

1 **Using logical constraints to validate statistical**
2 **information about COVID-19 in collaborative**
3 **knowledge graphs: the case of Wikidata**

4

5 Houcemeddine Turki¹, Dariusz Jemielniak², Mohamed Ali Hadj Taieb¹, Jose Emilio Labra
6 Gayo³, Mohamed Ben Aouicha¹, Mus'ab Banat⁴, Thomas Shafee⁵, Eric Prud'Hommeaux⁶, Tiago
7 Lubiana^{7,8}, Diptanshu Das⁹, Daniel Mietchen^{8,10,11}

8

9 ¹ Data Engineering and Semantics Research Unit, Faculty of Sciences of Sfax, University of
10 Sfax, Sfax, Tunisia

11 ² Department of Management in Networked and Digital Societies, Kozminski University,
12 Warsaw, Poland

13 ³ Web Semantics Oviedo (WESO) Research Group, University of Oviedo, Oviedo, Spain

14 ⁴ Faculty of Medicine, Hashemite University, Zarqa, Jordan

15 ⁵ La Trobe University, Melbourne, Victoria, Australia

16 ⁵ Swinburne University of Technology, Melbourne, Victoria, Australia

17 ⁶ World Wide Web Consortium, Cambridge, Massachusetts, United States of America

18 ⁷ Computational Systems Biology Laboratory, University of São Paulo, São Paulo, Brazil

19 ⁸ Ronin Institute, Montclair, New Jersey, United States of America

20 ⁹ Institute of Child Health (ICH), Kolkata, India

21 ⁹ Medica Superspecialty Hospital, Kolkata, India

22 ¹⁰ School of Data Science, University of Virginia, Charlottesville, Virginia, United States of
23 America

24 ¹¹ Biomedical Data & Bioethics, Fraunhofer Institute for Biomedical Engineering, Würzburg,
25 Germany

26

27 Corresponding Author:

28 Daniel Mietchen^{8,10,11}

29 Ronin Institute, 127 Haddon Pl, Montclair, New Jersey 07043, United States of America

30 Email address: daniel.mietchen@roninstitute.org

31 **Abstract**

32 Urgent global research demands real-time dissemination of precise data. Wikidata, a
33 collaborative and openly licensed knowledge graph available in RDF format, provides an ideal
34 forum for exchanging structured data that can be verified and consolidated using validation
35 schemas and bot edits. In this research paper, we catalog an automatable task set necessary to
36 assess and validate the portion of Wikidata relating to the COVID-19 epidemiology. These tasks
37 assess statistical data and are implemented in SPARQL, a query language for semantic
38 databases. We demonstrate the efficiency of our methods for evaluating structured non-relational
39 information on COVID-19 in Wikidata, and its applicability in collaborative ontologies and
40 knowledge graphs more broadly. We show the advantages and limitations of our proposed
41 approach by comparing it to the features of other methods for the validation of linked web data
42 as revealed by previous research.

43

44 **Introduction**

45 Since December 2019, the COVID-19 disease has spread to become a global pandemic. This
46 disease is caused by a zoonotic coronavirus called *SARS-CoV-2* (Severe Acute Respiratory
47 Syndrome CoronaVirus 2) and is characterized by the onset of acute pneumonia and respiratory
48 distress. The global impact, with more than 388 million infections and almost 5.7 million deaths
49 globally (as of February 4, 2022¹), is frequently compared to the 1918 Spanish Flu (Krishnan,
50 Ogunwole, & Cooper, 2020). Emerging mRNA vaccines entail serious distribution and storage
51 challenges, and no therapies are especially effective against late stages of the disease. As with all
52 zoonotic diseases, its abrupt introduction to humans demands an outsized effort for data
53 acquisition, curation, and integration to drive evidence-based medicine, predictive modeling, and
54 public health policy (Dong, Du, & Gardner, 2020; Xu, Kraemer, & Data Curation Group, 2020).

55 Agile data sharing and computer-supported reasoning about the COVID-19 pandemic and
56 SARS-CoV-2 virus allow us to quickly understand more about the disease's epidemiology,
57 pathogenesis, and physiopathology. This understanding can then inform the required clinical,
58 scholarly, and public health measures to fight the condition and handle its nonmedical
59 ramifications (Heymann, 2020; Mietchen & Li, 2020; RDA COVID-19 Working Group, 2020).
60 Consequently, initiatives have rapidly emerged to create datasets, web services, and tools to
61 analyze and visualize COVID-19 data. Examples include Johns Hopkins University's COVID-19
62 dashboard (Dong, Du, & Gardner, 2020) and the Open COVID-19 Data Curation Group's
63 epidemiological data (Xu, Kraemer, & Data Curation Group, 2020). Some of these resources are
64 interactive and return their results based on combined clinical and epidemiological information,
65 scholarly information, and social network analysis (Cuan-Baltazar, et al., 2020; Ostaszewski, et
66 al., 2020; Kagan, Moran-Gilad, & Fire, 2020). However, a significant shortfall in interoperability

1 ¹ "[COVID-19 Dashboard](#) by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University
2 (JHU)". ArcGIS. Johns Hopkins University. Retrieved 4 February 2022.

67 is common: although these dashboards facilitate examination of their slice of the data, most of
68 them lack general integration with other sites or datasets. The lack of technical support for
69 interoperability is exacerbated by legal restrictions: despite being free to access, the majority of
70 such dashboards are provided under *All Rights Reserved* terms or licenses. Similarly, >84% of
71 the 142,665 COVID-19-related projects on the GitHub repository for computing projects are
72 under *All Rights Reserved*² terms (as of 4 February 2022). Restrictive licensing of data sets and
73 applications severely impedes their dissemination and integration, ultimately undermining their
74 value for the community of users and re-users. For complex and multifaceted phenomena such as
75 the COVID-19 pandemic, there is a particular need for a collaborative, free, machine-readable,
76 interoperable, and open approach to knowledge graphs that integrate the varied data.

77 Wikidata³ just fits the need as a CC0⁴ licensed, large-scale, multilingual knowledge graph
78 used to represent human knowledge in a structured format (Resource Description Framework or
79 RDF) (Vrandečić & Krötzsch, 2014; Turki, et al., 2019). It, therefore, has the advantage of being
80 inherently findable, accessible, interoperable, and reusable, i.e., FAIR (Waagmeester, et al.,
81 2021). It was initially developed in 2012 as an adjunct to Wikipedia, but has grown significantly
82 beyond its initial parameters. As of now, it is a centralized, cross-disciplinary meta-database and
83 knowledge base for storing structured information in a format optimized to be easily read and
84 edited by both machines and humans (Erxleben, Günther, Krötzsch, Mendez, & Vrandečić,
85 2014). Thanks to its flexible representation of facts, Wikidata can be automatically enriched
86 using information retrieved from multiple public domain sources or inferred from synthesized
87 data (Turki, et al., 2019). This database includes a wide variety of pandemic-related information,
88 including clinical knowledge, epidemiology, biomedical research, software development,
89 geographic, demographic, and genetics data. It can consequently be a vital large-scale reference
90 database to support research and medicine during the COVID-19 pandemic (Turki, et al., 2019;
91 Waagmeester, et al., 2021).

92 The key hurdle to overcome for projects such as Wikidata is that several of their features can
93 make them at-risk of inconsistent structure or coverage: 1) collaborative projects use
94 decentralized contributions rather than central oversight, 2) large-scale projects operate at a scale
95 where manual checking is not possible, and 3) interdisciplinary projects regulate the acquisition
96 of data to integrate a wide variety of data sources. To maximize the usability of the data, it is
97 therefore important to minimize inconsistencies in its structure and coverage. As a result,
98 methods of evaluating the existing knowledge graphs and ontologies, integral to knowledge
99 graph maintenance and development, are of crucial importance. Such an evaluation is
100 particularly relevant in the case of collaborative semantic databases, such as Wikidata.

3 2 120,109 of 142,665 as of 4 February 2022: <https://github.com/search?q=covid-19+OR+covid19+OR+coronavirus+OR+cord19+OR+cord-19>

5 3 <https://www.wikidata.org/>

6 4 CC0 is a rights waiver similar to Creative Commons licenses, used to publish material into the public domain. It
7 waives as much copyright as possible within a given jurisdiction. Further information can be found at
8 <https://creativecommons.org/publicdomain/zero/1.0/> .

101 Knowledge graph evaluation is, therefore, necessary to assess the quality, correctness, or
102 completeness of a given knowledge graph against a set of predetermined criteria (Amith, He,
103 Bian, Lossio-Ventura, & Tao, 2018). There are several possible approaches to evaluating a
104 knowledge graph based on external information (so-called extrinsic evaluation), including
105 comparing its structure to a paragon ontology, comparing its coverage to source data, applying it
106 to a test problem and judging the outcomes, and manual expert review of its ontology (Brank,
107 Grobelnik, & Mladenic, 2005). Different systematic approaches have been proposed for the
108 comparison of ontologies and knowledge graphs, including NLP techniques, machine learning,
109 association rule mining, and other methods (Lozano-Tello & Gomez-Perez, 2004; Degbelo,
110 2017; Paulheim, 2017). The criteria for evaluating ontologies typically include *Accuracy*, which
111 determines if definitions, classes, properties, and individual entries in the evaluated ontology are
112 correct; *Completeness*, referring to the scope of coverage of a given knowledge domain in the
113 evaluated ontology; *Adaptability*, determining the range of different anticipated uses of the
114 evaluated ontology (versatility); and *Clarity*, determining the effectiveness of communication of
115 intended meanings of defined terms by the evaluated ontology (Vrandečić, 2009; Obrst,
116 Ceusters, Mani, Ray, & Smith, 2007; Raad & Cruz, 2015; Amith, et al., 2018). However,
117 extrinsic methods are not the only ones that are used for evaluating such a set of criteria.
118 Knowledge graphs can be also assessed through an intrinsic evaluation that assesses the structure
119 of the analyzed knowledge graph thanks to the inference of internal description logics and
120 consistency rules (Amith, et al., 2018).

121 In this research paper, we emphasize the use of intrinsic methods to evaluate knowledge
122 graphs by presenting our approach to quality assurance checks and corrections of statistical
123 semantic data in Wikidata, mainly in the context of COVID-19 epidemiological information.
124 This consists of a catalog of automatable tasks based on logical constraints expected of the
125 knowledge graph. Most of these constraints were not explicitly available in the RDF validation
126 resources of Wikidata before the pandemic and are designed in this work to support new types of
127 COVID-19 information in the assessed knowledge graph, particularly epidemiological data. Our
128 approach is built upon the outcomes of previous outbreaks such as the Zika pandemic (Ekins et
129 al., 2015) and aims to pave the way towards handling future outbreaks. We implement these
130 constraints with SPARQL and test them on Wikidata using the SPARQL endpoint of this
131 knowledge graph, available at <https://query.wikidata.org>. SPARQL was officially created in 2008
132 as a query language and protocol to search, add, modify or delete RDF data available over the
133 Internet. Its name is a recursive acronym that stands for "SPARQL Protocol and RDF Query
134 Language". SPARQL⁵ allows a query to be composed of triple patterns, conjunctions,
135 disjunctions, and optional patterns and can consequently be used to retrieve contextualized
136 information from knowledge graphs without having to retrieve and process the ontological
137 database. We introduce the value of Wikidata as a multipurpose collaborative knowledge graph

9 5 An open license SPARQL textbook available in multiple languages can be found at
10 <https://en.wikibooks.org/wiki/SPARQL>.

138 for the flexible and reliable representation (Section 2) and validation (Section 3) of COVID-19
139 knowledge. Furthermore, we cover the use of SPARQL to query this knowledge graph (Section
140 4). Then, we demonstrate how statistical constraints can be implemented using SPARQL and
141 applied to verify epidemiological data related to the COVID-19 pandemic (Section 5). Finally,
142 we compare our constraint-based approach with other RDF validation methods through the
143 analysis of the main outcomes of previous research papers related to knowledge graph validation
144 (Section 6) and conclude future directions (Section 7).

145 **Wikidata as a collaborative knowledge graph**

146 Wikidata currently serves as a semantic framework for a variety of scientific initiatives, such as
147 GeneWiki (Burgstaller-Muehlbacher, et al., 2016), allowing different teams of scholars to upload
148 valuable academic data into a collective and standardized pool. Its versatility and
149 interconnectedness are making it a standard for interdisciplinary data integration and
150 dissemination across fields as diverse as linguistics, information technology, film studies, and
151 medicine (Turki, et al., 2019; Mitraka, et al., 2015; Mietchen, et al., 2015; Waagmeester,
152 Schriml, & Su, 2019, Turki, Vrandečić, Hamdi, & Adel, 2017; Wasi, Sachan, & Darbari, 2020;
153 Heftberger, Höper, Müller-Birn, & Walkowski, 2020), although its popularity and recognition
154 across fields still vary significantly (Mora-Cantalops, et al., 2019). It contains concepts, linked
155 by their taxonomic relations, allowing embedding and creating instances of subclasses of
156 classified data and links between them. Its multilingual nature enables fast-updating dynamic
157 data reuse across different language versions of a resource such as Wikipedia (Müller-Birn,
158 Karran, Lehmann, & Luczak-Rösch, 2015), with fewer inconsistencies from local culture
159 (Miquel-Ribé & Laniado, 2018) or language biases (Kaffee, et al., 2017; Jemielniak &
160 Wilamowski, 2017).

161 The data structure employed by Wikidata is intended to be highly standardized, whilst
162 maintaining the flexibility to be applied across highly diverse use-cases. There are mainly two
163 essential components: Items, which represent objects, concepts, or topics; and properties, which
164 describe how one item relates to another. A statement, therefore, consists of a subject item (*S*), a
165 property that describes the nature of the statement (*P*), and an object (*O*) that can be an item, a
166 value, an external ID, or a string, etc. While items can be freely created, new properties require
167 community discussion and vote, with about 9500 properties⁶ currently available. Statements can
168 be further modified by any number of qualifiers to make them more specific, and be supported by
169 references to indicate the source of the information. Thus, Wikidata forms a continuously
170 growing, single, unified network graph, with 96M items forming the nodes, and 1360M
171 statements forming the edges. A live SPARQL endpoint and query service, regular RDF dumps,
172 as well as linked data APIs and visualization tools, establish a backbone of Wikidata uses
173 (Malyshev, Krötzsch, González, Gonsior, & Bielefeldt, 2018; Nielsen, Mietchen, & Willighagen,
174 2017).

