# Wikidata and COVID-19

**Creating a collaborative knowledge graph from CORD-19 scholarly publications**

Houcemeddine Turki, *University of Sfax, Tunisia*

جامعة صفاقس
University of Sfax

W3C®

AI2 Allen Institute for AI

# Team Members
## WikiProject COVID-19

Houcemeddine Turki, *University of Sfax, Tunisia*
Thomas Shafee, *La Trobe University, Australia*
Daniel Mietchen, *University of Virginia, United States of America*
Tiago Lubiana, *University of São Paolo, Brazil*
Dariusz Jemielniak, *Kozminski University, Poland*
Jose Emilio Labra Gayo, *University of Oviedo, Spain*
Eric Prud'Hommeaux, *World Wide Web Consortium, United States of America*
Mohamed Ali Hadj Taieb, *University of Sfax, Tunisia*
Mohamed Ben Aouicha, *University of Sfax, Tunisia*
Mus'ab Banat, *Hashemite University, Jordan*
Diptanshu Das, *Institute of Child Health, India*

## University of Sfax

- Top 2 in Tunisia

- Ranked among the first African universities in Computer Science (1st in Leiden CWTS, 2nd in URAP)

- Research Group specialized in Semantic Technologies and Biomedical Data Science (Data Engineering and Semantics)

- Technopark dealing with computer science research
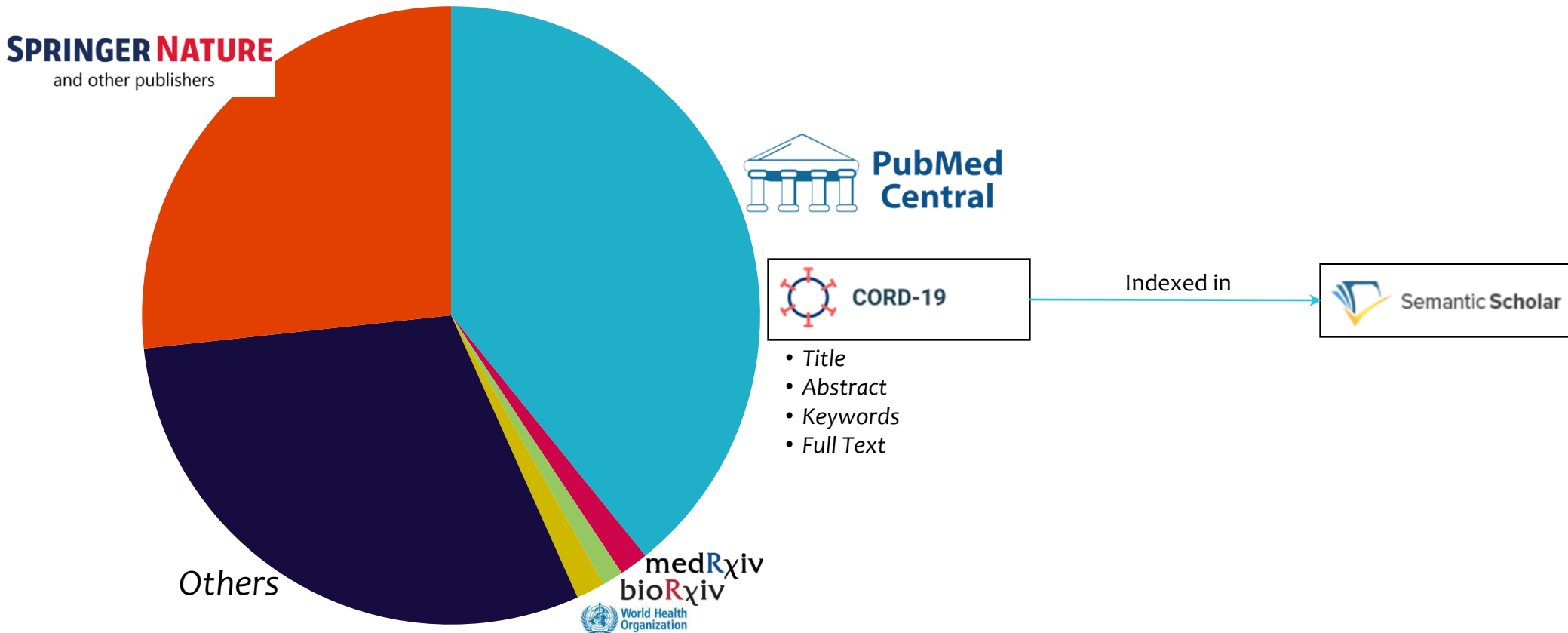
# About Us

WikiProject COVID-19
University of Sfax

# Introduction

Insights on Covid-19 Open Research Dataset (CORD-19)

# CORD-19:  COVID-19 Open Research Dataset

An initiative of *AI2*, *CZI*, *MSR*, *Georgetown*, *NIH*, and *The White House*



PubMed Central

CORD-19

Indexed in

Semantic Scholar

- *Title*
- *Abstract*
- *Keywords*
- *Full Text*

SPRINGER NATURE
and other publishers

*Others*

medRχiv
bioRχiv

World Health Organization

# CORD-19

**The sum of all human knowledge about COVID-19**

## A huge amount of raw texts

- Difficult to study by humans

- Hard to process by computer programs

- Knowledge updated every day

# Proposed Solution

A knowledge graph for COVID-19 information

- A fully structured semantic database in the form of RDF triples.

