal.pcbi.1007617.
- [55] M. Ben Aouicha and M. A. Hadj Taieb. Computing semantic similarity between biomedical concepts using new information content approach. *Journal of biomedical informatics*, 59 (2016), pp. 258-275. doi:10.1016/j.jbi.2015.12.007.
- [56] J. J. Lastra-Díaz, J. Goikoetxea, M. A. Hadj Taieb, A. García-Serrano, M. Ben Aouicha, and E. Agirre. A reproducible survey on word embeddings and ontology-based methods for word similarity: linear combinations outperform the state of the art. *Engineering Applications of Artificial Intelligence*, 85 (2019), pp. 645-665. doi:10.1016/j.engappai.2019.07.010.
- [57] M. A. Hadj Taieb, T. Zesch, and M. Ben Aouicha. A survey of semantic relatedness evaluation datasets and procedures. *Artificial Intelligence Review*, 53.6 (2020), pp. 4407-4448. doi:10.1007/s10462-019-09796-3.
- [58] Y. Zhang, H. Lin, Z. Yang, J. Wang, S. Zhang, Y. Sun, and L. Yang. A hybrid model based on neural networks for biomedical relation extraction. *Journal of biomedical informatics*, 81 (2018), 83-92. doi:10.1016/j.jbi.2018.03.011.
- [59] H. Turki. Citation analysis is also useful to assess the eligibility of biomedical research works for inclusion in living systematic reviews. *Journal of clinical epidemiology*, 97 (2018), pp. 124-125. doi:10.1016/j.jclinepi.2017.11.002.
- [60] P. Mayr, A. Scharnhorst, B. Larsen, P. Schaer, and P. Mutschke. Bibliometric-enhanced information retrieval, in: *European Conference on Information Retrieval*, Springer, Cham., 2014, pp. 798-801. doi:10.1007/978-3-319-06028-6\_99.
- [61] H. Turki, M. A. Hadj Taieb, and M. Ben Aouicha. MeSH qualifiers, publication types and relation occurrence frequency are also useful for a better sentence-level extraction of biomedical relations. *Journal of biomedical informatics*, 83 (2018), pp. 217-218. doi:10.1016/j.jbi.2018.05.011.
- [62] A. Piad-Morffis, Y. Gutiérrez, and R. Muñoz. A corpus to support ehealth knowledge discovery technologies. *Journal of biomedical informatics*, 94 (2019), 103172. doi:10.1016/j.jbi.2019.103172.
- [63] T. Pellissier Tanon, C. Bourgaux, and F. Suchanek, F. Learning how to correct a knowledge base from the edit history, in: *The World Wide Web Conference*, ACM, 2019, pp. 1465-1475. doi:10.1145/3308558.3313584.
- [64] T. Hanika, M. Marx, and G. Stumme. Discovering implicational knowledge in Wikidata, in: *International Conference on Formal Concept Analysis*, Springer, Cham., 2019, pp. 315-323. doi:10.1007/978-3-030-21462-3\_21.