11 ⁶ For an updated list of available Wikidata properties, please see <https://tools.wmflabs.org/hay/propbrowse/>.

175 Importantly, Wikidata is based on free and open-source philosophy and software and is a
176 database that anyone can edit, similarly to the very popular online encyclopedia, Wikipedia
177 (Jemielniak, 2014). As a result, the emerging ontologies are created entirely collaboratively,
178 without centralized coordination (Piscopo & Simperl, 2018), and developed in a community-
179 driven fashion (Samuel, 2017). This approach allows for the dynamic development of areas of
180 interest for the user community but poses challenges, e.g., to systematize and proportionate class
181 completeness across topics (Luggen, Difallah, Sarasua, Demartini, & Cudré-Mauroux, 2019).
182 Also, since the edit history is available to anyone, tracing human and non-human contributions,
183 as well as detecting and reverting vandalism is available by design and relies on community
184 management (Pellissier Tanon & Suchanek, 2019) as well as on software tools like ORES
185 (Sarabadani, et al., 2017) or the Item Quality Evaluator⁷.

186 Other ontological databases and knowledge graphs exist (Färber, Bartscherer, Menne, &
187 Rettinger, 2018; Pillai, Soon, & Haw, 2019). However, much like the factors that led Wikipedia
188 to rise to be a dominant encyclopedia (Shafee et al., 2017; Jemielniak & Wilamowski, 2017),
189 Wikidata's close connection to Wikimedia volunteer communities and wide readership provided
190 by Wikipedia have quickly given it a competitive edge. The system, therefore, aims to combine
191 the wisdom of the crowds with advanced algorithms. For instance, Wikidata editors are assisted
192 by a property suggesting system, proposing additional properties to be added to entries
193 (Zangerle, Gassler, Pichl, Steinhäuser, & Specht, 2016). Wikidata has subsequently exhibited the
194 highest growth rate of any Wikimedia project and was the first amongst them to pass one billion
195 contributions (Waagmeester, et al., 2020).

196 As a collaborative venture, its governance model is similar to Wikipedia (Lanamäki &
197 Lindman, 2018), but with some important differences. Wide permissions to edit Wikidata are
198 manually granted to approved bots and to Wikimedia accounts that are at least 4 days old and
199 have made at least 50 edits using manual modifications or semi-automated tools for editing
200 Wikidata⁸. These accounts are supervised by a limited number of experienced administrators to
201 prevent misleading editing behaviors (such as vandalism, harassment, and abuse) and to ensure a
202 sustainable consistency of the information provided by Wikidata⁹. As such, Wikidata is highly
203 relevant to the computer-supported collaborative work (CSCW) field, yet the number of studies
204 of Wikidata from this perspective is still very limited (Sarasua et al., 2019). To understand the
205 value of using SPARQL to validate the usage of relation types in collaborative ontologies and
206 knowledge graphs, it is important to understand the main distinctive features of Wikidata as a
207 collaborative project. Much as Wikidata is developed collaboratively by an international
208 community of editors, it is also designed to be language-neutral. As a result, it is quite possible to
209 contribute to Wikidata with only a limited command of English and to effectively collaborate

12 7 <https://item-quality-evaluator.toolforge.org/>

13 8 For an overview of the semi-automated editing tools for Wikidata, please see
14 <https://www.wikidata.org/wiki/Wikidata:Tools> .

15 9 Further information about the rights and governance of users in Wikidata is shown at
16 https://www.wikidata.org/wiki/Wikidata:User_access_levels .

210 whilst sharing no common human language - an aspect unique even in the already rich
211 ecosystem of collaborative projects (Jemielniak & Przegalinska, 2020). It may well be a corner
212 stone towards the creation of other language-independent cooperative knowledge creation
213 initiatives, such as Wikifunctions, which is an abstract, language-agnostic Wikipedia currently
214 developed and based on Wikidata (Vrandečić, 2021).

215 It is also possible to build Wikipedia articles, especially in underrepresented languages, based
216 on Wikidata data only, and create article placeholders to stimulate encyclopedia articles' growth
217 (Kaffee et al., 2018). This stems from combining concepts that are relatively easily inter-
218 translatable between languages (e.g., professions, causes of death, and capitals) with language-
219 agnostic data (e.g., numbers, geographical coordinates, and dates). As a result, Wikidata is a
220 paragon example of not only cross-cultural cooperation but also human-bot collaborative efforts
221 (Piscopo, 2018; Farda-Sarbas, et al., 2019). Given the large-scale crowdsourcing efforts in
222 Wikidata and the use of bots and semi-automated tools to mass edit Wikidata, its current volume
223 is higher than what can be reviewed and curated by administrators manually. It is quite intuitive:
224 as the general number of edits created by bots grows, so grows the number of administrative
225 tasks to be automated. Automation may include simplifying alerts, fully and semi-automated
226 reverts, better user tracking, or automated corrections. However, the creation of automated
227 methods for the verification and validation of the ontological statements it contains is required
228 most.

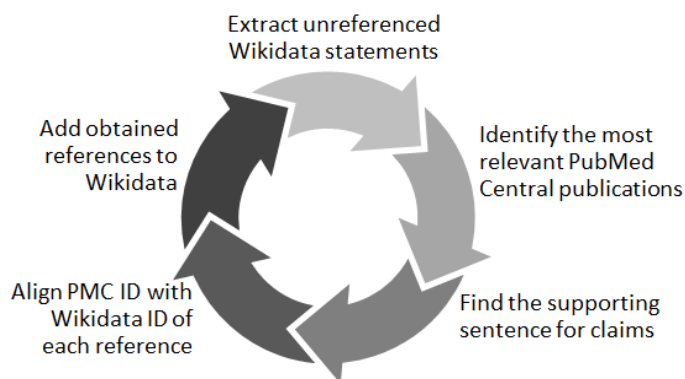
229

230 **Knowledge graph validation of Wikidata**

231 As Wikidata properties are assigned labels, descriptions, and aliases in multiple languages (Red
232 in Fig. 2), multilingual information of these properties can be used alongside the labels,
233 descriptions, and aliases of Wikidata items to verify and find sentences supporting biomedical
234 statements in scholarly outputs (Zhang, et al., 2019). Such a process can be based on various
235 natural language processing techniques, including word embeddings (Zhang, et al., 2019; Chen,
236 et. al., 2020) and semantic similarity (Ben Aouicha & Hadj Taieb, 2016). These techniques are
237 robust enough to achieve an interesting level of accuracy, and some of them can achieve better
238 accuracy when the Wikidata classes of the subject and object of semantic relations are given as
239 inputs (Lastra-Díaz, et al., 2019; Hadj Taieb, Zesch, & Ben Aouicha, 2020). The subjects and
240 objects of Wikidata relations can likewise be aligned to other biomedical semantic resources
241 such as MeSH and UMLS Metathesaurus (Turki, et al., 2019). Thus, benchmarks for relation
242 extraction based on one of the major biomedical ontologies can be converted into a Wikidata
243 friendly format and used to automatically enrich Wikidata with novel biomedical relations or to
244 automatically find statements supporting existing biomedical Wikidata relations (Zhang, et al.,
245 2018). Furthermore, MeSH keywords of scholarly publications can be converted into their
246 Wikidata equivalents, refined using citation and co-citation analysis (Turki, 2018), and used to
247 verify and add biomedical Wikidata relations, e.g., by applying deep learning-based

248 bibliometric-enhanced information retrieval techniques (Mayr, Scharnhorst, Larsen, Schaer, &
249 Mutschke, 2014; Turki, Hadj Taieb, & Ben Aouicha, 2018).

250 Another option of validating biomedical statements based on the labels and external
251 identifiers of their subjects, predicates, and objects in Wikidata can be the use of these labels and
252 external IDs to find whether the assessed Wikidata statements are available in other knowledge
253 resources (e.g., Disease Ontology) and in open bibliographic databases (e.g., PubMed). Several
254 tools have been successfully built using this principle such as the *Wikidata Integrator*¹⁰ that
255 extracts the Wikidata statements of a given gene, protein or cell line using SPARQL, compares
256 them with their equivalents in other structured databases like NCBI's Gene resources, Uniprot or
257 Cellosaurus and adjusts them if needed, *Mismatch Finder*¹¹ that identifies Wikidata statements
258 that are not available in external databases, *Structured Categories*¹² that uses SPARQL to
259 identify how the members of a Wikipedia Category are described using Wikidata statements and
260 to reveal whether a statement is missing or mistakenly edited for the definition of category items
261 (Turki, Hadj Taieb, & Ben Aouicha, 2021), and *RefB*¹³ (Fig. 1) that extracts biomedical Wikidata
262 statements not supported by references using SPARQL and identifies the sentences supporting
263 them in scholarly publications using the PubMed Central search engine and a variety of
264 techniques such as concept proximity analysis.



265 **Figure 1. RrefB workflow.** Process of RefB, a bot that adds scholarly references to biomedical Wikidata statements based on
266 PubMed Central [Source: [https://w.wiki/an\\$](https://w.wiki/an$), License: CC BY 4.0]. The source code of RefB is available at
267 <https://github.com/Data-Engineering-and-Semantics/refb/>.
268

269
270 In addition to their multilingual set of labels and descriptions, Wikidata properties are assigned
271 object types using wikibase:propertyType relations (Blue in Fig. 2). These relations allow the

17 10 Wikidata Integrator is a bot framework for automatically curating genetic information provided by Wikidata
18 (<https://github.com/SuLab/WikidataIntegrator>). For Wikidata bots using this framework, refer to
19 https://www.wikidata.org/wiki/Wikidata:WikiProject_Gene_Wiki#Bot_accounts. The framework has been adapted
20 to various specific contexts, e.g., the curation of cell lines indexed in Cellosaurus, as per <https://github.com/caliphosib/cellosaurus-wikidata-bot>.

21 22 11 https://www.wikidata.org/wiki/Wikidata:Mismatch_Finder

23 12 https://www.wikidata.org/wiki/Wikidata:Structured_Categories

24 13 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)),

25 *Source code* at <https://github.com/Data-Engineering-and-Semantics/refb/>, *Wikidata edits* at

26 [https://www.wikidata.org/wiki/Special:Contributions/RefB_\(WikiCred\)](https://www.wikidata.org/wiki/Special:Contributions/RefB_(WikiCred)).

272 assignment of appropriate objects to statements, so that non-relational statements cannot have a
 273 Wikidata item as an object, while objects of relational statements are not allowed to have data
 274 types like a value or a URL (Vrandečić & Krötzsch, 2014).

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 symptoms headache	edit
	0 references	+ 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, symptom	
	relation: instance or subclass of	
	0 references	+ add reference
		+ add value
	type constraint	edit
	class: physiological condition, fictional medical condition	
	relation: instance or subclass of	
	0 references	+ add reference
		+ add value
	citation needed constraint	edit
	0 references	+ add reference
		+ add value

275
 276 **Figure 2. Example of a Wikidata property and its annotations.** Wikidata page of a clinical property [Source:
 277 <https://w.wiki/aeF>, Derived from: <https://w.wiki/aeG>, License: CC0]. It includes the labels, descriptions, and aliases of the
 278 property in multiple languages (Red), the object data type (Blue), statements where the property is the subject (Green) as well as
 279 property constraints (Brown).

280

281 Just like a Wikidata item, a property can be described by statements (Green in Fig. 2). The
282 predicates of these statements link a property to its class (*instance of* [P31]), to its corresponding
283 Wikidata item (*subject item of this property* [P1629]), to example usages (*Wikidata property*
284 *example* [P1855]), to equivalents in other IRIs¹⁴ (*equivalent property* [P1628]), to Wikimedia
285 categories that track its usage on a given wiki (*property usage tracking category* [P2875]), to its
286 inverse property (*inverse property* [P1696]), or to its proposal discussion (*property proposal*
287 *discussion* [P3254]), etc. These statements can be interesting for various knowledge graph
288 validation purposes. The class, the usage examples, and the proposal discussion of a Wikidata
289 property can be useful through the use of several natural language processing techniques,
290 particularly semantic similarity, to provide several features of the use of the property such as its
291 domain of application (e.g., the subject or object of a statement using a Wikidata property related
292 to medicine should be a medical item) and consequently to eliminate some of the erroneous use
293 by screening the property usage tracking category. The class of the Wikidata item corresponding
294 to the property can be used to identify the field of work of the property and thus flag some
295 inappropriate applications. In addition, the external identifiers of such an item can be used for the
296 verification of biomedical relations by their identification within the semantic annotations of
297 scholarly publications built using the *SAT+R* (Subject, Action, Target, and Relations) model
298 (Piad-Morffis, Gutiérrez, & Muñoz, 2019). The inverse property relations can identify missing
299 Wikidata statements, which are implied by the presence of inverse statements in other Wikidata
300 resources.

301 Despite the importance of these statements defining properties, *property constraint* [P2302]
302 relations (Brown in Fig. 2) are the semantic relations that are primarily used for the validation of
303 the usage of a property. In essence, they define a set of conditions for the use of a property,
304 including several heuristics for the type and format of the subject or the object, information about
305 the characteristics of the property, and several description logics for the usage of the property as
306 shown in Table 1. Property constraints are either manually added by Wikidata users or inferred
307 with high accuracy from the knowledge graph of Wikidata or the history of human changes to
308 Wikidata statements (Pellissier Tanon, Bourgaux, & Suchanek, 2019; Hanika, et al., 2019).