- [65] E. Prud'hommeaux, J. E. Labra Gayo, and H. Solbrig. Shape Expressions: An RDF Validation and Transformation Language, in: *Proceedings of the 10th International Conference on Semantic Systems*, ACM, 2014, pp. 32-40. doi:10.1145/2660517.2660523
- [66] H. Knublauch and D. Kontokostas. Shapes Constraint Language (SHACL), W3C Recommendation 20 July 2017. W3C Recommendation, #w3c# (2017). Retrieved from <https://www.w3.org/TR/2017/REC-shacl-20170720/>.
- [67] J. E. Labra Gayo, E. Prud'Hommeaux, I. Boneva, and D. Kontokostas. Validating RDF data. *Synthesis Lectures on Semantic Web: Theory and Technology*, 7.1 (2017), pp. 1-328. doi:10.2200/s00786ed1v01y201707wbe016.
- [68] F. Å. Nielsen, K. Thornton, and J. E. Labra-Gayo. Validating Danish Wikidata lexemes, in: *15th International Conference on Semantic Systems, SEMPDS 2019*, Karlsruhe: CEUR-WS, 2019.
- [69] K. Thornton, H. Solbrig, G. S. Stupp, J. E. Labra Gayo, D. Mietchen, E. Prud'Hommeaux, and A. Waagmeester. Using Shape Expressions (ShEx) to share RDF data models and to guide curation with rigorous validation, in: *European Semantic Web Conference*, Springer, 2019, pp. 606-620. doi:10.1007/978-3-030-21348-0\_39.
- [70] J. Pérez, M. Arenas, and C. Gutierrez. Semantics and complexity of SPARQL. *ACM Transactions on Database Systems (TODS)*, 34.3 (2009), p. 16. doi:10.1145/1567274.1567278.
- [71] D. Wiśniewski, J. Potoniec, A. Ławrynowicz, and C. M. Keet. Analysis of Ontology Competency Questions and their formalizations in SPARQL-OWL. *Journal of Web Semantics*, 59 (2019), 100534. doi:10.1016/j.websem.2019.100534.
- [72] S. Vasanthapriyan, J. Tian, and J. Xiang. An Ontology-Based Knowledge Framework for Software Testing, in: *Knowledge and Systems Sciences*, Springer Singapore, 2017, pp. 212-226. doi:10.1007/978-981-10-6989-5\_18.
- [73] R. Bansal and S. Chawla. Design and development of semantic web-based system for computer science domain-specific information retrieval. *Perspectives in Science*, 8 (2016), pp. 330-333. doi:10.1016/j.pisc.2016.04.067.
- [74] P. A. Martin. Evaluating Ontology Completeness via SPARQL and Relations-between-Classes Based Constraints, in: *11th International Conference on the Quality of Information and Communications Technology (QUATIC)*, IEEE, 2018, pp. 255-263. doi:10.1109/QUATIC.2018.00045.
- [75] Addshore, D. Mietchen, and E. Willighagen. *Wikidata Queries around the SARS-CoV-2 virus and pandemic*, NL: Zenodo, 2020. doi:10.5281/zenodo.3977414.
- [76] J. E. Labra Gayo and J. M. Alvarez Rodríguez. Validating statistical index data represented in RDF using SPARQL queries. *RDF Validation Workshop. Practical Assurances for Quality RDF Data* (2013). Cambridge, MA: World Wide Web Consortium. Retrieved from: <http://www.w3.org/2012/12/rdf-val>.
- [77] A. I. Walisadeera, A. Ginige, and G. N. Wikramanayake. Ontology Evaluation Approaches: A Case Study from Agriculture Domain, in: *Computational Science and Its Applications -- ICCSA 2016*, Springer International Publishing, 2016, pp. 318-333. doi:10.1007/978-3-319-42089-9\_23.
- [78] J. E. Labra-Gayo, H. García-González, D. Fernández-Alvarez, and E. Prud'hommeaux, (2019). Challenges in RDF validation, in: *Current Trends in Semantic Web Technologies: Theory and Practice*, Springer, Cham., 2019, pp. 121-151. doi:10.1007/978-3-030-06149-4\_6.
- [79] A. P. Kumar, A. Kumar, and V. N. Kumar. A comprehensive comparative study of SPARQL and SQL. *International Journal of Computer Science and Information Technologies*, 2.4 (2011), pp. 1706-1710.
- [80] A. Bonifati, W. Martens, and T. Timm. An analytical study of large SPARQL query logs. *Proceedings of the VLDB Endowment*, 11.2 (2017), 149-161. doi:10.14778/3149193.3149196.
- [81] B. DuCharme. *Learning SPARQL: querying and updating with SPARQL 1.1*. O'Reilly Media, Inc., 2013. ISBN:978-1449306595.
- [82] S. Harris, A. Seaborne, and E. Prud'hommeaux. SPARQL 1.1 query language. W3C recommendation, 21.10 (2013), p. 778.
- [83] P. Y. Hsu and D. S. Parker. Improving SQL with generalized quantifiers, in: *Proceedings of the Eleventh International Conference on Data Engineering*, IEEE, 1995, pp. 298-305. doi:10.1109/ICDE.1995.380381.
- [84] M. Marx and M. Krötzsch. SQID: Towards Ontological Reasoning for Wikidata, in: *Proceedings of the ISWC 2017 Posters & Demonstrations Track*, CEUR Workshop Proceedings, 2017.
- [85] J. M. Mortensen, E. P. Minty, M. Januszyk, T. E. Sweeney, A. L. Rector, N. F. Noy, and M. A. Musen. Using the wisdom of the crowds to find critical errors in biomedical ontologies: a study of SNOMED CT. *Journal of the American Medical Informatics Association*, 22 (2014), pp. 640-648. doi:10.1136/amiajnl-2014-002901.
- [86] K. J. Rothman, S. Greenland, and T. L. Lash. *Modern epidemiology*, Lippincott Williams & Wilkins, 2008. ISBN:978-1451190052.
- [87] Z. Y. Zu, M. D. Jiang, P. P. Xu, W. Chen, Q. Q. Ni, G. M. Lu, and L. J. Zhang. Coronavirus disease 2019 (COVID-19): a perspective from China. *Radiology* (2020), p. 200490. doi:10.1148/radiol.2020200490.
- [88] P. L. Delamater, E. J. Street, T. F. Leslie, Y. T. Yang, and K. H. Jacobsen. Complexity of the basic reproduction number (R0). *Emerging infectious diseases*, 25.1 (2019), p. 1. doi:10.3201/eid2501.171901.
- [89] J. A. Backer, D. Klinkenberg, and J. Wallinga. Incubation period of 2019 novel coronavirus (2019-nCoV) infections among travellers from Wuhan, China, 20-28 January 2020. *Eurosurveillance*, 25.5 (2020), p. 2000062. doi:10.2807/1560-7917.ES.2020.25.5.2000062.
- [90] Q. Li, X. Guan, P. Wu, X. Wang, L. Zhou, Y. Tong, et al. Early transmission dynamics in Wuhan, China, of novel coronavirus-infected pneumonia. *New England Journal of Medicine*, 382 (2020), pp. 1199-1207. doi:10.1056/NEJMoa2001316.
- [91] D. Altman, D. Machin, T. Bryant, and M. Gardner. *Statistics with confidence: confidence intervals and statistical guidelines*, John Wiley & Sons, 2013. ISBN:978-0-727-91375-3.
- [92] A. Ward and M. Murray-Ward. *Assessment in the classroom*, Wadsworth Publishing Company, 1999. ISBN:978-0534527044.
- [93] J. Bolleman, E. de Castro, D. Baratin, S. Gehant, B. A. Cuche, A. H. Auchincloss, et al. HAMAP as SPARQL rules—A portable annotation pipeline for genomes and proteomes. *GigaScience*, 9.2 (2020), p. giaa003. doi:10.1093/gigascience/giaa003.
- [94] F. Darari, W. Nutt, S. Razniewski, and S. Rudolph. Completeness and soundness guarantees for conjunctive SPARQL queries over RDF data sources with completeness statements. *Semantic Web*, 11.3 (2020), pp. 441-482. doi:10.3233/SW-190344.
- [95] M. Amith, F. Manion, C. Liang, M. Harris, D. Wang, Y. He, and C. Tao. Architecture and usability of OntoKeeper, an ontology evaluation tool. *BMC medical informatics and*