309

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

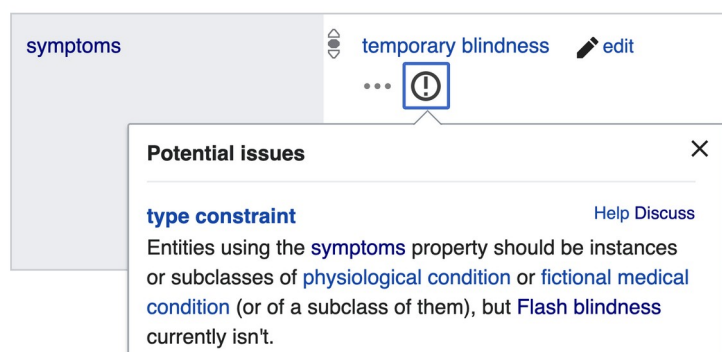
27 14 Internationalized Resource Identifier (IRI) is a standardized character string (e.g., a URL) that recognizes a given
28 item in a semantic resource

		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/meta-levels of its members
Q21510851	allowed qualifiers constraint	Constraint used to specify that only the listed qualifiers should be used. Novalue 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	Constraint used to specify that the subject and the object have to coincide or coexist at some point in 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 a 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

310 **Table 1. Constraint types for the usage of Wikidata properties.** Each property constraint is given with its Wikidata identifier,
311 an English label and an English description.

312
313 As shown in Fig. 2, a property constraint is defined as a relation where the property type is
314 featured as an object and the detailed conditions of the constraint to be applied on Wikidata
315 statements are integrated as qualifiers to the relation. When a property constraint is violated, the
316 corresponding statement is automatically included in a report of property constraint violations¹⁵
317 and is marked by an exclamation mark on the page of the subject item (Fig. 3) so that it can be
318 quickly processed and adjusted by the community or by Wikidata bots if applicable.



319 **Figure 3. Example of a property constraint violation indicated via the Wikidata user interface.** On the page of the Wikidata
320 item Q3603152 (flash blindness), a constraint violation is indicated by the encircled exclamation mark. Clicking on it reveals the
321 display of the popup with some further explanation. [File available on Wikimedia Commons: <https://w.wiki/ZuJ>, License: CC0].
322

323
324 Although these methods are important to verify and validate Wikidata, they are not the only ones that are
325 used for these purposes. In 2019, Wikidata announced the adoption of the Shape Expressions
326 language (ShEx) as part of the Mediawiki entity schemas extension¹⁶. ShEx was proposed
327 following an RDF validation workshop that was organized by W3C¹⁷ in 2014 as a concise, high-
328 level language to describe and validate RDF data (Prud'hommeaux, Labra Gayo, & Solbrig,
329 2014). This Mediawiki extension uses ShEx to store structure definitions (EntitySchemas or
330 Shapes) for sets of Wikidata entities that are selected by some query pattern (frequently the
331 involvement of said entities in a Wikidata class). This provides collaborative quality control

29 15 https://www.wikidata.org/wiki/Wikidata:Database_reports/Constraint_violations

30 16 <https://www.mediawiki.org/wiki/Extension:EntitySchema>

31 17 <https://www.w3.org/2012/12/rdf-val/report>

332 where the community can iteratively develop a schema and refine the data to conform to that
333 schema. For those familiar with XML, ShEx is analogous to XML Schema or RelaxNG. *SHACL*
334 (Shapes Constraint Language), another language used to constraint RDF data models, uses a flat
335 list of constraints, analogous to XML's Schematron. *SHACL* was adapted from SPIN (SPARQL
336 Inference Notation) by the W3C Data Shapes working group in 2014 and became a W3C
337 recommendation in 2017 (Knublauch & Kontokostas, 2017). However, ShEx was chosen to
338 represent EntitySchemas in Wikidata, as it has a compact syntax that makes it more human-
339 friendly, supports recursion, and is designed to support distributed networks of reusable schemas
340 (Labra Gayo, Prud'hommeaux, Boneva, & Kontokostas, 2017). Besides the possibility to infer
341 ShEx expressions from the screening of a large set of concerned items, they can be easily and
342 intuitively written by humans.

343 In Wikidata, ShEx-based EntitySchemas are assigned an identifier (a number beginning with
344 an E) as well as labels, descriptions, and aliases in multiple languages, so that they can be easily
345 identified by users. Entity schemas are defined using the ShEx-compact syntax¹⁸, which is a
346 concise, human-readable syntax. A schema usually begins with some prefix declarations similar
347 to those in SPARQL. An optional start definition declares the shape which will be used by
348 default. In the example (Fig. 4), the shape <app> will be used, and its declaration contains a list
349 of properties, possible values, and cardinalities. By default, shapes are open, which means that
350 other properties apart from the ones declared are allowed. In this example, the values of property
351 `wdt:P31` are declared to be either a COVID-19 dashboard (`wd:Q90790055`), a search engine
352 (`wd:Q91136116`), or a dataset (`wd:Q91137337`). The `EXTRA` directive indicates that there can be
353 additional values for property `wdt:P31` that differ from the specified ones. The value for
354 property `wdt:P1476` is declared to be zero or more literals. The cardinality indicators come from
355 regular expressions, where '?' means zero or one, '*'; means zero or more, and '+' means one
356 or more. While the values for the other properties are declared to be anything (the dot indicates
357 no constraint) zero or more times, except for the properties `wdt:P577` and `wdt:P7103` that are
358 marked as optional using the question mark. Further documentation about ShEx can be found at
359 <http://shex.io/> and in Labra Gayo et al. (2017).

360

```
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
         ] ;
}
```

32 ¹⁸ ShEx schemas can also be defined in RDF-based representations like Turtle or JSON-LD.

```

wdt:P1476 LITERAL * ; #title
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
}

```

361 **Figure 4. Entity Schema example.** EntitySchema for COVID-19 dashboards, search engines and datasets [Source:
362 <https://www.wikidata.org/wiki/EntitySchema:E205> . File available on Wikimedia Commons: <https://w.wiki/4rg5>, License:
363 CC0.].

364
365 Due to the ease of using ShEx to define EntitySchemas, it has been used successfully to validate
366 Danish lexemes in Wikidata (Nielsen, Thornton, & Labra-Gayo, 2019) and biomedical Wikidata
367 statements (Thornton, et al., 2019). During the COVID-19 pandemic, Wikidata’s data model of
368 every COVID-19-related class as well as of all major biomedical classes has been converted to
369 an EntitySchema, so that it can be used to validate the representation of COVID-19 Wikidata
370 statements (Waagmeester, et al., 2021). These EntitySchemas were successfully used to enhance
371 the development and the robustness of the semantic structure of the data model underlying the
372 COVID-19 knowledge graph in Wikidata and are accordingly made available at a subpage of
373 Wikidata’s [WikiProject COVID-19,](https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Data_models) accessible via
374 [https://www.wikidata.org/wiki/Wikidata:WikiProject COVID-19/Data models](https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Data_models). Significant efforts are
375 currently underway to further simplify the definition of EntitySchemas by making them more
376 intuitive and concise, enabling an increase of the usage of ShEx to validate semantic knowledge
377 in Wikidata (Samuel, 2021).

378 Beyond these interesting methods, validation constraints can be inferred and used to verify
379 semantic statements in a knowledge graph through the use of the full screening of RDF dumps or
380 the use of SPARQL queries. RDF dumps are particularly used for screening Wikidata items in a
381 class to identify common features of the assessed entities based on a set of formal rules (Marx &
382 Krötzsch, 2017; Hanika et al. 2019). These features involve common characteristics of the data
383 model of the concerned class with patterns of used Wikidata properties such as symmetry and are
384 later used to verify the completeness of the class and validate the statements related to the
385 evaluated class. The analysis of RDF dumps for Wikidata can be coupled to the federated
386 screening of the RDF dumps of other knowledge graphs such as DBpedia through the alignment
387 of the types of relational and non-relational statements to benefit from the positive and negative
388 rules already defined and verified from the other knowledge resources to enrich the validation

389 tools of Wikidata (Ahmadi & Papotti, 2021). Nowadays, efforts are provided to extend inference-
390 based methods for the validation of Wikidata through the development of probabilistic
391 approaches to identify when a statement is unlikely to be defined for an item allowing to enhance
392 the evaluation of the completeness of Wikidata as an open knowledge graph (Arnaout, et al.,
393 2021). As SPARQL has been designed to extract a searched pattern from a semantic graph
394 (Pérez, Arenas, & Gutierrez, 2009), it has been used to query the competency questions¹⁹, and to
395 evaluate ontologies and knowledge graphs in a context-sensitive way (Vasanthapriyan, Tian, &
396 Xiang, 2017; Bansal & Chawla, 2016; Martin, 2018). Indeed, a sister project presents how
397 SPARQL can be used to generate data visualizations²⁰ (Nielsen, Mietchen & Willighagen 2017;
398 Shorland, Mietchen & Willighagen, 2020). Validating RDF data portals using SPARQL queries
399 has been regularly proposed as an approach that gives great flexibility and expressiveness (Labra
400 Gayo & Alvarez Rodríguez, 2013). However, academic literature is still far from revealing a
401 consensus on methods and approaches to evaluate ontologies using this query language
402 (Walisadeera, Ginige, & Wikramanayake, 2016), and other approaches have been proposed for
403 validation (Thornton, et al., 2019; Labra-Gayo, et al., 2019). Currently, there is mostly an effort
404 to normalize how to define SPARQL queries, particularly for knowledge graph validation
405 purposes, to save runtime and ameliorate the completeness of the output of a query using a set of
406 heuristics and axioms (Salas & Hogan, 2022).

407 In Wikidata, the Wikidata Query Service (<https://query.wikidata.org>) allows querying the
408 knowledge graph using SPARQL (Malyshev, et al., 2018; Turki, et al., 2019). The required
409 Wikidata prefixes are already supported in the backend of the service and do not need to be
410 defined (Malyshev, et al., 2018). What the user needs to do is to formulate their SPARQL query
411 (Black in Fig. 5) and click on the Run button (Blue in Fig. 5). After a compilation period, the
412 results will appear (Green in Fig. 5) and can be downloaded in different formats (Brown in Fig.
413 5), including JSON, TSV, CSV, HTML, and SVG. Different modes for the visualization of the
414 query results can be chosen (Purple in Fig. 5), particularly table, charts (line, scatter, area,
415 bubble), image grid, map, tree, timeline, and graph. The query service also allows users to use a
416 query helper (Red in Fig. 5) that can generate basic SPARQL queries, and to get inspired by
417 sample queries (Yellow in Fig. 5), especially when they lack experience. It also allows users to
418 generate a short link for the query (Pink in Fig. 5) and code snippets to embed the query results
419 in web pages and computer programs (Brown in Fig. 5) (Malyshev, et al., 2018).

33 ¹⁹ Competency questions: A set of requirements ensuring consistency of a knowledge graph, constraints
34 determining what knowledge to be involved in a knowledge graph (Wiśniewski, Potoniec, Ławrynowicz, & Keet,
35 2019).

36 ²⁰ For SPARQL-based visualizations of COVID-19 information in Wikidata, see <https://speed.ieee.tn/>,
37 <https://egonw.github.io/SARS-CoV-2-Queries/>,
38 https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Queries, and
39 <https://scholia.toolforge.org/topic/Q84263196>.

Wikidata Query Service

Examples Help More tools

English

Query Helper

+ Filter main subject COVID-19 pandemic

+ Show

Limit 100

```

1 SELECT ?COVID_19_pandemic ?COVID_19_pandemicLabel WHERE {
2   SERVICE wikibase:label { bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en". }
3   ?COVID_19_pandemic wdt:P921 wd:Q81068910.
4 }
5 LIMIT 100
6

```

100 results in 819 ms

Code Download Link

COVID_19_pandemic	COVID_19_pandemicLabel
Q83388131	The continuing 2019-nCoV epidemic threat of novel coronaviruses to global health - The latest 2019 novel coronavirus outbreak in Wuhan, China
Q83460376	Homologous recombination within the spike glycoprotein of the newly identified coronavirus may boost cross-species transmission from snake to human

421
 422 **Figure 5. Web interface of the Wikidata Query Service.** It involves a query field (Black), a query builder (Red), a short link
 423 button (Pink), a Run button (Blue), a visualization mode button (Purple), a download button (Brown), an embedding code
 424 generation button (Grey), a results field (green), and a sample query button (Yellow). [Source: <https://w.wiki/aeH>, Derived from:
 425 <https://query.wikidata.org>, License: CC0].

426

427

428 **Constraint-driven heuristics-based validation of epidemiological data**

429 The characterization of epidemiological data is possible using a variety of statistical measures
430 that show the acuteness, the dynamics, and the prognosis of a given disease outbreak. These
431 measures include the simple cumulative count of cases (P1603 [199569 statements, Orange in
432 Fig. 6], noted *c*, as defined before), deaths (P1120 [243250 statements²¹, Black in Fig. 6], noted
433 *d*), recoveries (P8010 [36119 statements, Green in Fig. 6], noted *r*), clinical tests (P8011 [21249
434 statements, Blue in Fig. 6], noted *t*), and hospitalized cases (P8049 [5755 statements, Grey in
435 Fig. 6], noted *h*) as well as several measurements done by the synthesis of the values of simple
436 epidemiological counts such as case fatality rate (P3457 [51504 statements, Red in Fig. 6], noted
437 *m*), basic reproduction number (P3492, noted R_0), minimal incubation period in humans (P3488,
438 noted *mn*), and maximal incubation period in humans (P3487, noted *mx*) (Rothman, Greenland,
439 & Lash, 2008). For all these statistical data, every information should be coupled by a *point in*
440 *time* (P585, noted *Z*) qualifier defining the date of the stated measurement and by a
441 *Determination method* (P459, noted *Q*) qualifier identifying the measurement method of the
442 given information as these variables are subject to change over days or according to used
443 methods of computation.

40 21 As of August 8, 2020. For updated statistics, see <https://w.wiki/Z5m>.

2020 COVID-19 pandemic in Tunisia (Q87343682)

viral outbreak in Tunisia

 edit

2020 coronavirus outbreak in Tunisia

Statements

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

444

445

446

Figure 6. Sample statistical data available through Wikidata. The item about the COVID-19 pandemic in Tunisia is shown. [Adapted from: <https://www.wikidata.org/wiki/Q87343682>, Source: <https://w.wiki/uUr>, License: CC0].

447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463

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 2. Such comparisons are enabled using simple statistical conditions that are commonly used in epidemiology (Zu, et al., 2020). 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 (Zu, et al., 2020). 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 COVID-19 deaths are the consequence of severe infections by SARS-CoV-2 that can only be managed in hospitals (Rothman, Greenland, & Lash, 2008) and as a patient needs to undergo COVID-19 testing to be confirmed as a case of the disease (Zu, et al., 2020). V10 is built upon the assumption that c , d , r , h , and t values can be geographically aggregated (Rothman, Greenland, & Lash, 2008).

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> March 25, 20
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> COVID-19 Dashboard
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> no value <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 <i>c</i> statements having a value superior or equal to the number of clinical tests for a date <i>Z</i>	(<i>COVID-19 pandemic in X</i> <number of clinical tests> 4 <point in time> <i>March 25, 2020</i>) AND (<i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> <i>March 25, 2020</i>)
V8	Identify <i>c</i> statements having a value inferior to the number of deaths for a date <i>Z</i>	(<i>COVID-19 pandemic in X</i> <number of deaths> 10 <point in time> <i>March 25, 2020</i>) AND (<i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> <i>March 25, 2020</i>)
V9	Identify <i>c</i> statements having a value inferior to the number of recoveries for a date <i>Z</i>	(<i>COVID-19 pandemic in X</i> <number of recoveries> 10 <point in time> <i>March 25, 2020</i>) AND (<i>COVID-19 pandemic in X</i> <number of cases> 5 <point in time> <i>March 25, 2020</i>)
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> <i>March 25, 2020</i>) AND (<i>COVID-19 pandemic in Y</i> <number of cases> 5 <point in time> <i>March 25, 2020</i>) WHERE <i>X</i> is a district of <i>Y</i>

464 **Table 2. Tasks for the heuristics-based evaluation of epidemiological data using the Wikidata SPARQL endpoint.** Each
465 validation task is given with its identifier, a brief description of the heuristic validation criteria and an example where the data
466 does not fit them. See the section "Constraint-driven heuristics-based validation of epidemiological data" for definitions of the
467 epidemiological variables.

468
469 This task set has easily been applied using ten simple SPARQL queries that can be found in
470 Appendix A where <PropertyID> is the Wikidata property to be analyzed and has returned 5496
471 deficiencies in the COVID-19 epidemiological information (as of August 8, 2020) as shown in
472 Table 3. Among these mistaken statements, 2856 were *number of cases* statements, 2467 were
473 *number of deaths* statements, 189 were *number of recoveries* statements, 9 were *number of*
474 *clinical tests* statements, and 10 were *number of hospitalized cases* statements. This distribution
475 of the deficiencies among epidemiological properties is explained by the dominance of *number*
476 *of cases* and *number of deaths* statements on the COVID-19 epidemiological information. Most
477 of these mistakes are linked to a violation of the cumulative pattern of major variables. These
478 deficiencies can be removed using tools for the automatic enrichment of Wikidata like
479 QuickStatements (cf. Turki, et al., 2019) or adjusted one by one by active members of
480 WikiProject COVID-19.

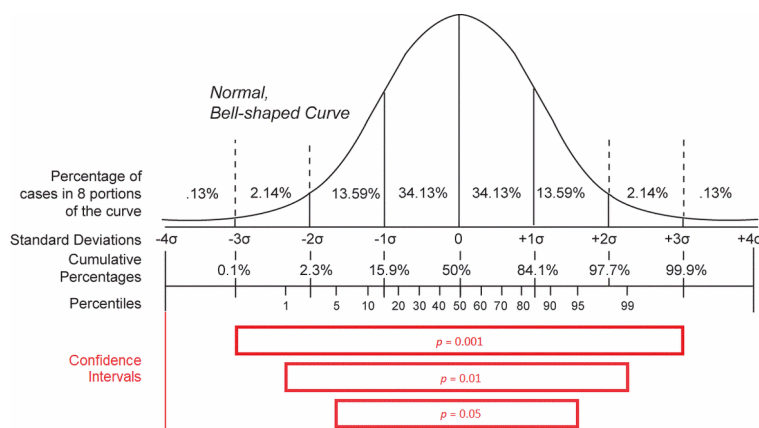
	<i>c</i>	<i>d</i>	<i>r</i>	<i>t</i>	<i>h</i>	Overall
V1	18	9	10	2	1	40
V2	2	91	6	0	0	99
V3	660	92	6	5		763
V4	2081	2247	149	1		4478
V5	0	0	0	0	0	0

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

481 **Table 3.** Matrix overview of data quality issues identified per validation task and epidemiological Wikidata property.
482 Rows represent validation tasks as defined in Table 2, columns the corresponding epidemiological Wikidata properties, and the
483 value in a given cell represents the number of deficient statements identified by the row's specific task for the column's
484 epidemiological Wikidata property on a given date (August 8, 2020).
485

486 Concerning the variables issued from the integration of basic epidemiological counts (m , R_0 , mn ,
487 and mx statements), they give a summary overview of the statistical behavior of the studied
488 infectious pandemic and that is why they can be useful to identify if the stated evolution of the
489 morbidity and mortality caused by the outbreak is reasonable (Delamater, et al., 2019). However,
490 the validation of these variables is more complicated due to the complexity of their definition
491 (Delamater, et al., 2019; Backer, Klinkenberg, & Wallinga, 2020; Li, et al., 2020). The basic
492 reproduction number (R_0) is meant to be a constant that characterizes the dissemination power of
493 an infectious disease. It is defined as the expected number of people (within a community with
494 no prior exposure to the disease) that can contract a disease via the same infected individual.
495 This variable should exceed the threshold of 1 to define a contagious disease (Delamater, et al.,
496 2019). Although R_0 can give an idea about the general behavior of an outbreak of a given
497 disease, any calculated value depends on the model used for its computation (e.g., *SIR Model*) as
498 well as the underlying data and is consequently a bit imprecise and variable from one study to
499 another (Delamater, et al., 2019). That is why it is not reliable to use this variable to evaluate the
500 accuracy of simple epidemiological counts for a given pandemic. The only heuristic that can be
501 applied to this variable is to verify if its value exceeds 1 for diseases causing large outbreaks.
502 The incubation period of a disease gives an overview of the silent time required by an infectious
503 agent to become active in the host organism and cause notable symptoms (Backer, Klinkenberg,
504 & Wallinga, 2020; Li, et al., 2020). This variable is very important, as it reveals how many days
505 an inactive case can spread the disease in the host's environment before the host is being
506 symptomatically identified. As a result, it can give an idea about the contagiousness of the
507 infectious disease and its basic reproduction number (R_0). However, the determination of the
508 incubation period - especially for a novel pathogen - is challenging, as a patient often cannot
509 identify with precision the day when they had been exposed to the disease, at least if they did not
510 travel to an endemic region or had not been in contact with a person they knew to be infected.
511 This factor was behind the measurement of falsely small incubation periods for COVID-19 at the

512 beginning of the COVID-19 epidemic in China (Backer, Klinkenberg, & Wallinga, 2020).
 513 Furthermore, the use of minimal (*mn*) and maximal (*mx*) incubation periods in Wikidata to
 514 epidemiologically describe a disease instead of the median incubation period is a source of a lack
 515 of accuracy of the extracted values (Backer, Klinkenberg, & Wallinga, 2020; Li, et al., 2020).
 516 Minimal and maximal incubation periods for a given disease are obtained in the function of the
 517 mean (\bar{X}) and standard deviation (σ) of the measures of the confidence interval of observed
 518 incubation periods in patients. Effectively, *mn* is equal to $\bar{X} - \frac{z * \sigma}{\sqrt{n}}$ and *mx* is equal to $\bar{X} + \frac{z * \sigma}{\sqrt{n}}$
 519 where *n* is the number of analyzed observations and *z* is a characteristic of the hypothetical
 520 statistical distribution and of the statistical confidence level adopted for the estimation (Altman,
 521 et al., 2013). As a consequence, *mn* and *mx* variables are modified according to the number of
 522 observations (*n*) with a smaller difference between the two variables for higher values of *n*. The
 523 two measures also vary according to the used statistical distribution and that is why different
 524 values of *mn* and *mx* were reported for COVID-19 when applying different distributions
 525 (Weibull, gamma, and log-normal distribution) using a confidence level of 0.95 on the same set
 526 of observed cases (Backer, Klinkenberg, & Wallinga, 2020). Similarly, the two variables can
 527 change according to the adopted confidence level (*p* - 1) when using the same statistical
 528 distribution where a higher confidence level is correlated with a higher difference between the
 529 calculated *mn* and *mx* values, as shown in Fig. 7 (Ward & Murray-Ward, 1999; Altman, et al.,
 530 2013). Given these reasons and despite the significant importance of the two measures, these two
 531 statistical variables cannot be used to evaluate statistical epidemiological counts for COVID-19
 532 due to their lack of precision and difficulty of determination.
 533



534
 535 **Figure 7. Distribution statistics.** Confidence intervals for different *p*-values (*p*) when using a normal distribution [Source:
 536 <https://w.wiki/aKT>, License: Public Domain] (after Ward & Murray-Ward, 1999).
 537

538 As for the reported case fatality rate (*m*), it is simply the quotient of the cumulative number of
 539 deaths (*d*) and the cumulative number of cases (*c*) as stated in official reports. It is consequently
 540 easy to validate for a given disease by comparing its values with simple reported counts of cases
 541 and deaths (Rothman, Greenland, & Lash, 2008). Here, two simple heuristics can be applied

542 using SPARQL queries as shown in Appendix B. As the number of deaths is less than or equal to
543 the number of cases of a given disease, m values should be set between 0 and 1. That is why
544 Task M1 is defined to extract m statements where $m > 1$ or $m < 0$. Also, as $m = d / c$ for a date Z ,
545 m values that are not close to the corresponding quotients of deaths by disease cases should be
546 identified as deficient and m values should be stated for a given date Z if mortality and morbidity
547 counts exist. Thus, Task M2 is created to extract m values where the absolute value of $(m - d/c)$
548 is superior to 0.001, and Task M3 is developed to identify (item, date) pairs where m statements
549 are missing and c and d statements are available in Wikidata. Absolute values for Task M2 are
550 obtained using SPARQL's ABS function, and deficient (item, date) pairs are eliminated in Task
551 M3 where $m > 1$ and $c < d$.

552 As a result of these three tasks, we interestingly identified 143 deficient m statements and
553 7116 missing m statements. 133 of the mistaken statements are identified thanks to Task M2 and
554 concern 25 Wikidata items and 31 distinct dates, and only 10 deficient statements related to 3
555 Wikidata items and 8 distinct dates are found using Task M1. These statements should be
556 verified against reference datasets to verify their values and to determine the reason behind their
557 deficiency. Such a reason can be the integration of the wrong case and death counts in Wikidata,
558 or a bug or inaccuracy within the source code of the bot making or updating such statements. The
559 verification process can be automatically done using an algorithm that compares Wikidata values
560 (c , d , and m statements) with their corresponding ones in other databases (using file or API
561 reading libraries) and subsequently adjusts statements using the Wikidata API directly or via
562 tools like QuickStatements (Turki et al., 2019). As for the missing m statements returned by M3,
563 they are linked to 395 disease outbreak items and to 205 distinct dates and concern 70%
564 (7116/10168) of the (case count, death count) pairs available in Wikidata. The outcome of M3
565 proves the efficiency of comparative constraints to enrich and assess the completeness of
566 epidemiological data available in a knowledge graph, particularly Wikidata, based on existing
567 information. Consequently, derivatives of Task M3 can build to infer d values based on c and m
568 statements or to find c values based on d and m statements. The missing statements found by
569 such tasks can be integrated in Wikidata using a bot based on Wikidata API and Wikidata Query
570 Service to ameliorate the completeness and integrity of available mortality data for epidemics,
571 mainly the COVID-19 pandemic (Turki, et al., 2019).

572

573 Discussion

574 The results presented here demonstrate the value of our statistical constraints-based validation
575 approach for knowledge graphs like Wikidata across a range of features (Tables 2 and 3). These
576 tasks successfully address most of the competency questions, particularly conceptual orientation
577 (*clarity*), coherence (*consistency*), strength (*precision*), and full coverage (*completeness*).
578 Combined with previous findings in the context of bioinformatics (Bolleman, et al., 2020; Marx
579 & Krötzsch, 2017; Darari, et al., 2020), this proves that the efficiency of rule-based approaches to
580 evaluate semantic information from scratch displays a similar accuracy as other available

581 ontology evaluation algorithms (Amith, et al., 2019; Zhang & Bodenreider 2010). The efficiency
582 of these constraint-based assessment methods can be further enhanced by using machine learning
583 techniques to perform imputations and adjustments on deficient data (Bischof, et al., 2020). The
584 scope of rule-based methods can be similarly expanded to cover other competency questions
585 such as non-redundancy (*conciseness*) through the proposal of other logical constraints to tackle
586 them, such as a condition to find taxonomic relations to trim in a knowledge graph (examples can
587 be found at https://www.wikidata.org/wiki/Wikidata:Database_evaluation). The main limitation of
588 applying the logical constraints using SPARQL in the context of Wikidata is that the runtime of a
589 query that infers or verifies a complex condition or that analyzes a huge amount of class items or
590 property use cases can exceed the timeout limit of the used endpoint (Malyshev, et al., 2018;
591 Chah & Andritsos, 2021). Here, the inference of logical constraints and the identification of
592 inconsistent semantic information through the analysis of full dumps of Wikidata can be more
593 efficient, although this comes with advanced storage and processing requirements (Chah &
594 Andritsos, 2021).