- decision making, 19.4 (2019), p. 152. doi:10.1186/s12911-019-0859-z.
- [96] G. Q. Zhang and O. Bodenreider. Large-scale, exhaustive lattice-based structural auditing of SNOMED CT, in: *AMIA Annual Symposium Proceedings*, American Medical Informatics Association, 2010, Vol. 2010, p. 922.
- [97] S. Bischof, A. Harth, B. Kämpgen, A. Polleres, and P. Schneider. Enriching integrated statistical open city data by combining equational knowledge and missing value imputation. *Journal of Web Semantics*, 48 (2018), pp. 22-47. doi:10.1016/j.websem.2017.09.003.
- [98] T. J. Lampoltshammer, and T. Heistracher. Ontology evaluation with Protégé using OWLET. *Infocommunications Journal*, 6.2 (2014), pp. 12-17.
- [99] D. Schober, I. Tudose, V. Svatek, and M. Boeker. OntoCheck: verifying ontology naming conventions and metadata completeness in Protégé 4. *Journal of Biomedical Semantics*, 3.Suppl 2 (2012), p. S4. doi:10.1186/2041-1480-3-S2-S4.
- [100] J. Geller, Z. He, Y. Perl, C. P. Morrey, and J. Xu. Rule-based support system for multiple UMLS semantic type assignments. *Journal of biomedical informatics*, 46.1 (2013), pp. 97-110.
- [101] S. G. Jantzen, B. J. Sutherland, D. R. Minkley, and B. F. Koop. GO Trimming: Systematically reducing redundancy in large Gene Ontology datasets. *BMC research notes*, 4.1 (2011), p. 267. doi:10.1186/1756-0500-4-267.
- [102] M. Halper, H. Gu, Y. Perl, and C. Ochs. Abstraction networks for terminologies: supporting management of “big knowledge”. *Artificial intelligence in medicine*, 64.1 (2015), pp. 1-16. doi:10.1016/j.artmed.2015.03.005.
- [103] A. Rector and L. Iannone. Lexically suggest, logically define: Quality assurance of the use of qualifiers and expected results of post-coordination in SNOMED CT. *Journal of biomedical informatics*, 45.2 (2012), pp. 199-209. doi:10.1016/j.jbi.2011.10.002.
- [104] L. Luo, J. L. Mejino Jr, and G. Q. Zhang. An analysis of FMA using structural self-bisimilarity. *Journal of biomedical informatics*, 46.3 (2013), pp. 497-505. doi:10.1016/j.jbi.2013.03.005.
- [105] H. Sebei, M. A. Hadj Taieb, M. Ben Aouicha. Review of social media analytics process and big data pipeline. *Social Network Analysis and Mining*, 8.1 (2018), p. 30. doi:10.1007/s13278-018-0507-0.
- [106] A. L. Rector, S. Brandt, and T. Schneider. Getting the foot out of the pelvis: modeling problems affecting use of SNOMED CT hierarchies in practical applications. *Journal of the American Medical Informatics Association*, 18.4 (2011), pp. 432-440. doi:10.1136/amiajnl-2010-000045.
- [107] A. Gangemi, C. Catenacci, M. Ciaramita, and J. Lehmann. A theoretical framework for ontology evaluation and validation, in: *SWAP*, 2005, Vol. 166, p. 16.
- [108] D. M. Bean, H. Wu, E. Iqbal, O. Dzahini, Z. Ibrahim, M. Broadbent, et al. Knowledge graph prediction of unknown adverse drug reactions and validation in electronic health records. *Scientific reports*, 7.1 (2017), pp. 1-11. doi:10.1038/s41598-017-16674-x.
- [109] A. Abd-Alrazaq, D. Alhuwail, M. Househ, M. Hamdi, and Z. Shah. Top concerns of tweeters during the COVID-19 pandemic: infoveillance study. *Journal of medical Internet research*, 22.4 (2020), p. e19016. doi:10.2196/19016.
- [110] D. Vanderkam, R. Schonberger, H. Rowley, and S. Kumar. *Nearest neighbor search in google correlate*, Google Inc., 2013.
- [111] J. Li, Y. Sun, R. J. Johnson, D. Sciaky, C. H. Wei, R. Leaman, et al. BioCreative V CDR task corpus: a resource for chemical disease relation extraction. *Database*, 2016 (2016). p. baw068. doi:10.1093/database/baw068.
- [112] M. S. Park, Z. He, Z. Chen, S. Oh, and J. Bian. Consumers’ use of UMLS concepts on social media: diabetes-related textual data analysis in blog and social Q&A sites. *JMIR medical informatics*, 4.4 (2016), p. e41.
- [113] M. Jalalifard, Y. Norouzi, and A. Isfandiyari-Moghaddam. Analyzing web citations availability and half-life in medical journals. *Aslib Proceedings*, 65.3 (2013), pp. 242. doi:10.1108/00012531311330638.
- [114] H. Turki, M. A. Hadj Taieb, M. Ben Aouicha, and A. Abraham. Nature or Science: what Google Trends says. *Scientometrics*, 124.2 (2020), pp. 1367-1385. doi:10.1007/s11192-020-03511-8.
- [115] D. Lee, R. Cornet, F. Lau, and N. De Keizer. A survey of SNOMED CT implementations. *Journal of biomedical informatics*, 46.1 (2013), pp. 87-96.
- [116] A. Melo and H. Paulheim. Automatic detection of relation assertion errors and induction of relation constraints. *Semantic Web*, 11.5 (2020), pp. 801-830. doi:10.3233/SW-200369.
- [117] A. Gyrard, S. K. Datta, and C. Bonnet. A survey and analysis of ontology-based software tools for semantic interoperability in IoT and WoT landscapes, in: *2018 IEEE 4th World Forum on Internet of Things (WF-IoT)*, IEEE, 2018, pp. 86-91. doi:10.1109/WF-IoT.2018.8355091.
- [118] J. C. Dos Reis, C. Pruski, M. Da Silveira, and C. Reynaud-Delaître. Understanding semantic mapping evolution by observing changes in biomedical ontologies. *Journal of biomedical informatics*, 47 (2014), pp. 71-82. doi:10.1016/j.jbi.2013.09.006