595 These evaluation assignments covered by our approach can be done by other rule-based
596 (*structure-based* and *semantic-based*) ontology evaluation methods. Structure-based methods
597 verify whether a knowledge graph is defined according to a set of formatting constraints, and
598 semantic-based methods check whether concepts and statements of a knowledge graph meet
599 logical conditions (Amith, et al., 2018). Some of these methods are software tools, particularly
600 Protégé extensions such as OWLET (Lampoltshammer & Heistracher, 2014) and OntoCheck
601 (Schober, et al., 2012). OWLET infers the JSON schema logics of a given knowledge graph,
602 converts them into OWL-DL axioms, and uses the semantic rules to validate the assessed
603 ontological data (Lampoltshammer & Heistracher, 2014). OntoCheck screens an ontology to
604 identify structural conventions and constraints for the definition of the analyzed relational
605 information and consequently to homogenize the data structure and quality of the ontology by
606 eliminating typos and pattern violations (Schober, et al., 2012). Here, the advantage of applying
607 constraints using SPARQL is that its runtime is faster, as it does not require the download of the
608 full dumps of the evaluated knowledge graph (Malyshev, et al., 2018). The benefit of our method
609 and other structure-based and semantic-based web-based tools for knowledge graph validation
610 like OntoKeeper (Amith, et al., 2019) and adviseEditor (Geller, et al., 2013), when compared to
611 software tools, is that the maximal size of the knowledge graphs that can be assessed by web
612 services is larger than the one that can be evaluated by software tools because the latter depends
613 on the requirements and capacities of the host computer (Lampoltshammer & Heistracher, 2014;
614 Schober, et al., 2012). These drawbacks of other structure-based tools can indeed be solved
615 through the simplification of the knowledge graph by reducing redundancies using techniques
616 like ontology trimming (Jantzen, et al., 2011) or through the construction of an abstraction
617 network to decrease the complexity of the analyzed knowledge graph (Amith, et al., 2018;
618 Halper, et al., 2015). However, knowledge graph simplification processes are time-consuming,

619 and resulting time gain can consequently be insignificant (Jantzen, et al., 2011; Amith, et al.,
620 2018; Halper, et al., 2015).

621 Such tasks can be also solved using data-driven ontology evaluation methods. These
622 techniques process texts in natural languages to validate the concepts and statements of a
623 knowledge graph and currently include intrinsic (*lexical-based*) and extrinsic (*cross-validation*,
624 *big data-based*, and *corpus-based*) methods (Amith, et al., 2018). Lexical-based methods use
625 rules implemented in SQL or SPARQL to retrieve items and glosses corresponding to a concept
626 and their semantic relations (mostly *subclass of* statements) (Rector & Iannone, 2012; Luo,
627 Mejino Jr, & Zhang, 2013). These items are then compared against a second set of rules to
628 identify inconsistencies in their labels, descriptions, or semantic relations (Amith, et al., 2018).
629 The output can then be analyzed using natural language processing techniques such as hamming
630 distance measures (Luo, Mejino Jr, & Zhang, 2013), semantic annotation tools (Rector &
631 Iannone, 2012), and semantic similarity measures (Amith, et al., 2018) to comparatively identify
632 deficiencies in the semantic representation, labelling, and symmetry of the assessed knowledge
633 graph. Conversely, extrinsic data-based methods extract the usage and linguistic patterns from
634 raw text corpuses such as bibliographic databases and clinical records (*Corpus-based methods*)
635 or from gold standard semantic resources like large ontologies and knowledge graphs (*Cross-*
636 *validation methods*) or social media posts and interactions, Internet of Things data or web service
637 statistics (*Big data-based methods*) (Amith, et al., 2018; Sebei, Hadj Taieb, & Ben Aouicha,
638 2018; Rector, Brandt, & Schneider, 2011; Gangemi, et al., 2005) using structure-based and
639 semantic-based ontology evaluation methods as explained above (Rector, Brandt, & Schneider,
640 2011) as well as a range of techniques including machine learning (Bean, et al., 2017; Zhang, et
641 al., 2018), topic modeling using Latent Dirichlet Analysis (Abd-Alrazaq, et al., 2020), word
642 embeddings (Zhang, et al., 2019), statistical correlations (Vanderkam, et al., 2013) and semantic
643 annotation methods (Li, et al., 2016). The returned features of the analyzed resources are
644 compared to the ones of the analyzed knowledge graph to assess the accuracy and completeness
645 of the definition and use of concepts and properties (Amith, et al., 2018).

646 When compared to our proposed approach, lexical-based methods have the advantage to
647 identify and adjust characteristics of a knowledge graph item based on its natural language
648 information of a knowledge graph item, particularly terms and glosses (Rector & Iannone, 2012;
649 Luo, Mejino Jr, & Zhang, 2013). The drawback of using semantic similarity, word embeddings,
650 and topic modeling techniques in such approaches is that these techniques are sensitive to the
651 used parameters, to input characteristics, and to the chosen models of computation and can
652 consequently give different results according to the context of determination (Lastra-Díaz, et al.,
653 2019; Hadj Taieb, Zesch, & Ben Aouicha, 2020). The current role of constraints in the extraction
654 of lexical information and respective semantic relations (Rector & Iannone, 2012; Luo, Mejino
655 Jr, & Zhang, 2013) proves that the scope of constraint-based validation should not only be
656 restricted to rule-based evaluation but also to lexical-based evaluation. Yet, the function of
657 logical conditions should be expanded to refine the list of pairs (lexical information, semantic

658 relation) to more accurately identify deficient and missing semantic relations and defective
659 lexical data and to support multilingual lexical-based methods. This would build on the many
660 SPARQL functions that analyze strings in knowledge graphs²² such as STRLEN (length of a
661 string), STRSTARTS (verification of a substring beginning a given string), STRENDS
662 (verification of a substring finishing a given string), and CONTAINS (verification of a substring
663 included in a given string) (DuCharme, 2013; Harris, Seaborne, & Prud'hommeaux, 2013).

664 As for the extrinsic data-driven methods, they are mainly based on large-scale resources that
665 are regularly curated and enriched. Raw-text corpora are mainly composed of scholarly
666 publications (Raad & Cruz, 2015) and blog posts (Park, et al., 2016). Information in scholarly
667 publications is ever-changing according to the dynamic advances in scholarly knowledge,
668 particularly medical data (Jalalifard, Norouzi, & Isfandyari-Moghaddam, 2013). This expansion
669 of scientific information in scholarly publications is highly recognized in the context of COVID-
670 19 where detailed information about COVID-19 disease and the SARS-CoV-2 virus is published
671 within less than six months (Kagan, Moran-Gilad, & Fire, 2020). Big data is the set of real-time
672 statistical and textual information that is generated by web services including search engines and
673 social media and by the Internet of Things objects including sensors (Sebei, Hadj Taieb, & Ben
674 Aouicha, 2018). This data is characterized by its value, variety, variability, velocity, veracity,
675 and volume (Sebei, Hadj Taieb, & Ben Aouicha, 2018) and can be consequently used to track the
676 changes of the community knowledge and consciousness over time (Abd-Alrazaq, et al., 2020;
677 Turki, et al., 2020). Large semantic resources are ontologies and knowledge graphs that are built
678 and curated by a community of specialists and that are regularly verified, updated, and enriched
679 using human efforts and computer programs (Lee, et al., 2013). These resources represent broad
680 and reliable information about a given specialty through machine learning techniques (Zhang, et
681 al., 2018) and the crowdsourcing of scientific efforts (Mortensen, et al., 2014) and can be
682 consequently compared to other semantic databases for validation purposes. Examples of these
683 resources are the COVID-19 Disease Map (Ostaszewski, et al., 2020) and SNOMED-CT²³ (Lee,
684 et al., 2013).

685 Large-scale knowledge graphs are dynamic corpora. Changes in the logical and semantic
686 conditions for the definition of knowledge in a particular domain need to be identified to adjust
687 the assessed knowledge graph accordingly. Rule-based and lexical-based approaches (especially
688 constraints-based methods) are therefore less simple to apply than extrinsic data-driven methods
689 (Amith, et al., 2018). Nonetheless, the growing and changing nature of gold-standard resources
690 require continuous human efforts and an advanced software architecture to maintain (e.g.,
691 structure-based and semantic-based methods), process (e.g., *word embeddings* and *latent*
692 *Dirichlet analysis*), and store (e.g., *Hadoop* and *MapReduce*) these reference resources
693 (Mortensen, et al., 2014; Le, et al., 2013; Sebei, Hadj Taieb, & Ben Aouicha, 2018). This

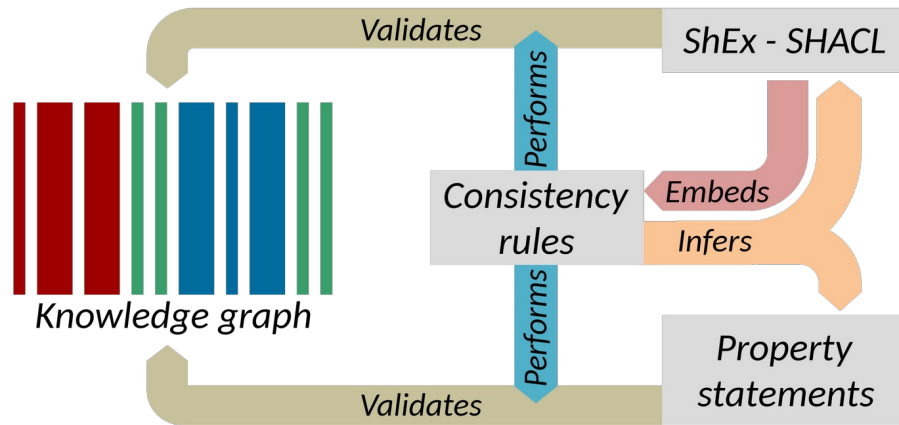
41 22 Detailed information about string functions in SPARQL can be found at [https://www.w3.org/TR/sparql11-
43 query/#func-strings](https://www.w3.org/TR/sparql11-
42 query/#func-strings).

43 23 Systematized Nomenclature Of Medicine - Clinical Terms

694 architecture has advanced hardware requirements and its results are subject to change according
695 to the used parameters (Sebei, Hadj Taieb, & Ben Aouicha, 2018).

696 These tasks are in line with the usage of Shape Expressions as well as property constraints
697 and relations for the validation of data quality and completeness of the semantic information of
698 class items in knowledge graphs as shown in the “Knowledge graph validation of Wikidata”
699 section. A ShEx ShapeMap is a pair of a triple pattern for selecting entities to validate and a
700 shape against which to validate them. This allows for the definition of the properties to be used
701 for the items of a given class (Prud'hommeaux, Labra Gayo, & Solbrig, 2014; Waagmeester, et
702 al., 2021) and property constraints and relations based on the meta-ontology (i.e., data skeleton)
703 of Wikidata. Expressions written in shape-based property usage validation languages for RDF
704 (e.g., *SHACL*) can be used to state conditions and formatting restrictions for the usage of
705 relational and non-relational properties (Erxleben, et al., 2014; Thornton, et al., 2019; Gangemi,
706 et al., 2005). SPARQL can be more efficient in inferring such information than the currently
707 existing techniques that screen all the items and statements of a knowledge graph one by one to
708 identify the conditions for the usage of properties (e.g., *SQID*) mainly because SPARQL is
709 meant to directly extract information according to a pattern without having to evaluate all the
710 conditions against all items of a knowledge graph (Marx & Krötzsch, 2017; Hanika, et al., 2019;
711 Pérez, Arenas, & Gutierrez, 2009).

712 The separate execution of value-based constraints is common in the quality control of XML
713 data. Typically, structural constraints are managed by RelaxNG or XML Schemas, while value-
714 based constraints are captured as Schematron. Much as Schematron rules are typically embedded
715 in RelaxNG, the consistency constraints presented above can be embedded in Shape Expressions
716 Semantic Actions or in SHACL-SPARQL as shown in Fig. 8 (Melo & Paulheim, 2020). These
717 supplement structural schema languages with mechanisms to capture value-based constraints and
718 in doing so, provide context for the enforcement of those constraints. The implementation of
719 value-based constraints shown in the “Constraint-driven heuristics-based validation of
720 epidemiological data” section can likewise be implemented in a shapes language (Labra-Gayo, et
721 al., 2019). Parsing the rules in the Table 2 would allow the mechanical generation or
722 augmentation of shapes, providing flexibility for how the rules are expressed while still
723 exploiting the power of shape languages for validation. More generally, ontology-based and
724 knowledge graph-based software tools have the potential to provide wide data and platform
725 interoperability, and thus their semantic interoperability is relevant for a range of downstream
726 applications such as IoT and WoT technologies (Gyrard, Datta, & Bonnet, 2018).



727
 728 **Figure 8. Key elements of data quality workflows on Wikidata.** Interactions between consistency rules, property statements,
 729 and RDF validation languages [Source: <https://w.wiki/ao5>, License: CC BY 4.0]
 730

731 **Conclusion**

732 In this paper, we investigate how to best assess COVID-19 epidemiological knowledge in
 733 collaborative ontologies and knowledge graphs based on the example of Wikidata using
 734 statistical constraints. Collaborative databases produced through the cumulative edits of
 735 thousands of users can generate huge amounts of structured information (Turki, et al., 2019) but
 736 as a result of their rather uncoordinated development, they often result in uneven coverage of
 737 crucial information and inconsistent expression of that information. The resulting gaps are a
 738 significant problem (conflicting values, reasoning deficiencies, and missing statements).
 739 Avoiding, identifying, and closing these gaps is therefore of top importance. We presented a
 740 standardized methodology for auditing key aspects of data quality and completeness for these
 741 resources²⁴.

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

749 This work has shown the state of the knowledge graph as a snapshot in time. Future work will
 750 extend this to investigate how the knowledge base evolves as more biomedical knowledge is
 751 integrated into it over time. This will require incorporating the edit history in the SPARQL
 752 endpoint APIs of knowledge graphs (Pellissier Tanon & Suchanek, 2019, Dos Reis, Pruski, Da
 753 Silveira, & Reynaud-Delaître, 2014) to dynamically visualize time-resolved SPARQL queries.

44 ²⁴ This method can be adapted to meet the needs of the user. For instance, the SPARQL queries can be slightly
 45 adjusted to assess other patterns in collaborative ontologies such as the usage of classes.

754 We will also couple the information inferred using this method²⁵ with Shape Expressions and the
755 explicit constraints of relation types to provide a more effective enrichment, refinement, and
756 adjustment of collaborative ontologies and knowledge graphs with statistical data. This will be
757 an excellent infrastructure to enable the support of non-relational information. We look forward
758 to extending our proposed approach to allow knowledge graphs to handle non-relational
759 statements about future epidemics and other disasters such as earthquakes.