## 11. Supplementary Data

### 11.1. Appendix A: SPARQL queries for the inference of the usage constraints of relation types in Wikidata

Task	SPARQL query
T1	<pre> SELECT ?Cs ?CsLabel ?Co ?CoLabel (COUNT(*) AS ?count) WHERE {   ?S wdt:&lt;PropertyID&gt; ?O.   ?S wdt:P31 ?Cs.   ?O wdt:P31 ?Co.   FILTER NOT EXISTS {     { ?Cs wdt:P31 wd:Q24017414 } UNION { ?Co wdt:P31 wd:Q24017414 }.   }   SERVICE wikibase:label { bd:serviceParam wikibase:language "en". } } GROUP BY ?Cs ?CsLabel ?Co ?CoLabel ORDER BY DESC(?count) LIMIT 5 </pre>
T2	<pre> SELECT ?Cs ?CsLabel ?P1 ?Co ?CoLabel (COUNT(*) AS ?count) WHERE {   ?S wdt:&lt;PropertyID&gt; ?O.   ?O ?P1 ?S.   ?S wdt:P31 ?Cs.   ?O wdt:P31 ?Co.   FILTER NOT EXISTS {     { ?Cs wdt:P31 wd:Q24017414 } UNION { ?Co wdt:P31 wd:Q24017414 }.   }   SERVICE wikibase:label { bd:serviceParam wikibase:language "en". } } GROUP BY ?Cs ?CsLabel ?P1 ?Co ?CoLabel ORDER BY DESC(?count) LIMIT 5 </pre>
T3	<pre> SELECT ?S ?SLabel ?O ?OLabel ?statement ?p ?ref ?refLabel WHERE {   ?S p:&lt;PropertyID&gt; ?statement.   ?statement ps:&lt;PropertyID&gt; ?O   ?S wdt:&lt;PropertyID&gt; ?O.   ?S wdt:P31 wd:&lt;SubjectID&gt;.   ?O wdt:P31 wd:&lt;ObjectID&gt;.   ?statement prov:wasDerivedFrom [?p ?ref] .   FILTER NOT EXISTS { ?O ?P1 ?S. }   SERVICE wikibase:label { bd:serviceParam wikibase:language "en". } } </pre>
T4	<p><b>G1:</b> Statements where the subject is not an instance of the most used subject class:</p> <pre> SELECT ?S ?SLabel ?O ?OLabel WHERE {   ?S wdt:&lt;PropertyID&gt; ?O. </pre>

	<pre> FILTER NOT EXISTS {   ?S wdt:P31* wd:&lt;SubjectID&gt;. } SERVICE wikibase:label { bd:serviceParam wikibase:language "en". } } </pre> <p><b>G2:</b> Statements where the object is not an instance of the most used object class:</p> <pre> SELECT ?S ?SLabel ?O ?OLabel WHERE {   ?S wdt:&lt;PropertyID&gt; ?O.   FILTER NOT EXISTS {     ?O wdt:P31* wd:&lt;ObjectID&gt;.   }   SERVICE wikibase:label { bd:serviceParam wikibase:language "en". } } </pre>
T5	<pre> SELECT ?p (COUNT(?p) AS ?count) WHERE {   ?S p:&lt;PropertyID&gt; ?statement.   ?statement ps:&lt;PropertyID&gt; ?O.   ?S wdt:&lt;PropertyID&gt; ?O.   ?statement prov:wasDerivedFrom [?p ?ref] . } GROUP BY ?p ORDER BY DESC(?count) </pre>