760

761 **Author statements**

762 **Data availability:** All the SPARQL queries used in this research work are provided in the
763 appendices. The Internet Archive links of the URLs cited by this paper are made available at
764 [https://web.archive.org/save/https://www.wikidata.org/w/index.php?](https://web.archive.org/save/https://www.wikidata.org/w/index.php?title=User:Daniel_Mietchen/sandbox&oldid=1580603965)
765 [title=User:Daniel_Mietchen/sandbox&oldid=1580603965](https://web.archive.org/save/https://www.wikidata.org/w/index.php?title=User:Daniel_Mietchen/sandbox&oldid=1580603965).

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

772

773 **Acknowledgements**

774 The work done by Houcemeddine Turki, Mohamed Ali Hadj Taieb, and Mohamed Ben Aouicha
775 was supported by the Ministry of Higher Education and Scientific Research in Tunisia
776 (MoHESR) in the framework of Federated Research Project PRFCOV19-D1-P1, by Wikimedia
777 Foundation through a rapid grant, and by WikiCred Grants Initiative of Craig Newmark
778 Philanthropies, Facebook, and Microsoft. The work done by Jose Emilio Labra Gayo was
779 partially funded by the Spanish Ministry of Economy and Competitiveness (Society challenges:
780 TIN2017-88877-R). The work done by Daniel Mietchen was supported in part by the Alfred P.
781 Sloan Foundation under grant numbers G-2019-11458 and G-2021-17106. The work done by
782 Dariusz Jemielniak was funded by the Polish National Science Center grant no
783 2019/35/B/HS6/01056. We thank the Wikidata community, Olivier Corby (Université Côte
784 d'Azur, France), Odile Papini (Aix-Marseille Université, France), Egon Willighagen (Maastricht
785 University, Netherlands), and Mahir Morshed (University of Illinois at Urbana-Champaign,
786 United States of America) for useful comments and discussions about the topic of this research
787 paper. This research paper is published on behalf of the WikiProject COVID-19 members: Jan
788 Ainali, Susanna Ånäs, Erica Azzellini, Mus'ab Banat, Mohamed Ben Aouicha, Alessandra
789 Boccone, Jane Darnell, Diptanshu Das, Lena Denis, Rich Farmbrough, Daniel Fernández-
790 Álvarez, Konrad Foerstner, Jose Emilio Labra Gayo, Mauricio V. Genta, Mohamed Ali Hadj
791 Taieb, James Hare, Alejandro González Hevia, David Hicks, Toby Hudson, Netha Hussain,
792 Jinoy Tom Jacob, Dariusz Jemielniak, Krupal Kasyap, Will Kent, Samuel Klein, Jasper J.
793 Koehorst, Martina Kutmon, Antoine Logean, Tiago Lubiana, Andy Mabbett, Kimberli
794 Mäkäräinen, Tania Maio, Bodhisattwa Mandal, Nandhini Meenakshi, Daniel Mietchen, Nandana

46 ²⁵ This information can be represented in the form of RDF triples where the subject is the studied relation type and
47 integrated into Wikidata.

795 Mihindikulasooriya, Mahir Morshed, Peter Murray-Rust, Minh Nguyễn, Finn Årup Nielsen,
796 Mike Nolan, Shay Nowick, Julian Leonardo Paez, João Alexandre Peschanski, Alexander Pico,
797 Lane Rasberry, Mairelys Lemus-Rojas, Diego Saez-Trumper, Magnus Säljö, John Samuel, Peter
798 J. Schaap, Jodi Schneider, Thomas Shafee, Nick Sheppard, Adam Shorland, Ranjith Siji, Michal
799 Josef Špaček, Ralf Stephan, Andrew I. Su, Hilary Thorsen, Houcemeddine Turki, Lisa M.
800 Verhagen, Denny Vrandečić, Andra Waagmeester, and Egon Willighagen.
801

802 **References**

- 803 Abd-Alrazaq, A., Alhuwail, D., Househ, M., Hamdi, M., & Shah, Z. (2020). Top concerns of tweeters
804 during the COVID-19 pandemic: infoveillance study. *Journal of medical Internet research*, 22(4),
805 e19016. doi:10.2196/19016.
- 806 Ahmadi, N., & Papotti, P. (2021, April). Wikidata Logical Rules and Where to Find Them. In *Companion*
807 *Proceedings of the Web Conference 2021* (pp. 580-581). doi:10.1145/3442442.3452343.
- 808 Altman, D., Machin, D., Bryant, T., & Gardner, M. (Eds.). (2013). *Statistics with confidence: confidence*
809 *intervals and statistical guidelines*. John Wiley & Sons. ISBN:978-0-727-91375-3.
- 810 Amith, M., He, Z., Bian, J., Lossio-Ventura, J. A., & Tao, C. (2018). Assessing the practice of biomedical
811 ontology evaluation: Gaps and opportunities. *Journal of Biomedical Informatics*, 80, 1-13.
812 doi:10.1016/j.jbi.2018.02.010.
- 813 Amith, M., Manion, F., Liang, C., Harris, M., Wang, D., He, Y., & Tao, C. (2019). Architecture and
814 usability of OntoKeeper, an ontology evaluation tool. *BMC medical informatics and decision*
815 *making*, 19(4), 152. doi:10.1186/s12911-019-0859-z.
- 816 Arnaout, H., Razniewski, S., Weikum, G., & Pan, J. Z. (2021, April). Negative knowledge for open-world
817 Wikidata. In *Companion Proceedings of the Web Conference 2021* (pp. 544-551).
818 doi:10.1145/3442442.3452339.
- 819 Backer, J. A., Klinkenberg, D., & Wallinga, J. (2020). Incubation period of 2019 novel coronavirus
820 (2019-nCoV) infections among travellers from Wuhan, China, 20–28 January 2020.
821 *Eurosurveillance*, 25(5), 2000062. doi:10.2807/1560-7917.ES.2020.25.5.2000062.
- 822 Bansal, R., & Chawla, S. (2016). Design and development of semantic web-based system for computer
823 science domain-specific information retrieval. *Perspectives in Science*, 8, 330–333.
824 doi:10.1016/j.pisc.2016.04.067.
- 825 Bean, D. M., Wu, H., Iqbal, E., Dzahini, O., Ibrahim, Z. M., Broadbent, M., et al. (2017). Knowledge
826 graph prediction of unknown adverse drug reactions and validation in electronic health records.
827 *Scientific reports*, 7(1), 1-11. doi:10.1038/s41598-017-16674-x.
- 828 Ben Aouicha, M., & Hadj Taieb, M. A. (2016). Computing semantic similarity between biomedical
829 concepts using new information content approach. *Journal of biomedical informatics*, 59, 258-
830 275. doi:10.1016/j.jbi.2015.12.007.
- 831 Bischof, S., Harth, A., Kämpgen, B., Polleres, A., & Schneider, P. (2018). Enriching integrated statistical
832 open city data by combining equational knowledge and missing value imputation. *Journal of Web*
833 *Semantics*, 48, 22-47. doi:10.1016/j.websem.2017.09.003.
- 834 Bolleman, J., de Castro, E., Baratin, D., Gehant, S., Cuche, B. A., Auchincloss, A. H., et al. (2020).
835 HAMAP as SPARQL rules—A portable annotation pipeline for genomes and proteomes.
836 *GigaScience*, 9(2), g1aa003. doi:10.1093/gigascience/g1aa003.

837 Brank, J., Grobelnik, M., & Mladenic, D. (2005). A survey of ontology evaluation techniques.
838 *Proceedings of the Conference on Data Mining and Data Warehouses (SiKDD 2005)* (pp. 166–
839 170). Ljubljana, Slovenia: Citeseer. [http://citeseerx.ist.psu.edu/viewdoc/summary?](http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.101.4788)
840 [doi=10.1.1.101.4788](http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.101.4788).

841 Burgstaller-Muehlbacher, S., Waagmeester, A., Mitra, E., Turner, J., Putman, T., Leong, J., et al.
842 (2016). Wikidata as a semantic framework for the Gene Wiki initiative. *Database, 2016*, baw015.
843 [doi:10.1093/database/baw015](https://doi.org/10.1093/database/baw015).

844 Chah, N., & Andritsos, P. (2021). WikiMetaData Studio: Dashboards From Data Profiling the Languages,
845 Properties, and Items of Wikidata. In *Proceedings of the 2nd Wikidata Workshop*
846 (*Wikidata@ISWC 2021*) (pp. 13:1-13:8). <http://ceur-ws.org/Vol-2982/paper-13.pdf>.

847 Chen, Q., Lee, K., Yan, S., Kim, S., Wei, C. H., & Lu, Z. (2020). BioConceptVec: Creating and
848 evaluating literature-based biomedical concept embeddings on a large scale. *PLoS computational*
849 *biology, 16*(4), e1007617. [doi:10.1371/journal.pcbi.1007617](https://doi.org/10.1371/journal.pcbi.1007617).

850 Cuan-Baltazar, J. Y., Muñoz-Perez, M. J., Robledo-Vega, C., Pérez-Zepeda, M. F., & Soto-Vega, E.
851 (2020). Misinformation of COVID-19 on the internet: infodemiology study. *JMIR public health*
852 *and surveillance, 6*(2), e18444. [doi:10.2196/18444](https://doi.org/10.2196/18444).

853 Darari, F., Nutt, W., Razniewski, S., & Rudolph, S. (2020). Completeness and soundness guarantees for
854 conjunctive SPARQL queries over RDF data sources with completeness statements. *Semantic*
855 *Web, 11*(3), 441-482. [doi:10.3233/SW-190344](https://doi.org/10.3233/SW-190344).

856 Degbelo, A. (2017). A Snapshot of Ontology Evaluation Criteria and Strategies. *Proceedings of the 13th*
857 *International Conference on Semantic Systems* (pp. 1–8). New York: ACM.
858 [doi:10.1145/3132218.3132219](https://doi.org/10.1145/3132218.3132219).

859 Delamater, P. L., Street, E. J., Leslie, T. F., Yang, Y. T., & Jacobsen, K. H. (2019). Complexity of the
860 basic reproduction number (R0). *Emerging infectious diseases, 25*(1), 1-4.
861 [doi:10.3201/eid2501.171901](https://doi.org/10.3201/eid2501.171901).

862 Dong, E., Du, H., & Gardner, L. (2020). An interactive web-based dashboard to track COVID-19 in real
863 time. *The Lancet infectious diseases, 20*(5), 533-534. [doi: 10.1016/S1473-3099\(20\)30120-1](https://doi.org/10.1016/S1473-3099(20)30120-1).

864 Dos Reis, J. C., Pruski, C., Da Silveira, M., & Reynaud-Delaître, C. (2014). Understanding semantic
865 mapping evolution by observing changes in biomedical ontologies. *Journal of biomedical*
866 *informatics, 47*, 71-82. [doi:10.1016/j.jbi.2013.09.006](https://doi.org/10.1016/j.jbi.2013.09.006)

867 DuCharme, B. (2013). *Learning SPARQL: querying and updating with SPARQL 1.1*. O'Reilly Media, Inc.
868 ISBN:978-1449306595.

869 Ekins, S., Mietchen, D., Coffee, M., Stratton, T. P., Freundlich, J. S., Freitas-Junior, L., et al. (2016).
870 Open drug discovery for the Zika virus. *F1000Research, 5*, 150.
871 [doi:10.12688/f1000research.8013.1](https://doi.org/10.12688/f1000research.8013.1).

872 Erxleben, F., Günther, M., Krötzsch, M., Mendez, J., & Vrandečić, D. (2014). Introducing Wikidata to
873 the Linked Data Web. *The Semantic Web – ISWC 2014* (pp. 50–65). Springer International
874 Publishing. [doi:10.1007/978-3-319-11964-9_4](https://doi.org/10.1007/978-3-319-11964-9_4).

875 Färber, M., Bartscherer, F., Menne, C., & Rettinger, A. (2018). Linked data quality of DBpedia, Freebase,
876 OpenCyc, Wikidata, and YAGO. *Semantic Web, 9*(1), 77–129. [doi:10.3233/SW-170275](https://doi.org/10.3233/SW-170275).

877 Farda-Sarbas, M., Zhu, H., Nest, M. F., & Müller-Birn, C. (2019). Approving automation: analyzing
878 requests for permissions of bots in wikidata. In *Proceedings of the 15th International Symposium*
879 *on Open Collaboration* (pp. 1-10). doi:10.1145/3306446.3340833.

880 Gangemi, A., Catenacci, C., Ciaramita, M., & Lehmann, J. (2005). A theoretical framework for ontology
881 evaluation and validation. In *SWAP* (Vol. 166, p. 16).
882 http://www.loa.istc.cnr.it/old/Papers/swap_final_v2.pdf.

883 Gyrard, A., Datta, S. K., & Bonnet, C. (2018). A survey and analysis of ontology-based software tools for
884 semantic interoperability in IoT and WoT landscapes. *2018 IEEE 4th World Forum on Internet of*
885 *Things (WF-IoT)*, (pp. 86–91). doi:10.1109/WF-IoT.2018.8355091.

886 Hadj Taieb, M. A., Zesch, T., & Ben Aouicha, M. (2020). A survey of semantic relatedness evaluation
887 datasets and procedures. *Artificial Intelligence Review*, 53(6), 4407-4448. doi:10.1007/s10462-
888 019-09796-3.

889 Halper, M., Gu, H., Perl, Y., & Ochs, C. (2015). Abstraction networks for terminologies: supporting
890 management of “big knowledge”. *Artificial intelligence in medicine*, 64(1), 1-16.
891 doi:10.1016/j.artmed.2015.03.005

892 Hanika, T., Marx, M., & Stumme, G. (2019). Discovering implicational knowledge in Wikidata. In
893 *International Conference on Formal Concept Analysis* (pp. 315-323). Springer, Cham.
894 doi:10.1007/978-3-030-21462-3_21.

895 Harris, S., Seaborne, A., & Prud'hommeaux, E. (2013). SPARQL 1.1 query language. *W3C*
896 *recommendation*, 21(10), 778.

897 Heftberger, A., Höper, J., Müller-Birn, C., & Walkowski, N.-O. (2020). Opening up Research Data in
898 Film Studies by Using the Structured Knowledge Base Wikidata. In H. Kremers, *Digital Cultural*
899 *Heritage* (pp. 401–410). Springer International Publishing. doi:10.1007/978-3-030-15200-0_27.

900 Heymann, D. L. (2020). Data sharing and outbreaks: best practice exemplified. *The Lancet*, 395 (10223),
901 469-470. doi: 10.1016/S0140-6736(20)30184-7.

902 Jalalifard, M., Norouzi, Y., & Isfandyari-Moghaddam, A. (2013). Analyzing web citations availability
903 and half-life in medical journals. *Aslib Proceedings*, 65(3), 242.
904 doi:10.1108/00012531311330638.

905 Jantzen, S. G., Sutherland, B. J., Minkley, D. R., & Koop, B. F. (2011). GO Trimming: Systematically
906 reducing redundancy in large Gene Ontology datasets. *BMC research notes*, 4(1), 267.
907 doi:10.1186/1756-0500-4-267.

908 Jemielniak, D. (2014). *Common knowledge?: An ethnography of Wikipedia*. Stanford: Stanford
909 University Press. ISBN:978-0804789448

910 Jemielniak, D., & Wilamowski, M. (2017). Cultural diversity of quality of information on Wikipedias.
911 *Journal of the Association for Information Science and Technology*, 68(10), 2460-2470.
912 doi:10.1002/asi.23901.

913 Jemielniak, D., & Przegalinska, A. (2020) *Collaborative Society*, Cambridge, MA: MIT Press. ISBN:978-
914 0262537919.

915 Kaffee, L. A., Piscopo, A., Vougiouklis, P., Simperl, E., Carr, L., & Pintscher, L. (2017). A glimpse into
916 babel: An analysis of multilinguality in wikidata. *Proceedings of the 13th International*
917 *Symposium on Open Collaboration* (p. 14). ACM. doi:10.1145/3125433.3125465.

918 Kaffee, L.-A., & Simperl, E. (2018). Analysis of Editors' Languages in Wikidata. *Proceedings of the 14th*
919 *International Symposium on Open Collaboration* (p. 21). ACM. doi:10.1145/3233391.3233965

920 Kagan, D., Moran-Gilad, J., & Fire, M. (2020). Scientometric trends for coronaviruses and other
921 emerging viral infections. *GigaScience*, 9(8), g1aa085. doi:10.1093/gigascience/g1aa085.

922 Knublauch, H., & Kontokostas, D. (2017, 6). *Shapes Constraint Language (SHACL)*, W3C
923 *Recommendation 20 July 2017*. W3C Recommendation, #w3c#. Retrieved from
924 <https://www.w3.org/TR/2017/REC-shacl-20170720/>

925 Krishnan, L., Ogunwole, S. M., & Cooper, L. A. (2020). Historical Insights on Coronavirus Disease 2019
926 (COVID-19), the 1918 Influenza Pandemic, and Racial Disparities: Illuminating a Path Forward.
927 *Annals of Internal Medicine*, 173(6), 474-481. doi:10.7326/M20-2223.

928 Labra Gayo, J. E., & Alvarez Rodríguez, J. M. (2013). Validating statistical index data represented in
929 RDF using SPARQL queries. *RDF Validation Workshop. Practical Assurances for Quality RDF*
930 *Data*. Cambridge: <http://www.w3.org/2012/12/rdf-val>.

931 Labra Gayo, J. E., Prud'Hommeaux, E., Boneva, I., & Kontokostas, D. (2017). Validating RDF data.
932 *Synthesis Lectures on Semantic Web: Theory and Technology*, 7(1), 1-328.
933 doi:10.2200/s00786ed1v01y201707wbe016.

934 Labra-Gayo, J. E., García-González, H., Fernández-Alvarez, D., & Prud'hommeaux, E. (2019).
935 Challenges in RDF validation. In *Current Trends in Semantic Web Technologies: Theory and*
936 *Practice* (pp. 121-151). Springer, Cham. doi:10.1007/978-3-030-06149-4_6.

937 Lastra-Díaz, J. J., Goikoetxea, J., Hadj Taieb, M. A., García-Serrano, A., Ben Aouicha, M., & Agirre, E.
938 (2019). A reproducible survey on word embeddings and ontology-based methods for word
939 similarity: linear combinations outperform the state of the art. *Engineering Applications of*
940 *Artificial Intelligence*, 85, 645-665. doi:10.1016/j.engappai.2019.07.010.

941 Lampoltshammer, T. J., & Heistracher, T. (2014). Ontology evaluation with Protégé using OWLET.
942 *Infocommunications Journal*, 6(2), 12-17. [https://www.researchgate.net/profile/Thomas-](https://www.researchgate.net/profile/Thomas-Lampoltshammer/publication/263692985_Ontology_evaluation_with_Protege_using_OWLET/links/00b4953bced7997952000000/Ontology-evaluation-with-Protege-using-OWLET.pdf)
943 [Lampoltshammer/publication/263692985_Ontology_evaluation_with_Protege_using_OWLET/](https://www.researchgate.net/profile/Thomas-Lampoltshammer/publication/263692985_Ontology_evaluation_with_Protege_using_OWLET/links/00b4953bced7997952000000/Ontology-evaluation-with-Protege-using-OWLET.pdf)
944 [links/00b4953bced7997952000000/Ontology-evaluation-with-Protege-using-OWLET.pdf](https://www.researchgate.net/profile/Thomas-Lampoltshammer/publication/263692985_Ontology_evaluation_with_Protege_using_OWLET/links/00b4953bced7997952000000/Ontology-evaluation-with-Protege-using-OWLET.pdf).

945 Lanamäki, A., & Lindman, J. (2018). Latent Groups in Online Communities: a Longitudinal Study in
946 Wikipedia. *Computer Supported Cooperative Work (CSCW)*, 27(1), 77-106. doi:10.1007/s10606-
947 017-9295-8.

948 Lee, D., Cornet, R., Lau, F., & De Keizer, N. (2013). A survey of SNOMED CT implementations.
949 *Journal of biomedical informatics*, 46(1), 87-96. doi:10.1016/j.jbi.2012.09.006.

950 Li, J., Sun, Y., Johnson, R. J., Sciaky, D., Wei, C. H., Leaman, R., et al. (2016). BioCreative V CDR task
951 corpus: a resource for chemical disease relation extraction. *Database*, 2016.
952 doi:10.1093/database/baw068.

953 Li, Q., Guan, X., Wu, P., Wang, X., Zhou, L., Tong, Y., et al. (2020). Early transmission dynamics in
954 Wuhan, China, of novel coronavirus-infected pneumonia. *New England Journal of Medicine*,
955 382, 1199-1207. doi:10.1056/NEJMoa2001316.

956 Lozano-Tello, A., & Gomez-Perez, A. (2004). Ontometric: A Method to Choose the Appropriate
957 Ontology. *Journal of Database Management (JDM)*, 15(2), 1-18.
958 <https://www.igi-global.com/article/ontometric-method-choose-appropriate-ontology/3308>.

959 Luggen, M., Difallah, D., Sarasua, C., Demartini, G., & Cudré-Mauroux, P. (2019). Non-parametric Class
960 Completeness Estimators for Collaborative Knowledge Graphs—The Case of Wikidata. *The*
961 *Semantic Web – ISWC 2019* (pp. 453–469). Springer International Publishing. doi:10.1007/978-
962 3-030-30793-6_26.

963 Luo, L., Mejino Jr, J. L., & Zhang, G. Q. (2013). An analysis of FMA using structural self-bisimilarity.
964 *Journal of biomedical informatics*, 46(3), 497-505. doi:10.1016/j.jbi.2013.03.005.

965 Malyshev, S., Krötzsch, M., González, L., Gonsior, J., & Bielefeldt, A. (2018). Getting the most out of
966 wikidata: Semantic technology usage in wikipedia’s knowledge graph. *International Semantic*
967 *Web Conference* (pp. 376-394). Springer, Cham. doi:10.1007/978-3-030-00668-6_23.

968 Martin, P. A. (2018). Evaluating Ontology Completeness via SPARQL and Relations-between-Classes
969 Based Constraints. *11th International Conference on the Quality of Information and*
970 *Communications Technology (QUATIC)*, (pp. 255–263). doi:10.1109/QUATIC.2018.00045.

971 Marx, M., & Krötzsch, M. (2017). SQID: Towards Ontological Reasoning for Wikidata. In *Proceedings*
972 *of the ISWC 2017 Posters & Demonstrations Track*. CEUR Workshop Proceedings.
973 <https://iccl.inf.tu-dresden.de/web/Inproceedings3169/en>.

974 Mayr, P., Scharnhorst, A., Larsen, B., Schaer, P., & Mutschke, P. (2014). Bibliometric-enhanced
975 information retrieval. *European Conference on Information Retrieval* (pp. 798-801). Springer,
976 Cham. doi:10.1007/978-3-319-06028-6_99.

977 Melo, A., & Paulheim, H. (2020). Automatic detection of relation assertion errors and induction of
978 relation constraints. *Semantic Web*, 11(5), 801-830. doi:10.3233/SW-200369.

979 Mietchen, D., Hagedorn, G., Willighagen, E., Rico, M., Gómez-Pérez, A., Aibar, E., Rafes, K., Germain,
980 C., Dunning, A., Pintscher, L., & Kinzler, D. (2015). Enabling open science: Wikidata for
981 research (Wiki4R). *Research Ideas and Outcomes*, 1, e7573. doi: 10.3897/rio.1.e7573.

982 Mietchen, D., & Li, J. (2020). Quantifying the Impact of Data Sharing on Outbreak Dynamics
983 (QIDSOD). *Research Ideas and Outcomes*, 6, e54770. doi: 10.3897/rio.6.e54770.

984 Miquel-Ribé, M., & Laniado, D. (2018). Wikipedia Culture Gap: Quantifying Content Imbalances Across
985 40 Language Editions. *Frontiers in Physics*, 6, 54. doi:10.3389/fphy.2018.00054.

986 Mitraka, E., Waagmeester, A., Burgstaller-Muehlbacher, S., Schriml, L. M., Su, A. I., & Good, B. M.
987 (2015). Wikidata: A platform for data integration and dissemination for the life sciences and
988 beyond. *bioRxiv*, 031971. doi:10.1101/031971.

989 Mora-Cantalops, M., Sánchez-Alonso, S., & García-Barriocanal, E. (2019). A systematic literature
990 review on Wikidata. *Data Technologies and Applications*, 53, 250–268. doi:10.1108/DTA-12-
991 2018-0110.

992 Mortensen, J. M., Minty, E. P., Januszyk, M., Sweeney, T. E., Rector, A. L., Noy, N. F., & Musen, M. A.
993 (2014). Using the wisdom of the crowds to find critical errors in biomedical ontologies: a study of
994 SNOMED CT. *Journal of the American Medical Informatics Association*, 22, 640-648.
995 doi:10.1136/amiajnl-2014-002901.

996 Müller-Birn, C., Karran, B., Lehmann, J., & Luczak-Rösch, M. (2015). Peer-production System or
997 Collaborative Ontology Engineering Effort: What is Wikidata? *Proceedings of the 11th*
998 *International Symposium on Open Collaboration* (pp. 20:1–20:10). New York: ACM.
999 doi:10.1145/2788993.2789836.

- 1000 Nielsen, F. Å., Mietchen, D., & Willighagen, E. (2017). Scholia, scientometrics and wikidata. In
1001 *European Semantic Web Conference* (pp. 237-259). Springer, Cham. doi:10.1007/978-3-319-
1002 70407-4_36.
- 1003 Nielsen, F. Å., Thornton, K., & Labra-Gayo, J. E. (2019). Validating Danish Wikidata lexemes. In *15th*
1004 *International Conference on Semantic Systems, SEMPDS 2019*. Karlsruhe: CEUR-WS.
- 1005 Obrst, L., Ceusters, W., Mani, I., Ray, S., & Smith, B. (2007). The Evaluation of Ontologies. *Semantic*
1006 *Web*, 139–158. doi:10.1007/978-0-387-48438-9_8.
- 1007 Ostaszewski, M., Mazein, A., Gillespie, M. E., Kuperstein, I., Niarakis, A., Hermjakob, H., et al. (2020).
1008 COVID-19 Disease Map, building a computational repository of SARS-CoV-2 virus-host
1009 interaction mechanisms. *Scientific data*, 7(1), 136. doi:10.1038/s41597-020-0477-8.
- 1010 Park, M. S., He, Z., Chen, Z., Oh, S., & Bian, J. (2016). Consumers' use of UMLS concepts on social
1011 media: diabetes-related textual data analysis in blog and social Q&A sites. *JMIR medical*
1012 *informatics*, 4(4), e41. doi:10.2196/medinform.5748.
- 1013 Paulheim, H. (2017). Knowledge graph refinement: A survey of approaches and evaluation methods.
1014 *Semantic Web*, 8(3), 489-508. doi:10.3233/SW-160218.
- 1015 Pellissier Tanon, T., & Suchanek, F. (2019). Querying the Edit History of Wikidata. *The Semantic Web:*
1016 *ESWC 2019 Satellite Events* (pp. 161–166). Springer International Publishing. doi:978-3-030-
1017 32327-1_32.
- 1018 Pellissier Tanon, T., Bourgaux, C., & Suchanek, F. (2019). Learning how to correct a knowledge base
1019 from the edit history. In *The World Wide Web Conference* (pp. 1465-1475).
1020 doi:10.1145/3308558.3313584.
- 1021 Pérez, J., Arenas, M., & Gutierrez, C. (2009). Semantics and complexity of SPARQL. *ACM Transactions*
1022 *on Database Systems (TODS)*, 34(3), 16. doi:10.1145/1567274.1567278.
- 1023 Piad-Morffis, A., Gutiérrez, Y., & Muñoz, R. (2019). A corpus to support ehealth knowledge discovery
1024 technologies. *Journal of biomedical informatics*, 94, 103172. doi:10.1016/j.jbi.2019.103172.
- 1025 Pillai, S., Soon, L.-K., & Haw, S.-C. (2019). Comparing DBpedia, Wikidata, and YAGO for Web
1026 Information Retrieval. *Intelligent and Interactive Computing* (pp. 525–535). Springer Singapore.
1027 doi:10.1007/978-981-13-6031-2_40.
- 1028 Piscopo, A., & Simperl, E. (2018). Who Models the World?: Collaborative Ontology Creation and User
1029 Roles in Wikidata. *Proceedings of the ACM on Human-Computer Interaction*, 2(CSCW), 141:1–
1030 141:18. doi:10.1145/3274410.
- 1031 Prud'hommeaux, E., Labra Gayo, J. E., & Solbrig, H. (2014). Shape Expressions: An RDF Validation and
1032 Transformation Language. In *Proceedings of the 10th International Conference on Semantic*
1033 *Systems* (pp. 32-40). doi:10.1145/2660517.2660523
- 1034 Raad, J., & Cruz, C. (2015). A survey on ontology evaluation methods. *Proceedings of the International*
1035 *Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge*
1036 *Management* (pp. 179-186). ACM. doi:10.5220/0005591001790186.
- 1037 RDA COVID-19 Working Group (2020). *RDA COVID-19; recommendations and guidelines, 5th release*
1038 *28 May 2020*. Research Data Alliance. doi:10.15497/RDA00046.
- 1039 Rector, A. L., Brandt, S., & Schneider, T. (2011). Getting the foot out of the pelvis: modeling problems
1040 affecting use of SNOMED CT hierarchies in practical applications. *Journal of the American*
1041 *Medical Informatics Association*, 18(4), 432-440. doi:10.1136/amiajnl-2010-000045.