### 11.2. Appendix B: SPARQL queries for the heuristics-based validation of epidemiological counts in Wikidata

Task	SPARQL query
V1	<pre> SELECT * WHERE {   ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196].   ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value; pq:P585 ?date].   FILTER(YEAR(?date) &lt; 2019) } </pre>
V2	<pre> SELECT * WHERE {   ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196].   ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value; pq:P459 ?method].   FILTER NOT EXISTS { ?method wdt:P279* wd:Q177719 } } </pre>
V3	<pre> SELECT * WHERE {   ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196].   ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value; pq:P585 ?date].   ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value1; pq:P585 ?date].   FILTER(?value &gt; ?value1)   FILTER(?date - ?date = -1) } </pre>
V4	<pre> SELECT * WHERE {   ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. </pre>

	<pre> ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value; pq:P585 ?date]. ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value1; pq:P585 ?date]. FILTER(?value = ?value1) FILTER(?date - ?datef = -2) FILTER NOT EXISTS { ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value2; pq:P585 ?date]. FILTER(?date = ?datep + 1) } } </pre>
V5	<pre> SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?value; pq:P585 ?date]. FILTER(?value &lt; 0) } </pre>
V6	<pre> SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:P8049 [ps:P8049 ?h; pq:P585 ?date]. ?x p:P1603 [ps:P1603 ?c; pq:P585 ?date]. FILTER(?h &gt; ?c) } </pre>
V7	<pre> SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:P8011 [ps:P8011 ?t; pq:P585 ?date]. ?x p:P1603 [ps:P1603 ?c; pq:P585 ?date]. FILTER(?c &gt;= ?t) } </pre>
V8	<pre> SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:P1603 [ps:P1603 ?c; pq:P585 ?date]. ?x p:P1120 [ps:P1120 ?d; pq:P585 ?date]. FILTER(?c &lt; ?d) } </pre>
V9	<pre> SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:P1603 [ps:P1603 ?c; pq:P585 ?date]. ?x p:P8010 [ps:P8010 ?r; pq:P585 ?date]. FILTER(?c &lt; ?r) } </pre>
V10	<pre> SELECT ?y ?date ((?count - ?c1) AS ?diff) WHERE { SELECT ?y ?c1 ?date (SUM(?c) AS ?count) WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?c; pq:P585 ?date]. ?x wdt:P361 ?y. } } </pre>

	<pre> ?y p:&lt;PropertyID&gt; [ps:&lt;PropertyID&gt; ?c1; pq:P585 ?date]. } GROUP BY ?y ?c1 ?date } ORDER BY DESC(?diff) </pre>
--	---

### 11.3. Appendix C: SPARQL queries for the validation of case fatality rate statements in Wikidata

Task	SPARQL query
M1	<pre> SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:P3457 [ps:P3457 ?value; pq:P585 ?date]. FILTER((?value &gt; 1)    (?value &lt; 0)) } </pre>
M2	<pre> SELECT ?x ?c ?d ?value ?date (ABS(?value - ?d / ?c) &gt; 0.001 AS ?diff) WITH { SELECT ?x { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. } } as %outbreaks WITH { SELECT ?x ?value ?date { INCLUDE %outbreaks. ?x p:P3457 [ps:P3457 ?value; pq:P585 ?date]. } } as %casefatalityrates WITH { SELECT ?x ?d ?date { INCLUDE %outbreaks. ?x p:P1120 [ps:P1120 ?d; pq:P585 ?date]. } } as %deaths WITH { SELECT ?x ?c ?date { INCLUDE %outbreaks. ?x p:P1603 [ps:P1603 ?c; pq:P585 ?date]. } } as %cases WHERE { INCLUDE %casefatalityrates. INCLUDE %deaths. INCLUDE %cases. } ORDER BY DESC(?diff) </pre>
M3	<pre> SELECT ?x ?c ?d ?date ((?d / ?c) AS ?m) WITH { SELECT ?x { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. } } </pre>

	<pre>    } as %outbreaks     WITH {       SELECT ?x ?d ?date {         INCLUDE %outbreaks.         ?x p:P1120 [ps:P1120 ?d; pq:P585 ?date].       }     } as %deaths     WITH {       SELECT ?x ?c ?date {         INCLUDE %outbreaks.         ?x p:P1603 [ps:P1603 ?c; pq:P585 ?date].       }     } as %cases     WHERE {       INCLUDE %deaths. INCLUDE %cases.       FILTER NOT EXISTS {?x p:P3457 [ps:P3457 ?value; pq:P585 ?date].}     }</pre>
--	---