1042 Rector, A., & Iannone, L. (2012). Lexically suggest, logically define: Quality assurance of the use of
1043 qualifiers and expected results of post-coordination in SNOMED CT. *Journal of biomedical*
1044 *informatics*, 45(2), 199-209. doi:10.1016/j.jbi.2011.10.002.

1045 Rothman, K. J., Greenland, S., & Lash, T. L. (2008). *Modern epidemiology*. Lippincott Williams &
1046 Wilkins. ISBN:978-1451190052.

1047 Salas, J., & Hogan, A. (2022). Semantics and Canonicalisation of SPARQL 1.1. *Semantic Web*.
1048 doi:10.3233/SW-212871.

1049 Samuel, J. (2017). Collaborative Approach to Developing a Multilingual Ontology: A Case Study of
1050 Wikidata. *Research Conference on Metadata and Semantics Research* (pp. 167-172). Springer.
1051 doi:10.1007/978-3-319-70863-8_16.

1052 Samuel, J. (2021, April). ShExStatements: Simplifying Shape Expressions for Wikidata. In *Companion*
1053 *Proceedings of the Web Conference 2021* (pp. 610-615). ACM. doi:10.1145/3442442.3452349.

1054 Sarabadani, A., Halfaker, A., & Taraborelli, D. (2017, April). Building automated vandalism detection
1055 tools for Wikidata. In *Proceedings of the 26th International Conference on World Wide Web*
1056 *Companion* (pp. 1647-1654). ACM. doi:10.1145/3041021.3053366.

1057 Sarasua, C., Checco, A., Demartini, G., Difallah, D., Feldman, M., & Pintscher, L. (2019). The evolution
1058 of power and standard Wikidata editors: comparing editing behavior over time to predict lifespan
1059 and volume of edits. *Computer Supported Cooperative Work (CSCW)*, 28(5), 843-882.
1060 doi:10.1007/s10606-018-9344-y.

1061 Schober, D., Tudose, I., Svatek, V., & Boeker, M. (2012). OntoCheck: verifying ontology naming
1062 conventions and metadata completeness in Protégé 4. *Journal of Biomedical Semantics*, 3(Suppl
1063 2), S4. doi:10.1186/2041-1480-3-S2-S4.

1064 Sebei, H., Taieb, M. A. H., & Aouicha, M. B. (2018). Review of social media analytics process and big
1065 data pipeline. *Social Network Analysis and Mining*, 8(1), 30. doi:10.1007/s13278-018-0507-0.

1066 Shafee, T., Masukume, G., Kipersztok, L., Das, D., Häggström, M., & Heilman, J. (2017). Evolution of
1067 Wikipedia's medical content: past, present and future. *J Epidemiol Community Health*, 71(11),
1068 1122-1129. doi:10.1136/jech-2016-208601.

1069 Shorland, A., Mietchen, D., & Willighagen, E. (2020). *Wikidata Queries around the SARS-CoV-2 virus*
1070 *and pandemic*. NL: Zenodo. doi:10.5281/zenodo.3977414.

1071 Thornton, K., Solbrig, H., Stupp, G. S., Labra Gayo, J. E., Mietchen, D., Prud'Hommeaux, E., &
1072 Waagmeester, A. (2019). Using Shape Expressions (ShEx) to share RDF data models and to
1073 guide curation with rigorous validation. *European Semantic Web Conference* (pp. 606-620).
1074 Springer. doi:10.1007/978-3-030-21348-0_39.

1075 Turki, H. (2018). Citation analysis is also useful to assess the eligibility of biomedical research works for
1076 inclusion in living systematic reviews. *Journal of clinical epidemiology*, 97, 124-125.
1077 doi:10.1016/j.jclinepi.2017.11.002.

1078 Turki, H., Hadj Taieb, M. A., & Ben Aouicha, M. (2018). MeSH qualifiers, publication types and relation
1079 occurrence frequency are also useful for a better sentence-level extraction of biomedical relations.
1080 *Journal of biomedical informatics*, 83, 217-218. doi:10.1016/j.jbi.2018.05.011.

1081 Turki, H., Shafee, T., Hadj Taieb, M. A., Ben Aouicha, M., Vrandečić, D., Das, D., & Hamdi, H. (2019).
1082 Wikidata: A large-scale collaborative ontological medical database. *Journal of Biomedical*
1083 *Informatics*, 99, 103292. doi:10.1016/j.jbi.2019.103292.

- 1084 Turki, H., Vrandečić, D., Hamdi, H., & Adel, I. (2017). Using WikiData as a Multi-lingual Multi-dialectal
1085 Dictionary for Arabic Dialects. *2017 IEEE/ACS 14th International Conference on Computer*
1086 *Systems and Applications (AICCSA)* (pp. 437–442). IEEE. doi:10.1109/AICCSA.2017.115.
- 1087 Turki, H., Hadj Taieb, M. A., Ben Aouicha, M., & Abraham, A. (2020). Nature or Science: what Google
1088 Trends says. *Scientometrics*, *124*(2), 1367-1385. doi:10.1007/s11192-020-03511-8.
- 1089 Turki, H., Hadj Taieb, M. A., & Ben Aouicha, M. (2021). Coupling Wikipedia Categories with Wikidata
1090 Statements for Better Semantics. In *Proceedings of the 2nd Wikidata Workshop*
1091 *(Wikidata@ISWC 2021)* (pp. 8:1-8:6). <http://ceur-ws.org/Vol-2982/paper-8.pdf>.
- 1092 Vanderkam, D., Schonberger, R., Rowley, H., & Kumar, S. (2013). *Nearest neighbor search in google*
1093 *correlate*. Google Inc. <https://research.google/pubs/pub41694/>.
- 1094 Vasanthapriyan, S., Tian, J., & Xiang, J. (2017). An Ontology-Based Knowledge Framework for
1095 Software Testing. *Knowledge and Systems Sciences* (pp. 212–226). Springer Singapore.
1096 doi:10.1007/978-981-10-6989-5_18.
- 1097 Vrandečić, D. (2009). Ontology Evaluation. In R. S. S. Staab, *Handbook on Ontologies* (pp. 293–313).
1098 Berlin, Heidelberg: Springer Berlin Heidelberg. doi:10.1007/978-3-540-92673-3_13.
- 1099 Vrandečić, D., & Krötzsch, M. (2014). Wikidata: a free collaborative knowledgebase. *Communications of*
1100 *the ACM*, *57*(10), 78-85. doi:10.1145/2629489.
- 1101 Vrandečić, D. (2021). Building a multilingual Wikipedia. *Communications of the ACM*, *64*(4), 38-41.
1102 doi:10.1145/3425778.
- 1103 Waagmeester, A., Schriml, L., & Su, A. I. (2019). Wikidata as a linked-data hub for Biodiversity data.
1104 *Biodiversity Information Science and Standards*, *3*, e35206. doi:10.3897/biss.3.35206.
- 1105 Waagmeester, A., Willighagen, E. L., Su, A. I., Kutmon, M., Gayo, J. E. L., Fernández-Álvarez, D., et al.
1106 (2021). A protocol for adding knowledge to Wikidata: aligning resources on human
1107 coronaviruses. *BMC biology*, *19*(1), 12:1-12:14. doi:10.1186/s12915-020-00940-y
- 1108 Waagmeester, A., Stupp, G., Burgstaller-Muehlbacher, S., Good, B. M., Malachi, G., Griffith, O. L., et al.
1109 (2020b). Wikidata as a knowledge graph for the life sciences. *eLife*, *9*, e52614.
1110 doi:10.7554/eLife.52614.
- 1111 Ward, A., & Murray-Ward, M. (1999). *Assessment in the classroom*. Wadsworth Publishing Company.
1112 ISBN:978-0534527044.
- 1113 Walisadeera, A. I., Ginige, A., & Wikramanayake, G. N. (2016). Ontology Evaluation Approaches: A
1114 Case Study from Agriculture Domain. *Computational Science and Its Applications -- ICCSA*
1115 *2016* (pp. 318–333). Springer International Publishing. doi:10.1007/978-3-319-42089-9_23.
- 1116 Wasi, S., Sachan, M., & Darbari, M. (2020). Document Classification Using Wikidata Properties.
1117 *Information and Communication Technology for Sustainable Development* (pp. 729–737).
1118 Singapore: Springer. doi:10.1007/978-981-13-7166-0_73.
- 1119 Wiśniewski, D., Potoniec, J., Ławrynowicz, A. & Keet, C. M. (2019). Analysis of Ontology Competency
1120 Questions and their formalizations in SPARQL-OWL. *Journal of Web Semantics*, *59*, 100534.
1121 doi:10.1016/j.websem.2019.100534.
- 1122 Xu, B., Kraemer, M. U., & Data Curation Group (2020). Open access epidemiological data from the
1123 COVID-19 outbreak. *The Lancet Infectious Diseases*, *20*(5), 534. doi:10.1016/S1473-
1124 3099(20)30119-5.

1125 Zangerle, E., Gassler, W., Pichl, M., Steinhauser, S., & Specht, G. (2016). An Empirical Evaluation of
1126 Property Recommender Systems for Wikidata and Collaborative Knowledge Bases. *Proceedings*
1127 *of the 12th International Symposium on Open Collaboration* (pp. 18:1–18:8). New York: ACM.
1128 doi:10.1145/2957792.2957804.

1129 Zhang, G. Q., & Bodenreider, O. (2010). Large-scale, exhaustive lattice-based structural auditing of
1130 SNOMED CT. In *AMIA Annual Symposium Proceedings* (Vol. 2010, p. 922). American Medical
1131 Informatics Association. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3041382/>.

1132 Zhang, Y., Lin, H., Yang, Z., Wang, J., Zhang, S., Sun, Y., & Yang, L. (2018). A hybrid model based on
1133 neural networks for biomedical relation extraction. *Journal of biomedical informatics*, *81*, 83-92.
1134 doi:10.1016/j.jbi.2018.03.011

1135 Zhang, Y., Chen, Q., Yang, Z., Lin, H., & Lu, Z. (2019). BioWordVec, improving biomedical word
1136 embeddings with subword information and MeSH. *Scientific data*, *6*(1), 52:1-52:9.
1137 doi:10.1038/s41597-019-0055-0.

1138 Zu, Z. Y., Jiang, M. D., Xu, P. P., Chen, W., Ni, Q. Q., Lu, G. M., & Zhang, L. J. (2020). Coronavirus
1139 disease 2019 (COVID-19): a perspective from China. *Radiology*, *296*(2), E15-E25.
1140 doi:10.1148/radiol.2020200490.

1141

1142 **Appendix A: SPARQL queries for the heuristics-based validation of**
1143 **epidemiological counts in Wikidata**

Task	SPARQL query
V1	<pre>SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:<PropertyID> [ps:<PropertyID> ?value; pq:P585 ?date]. FILTER(YEAR(?date) < 2019) }</pre>
V2	<pre>SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:<PropertyID> [ps:<PropertyID> ?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:<PropertyID> [ps:<PropertyID> ?value; pq:P585 ?datep]. ?x p:<PropertyID> [ps:<PropertyID> ?value1; pq:P585 ?date]. FILTER(?value > ?value1) FILTER(?datep - ?date = -1) }</pre>
V4	<pre>SELECT * WHERE { ?x p:P31 [ps:P31 wd:Q3241045; pq:P642 wd:Q84263196]. ?x p:<PropertyID> [ps:<PropertyID> ?value; pq:P585 ?datep]. ?x p:<PropertyID> [ps:<PropertyID> ?value1; pq:P585 ?datef]. FILTER(?value = ?value1) FILTER(?datep - ?datef = -2) FILTER NOT EXISTS { ?x p:<PropertyID> [ps:<PropertyID> ?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:<PropertyID> [ps:<PropertyID> ?value; pq:P585 ?date]. FILTER(?value < 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 > ?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 >= ?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 < ?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 < ?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:<PropertyID> [ps:<PropertyID> ?c; pq:P585 ?date]. ?x wdt:P361 ?y. ?y p:<PropertyID> [ps:<PropertyID> ?c1; pq:P585 ?date]. } } GROUP BY ?y ?c1 ?date } ORDER BY DESC(?diff)</pre>

1144 The SPARQL queries that were used for the Tasks defined in Table 2, to be run against the Wikidata
1145 Query Service available at <https://query.wikidata.org/> . Note that this query service has Wikidata-
1146 specific prefixes predefined, so they do not need to be re-stated in a query.

1147

1148

1149 **Appendix B: SPARQL queries for the validation of case fatality rate**

1150 **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 > 1) (?value < 0)) } </pre>
M2	<pre> SELECT ?x ?c ?d ?value ?date (ABS(?value - ?d / ?c) > 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]. } } 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>

1151 These SPARQL queries correspond to the Tasks M1, M2 and M3 that address heuristics
1152 concerning the case fatality rate m .