- Human-readable, Machine-readable

- Findable, Accessible, Interoperable and Reusable

- Flexible data model for the representation of COVID-19 information

- Screened using SPARQL

# Wikidata

A large-scale free knowledge base

- Available at https://www.wikidata.org

- Items and properties are assigned language-independent identifiers and labels and descriptions in multiple languages

- Items and properties are assigned statements in the form of RDF triples. These statements can be detailed using triple qualifiers and references.

- Statements can be relational (taxonomic or non-taxonomic) or non-relational ones (objects as values, external IDs, URLs and dates…).

- CC0 License (easily reusable but cannot include various datasets released under CC-BY and other licenses).

# Comparison

Collaborative multidisciplinary knowledge graphs vs. Specialized knowledge graphs

## Collaborative multidisciplinary knowledge graphs
### *Wikidata*

- Include already available non-COVID-19 items and statements
  - Allow correlation analysis between COVID-19 information and non-COVID-19 information
  - Data models already existing in part
- Edited and curated by a large community of editors

## Specialized knowledge graphs
### *CORD-19 NEKG*

- Does not include already available non-COVID-19 items and statements
  - Only deep analysis of COVID-19 information
  - Data models developed from scratch
- Edited and curated by a panel of experts
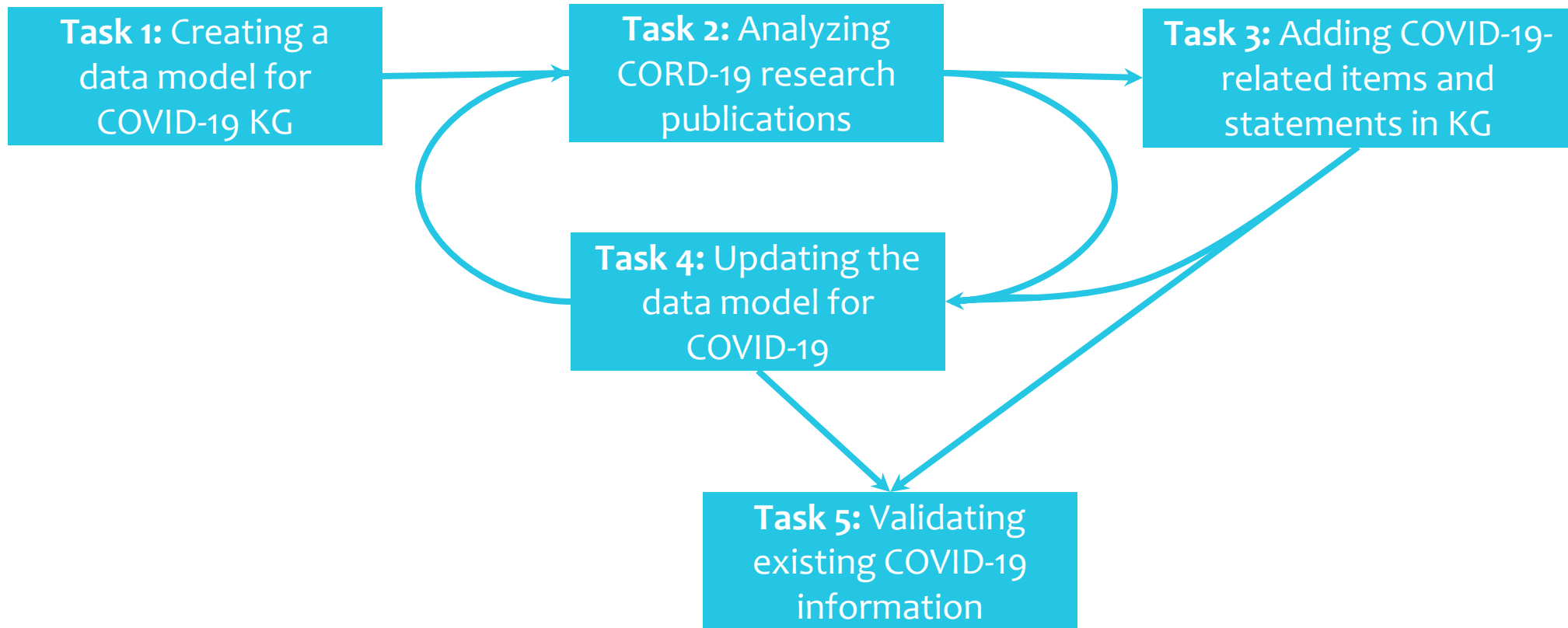
# Methods

Creating a large-scale COVID-19 knowledge graph in Wikidata

# Roadmap

Creation of knowledge graphs from CORD-19

**Task 1:** Creating a data model for COVID-19 KG

**Task 2:** Analyzing CORD-19 research publications

**Task 3:** Adding COVID-19-related items and statements in KG

**Task 4:** Updating the data model for COVID-19

**Task 5:** Validating existing COVID-19 information

# Creating a data model for COVID-19 KG

Task 1

- Defining the classes for COVID-19-related items (https://www.wikidata.org/wiki/Special:NewItem)

- Defining the structure of the items of each class using ShEx (https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Data_models)

- Adding and sustaining Wikidata properties to characterize COVID-19 statements (https://www.wikidata.org/wiki/Wikidata:Property_proposal)

- CORD-19 NEKG: https://github.com/Wimmics/CovidOnTheWeb/blob/master/doc/01-data-modeling.md (YouTube: https://youtu.be/oUk9PXGM2fY)

# Analyzing CORD-19 research publications

Task 2

- Human screening of CORD-19 research publications on Semantic Scholar and PubMed Central for the creation and adjustment of COVID-19-related statements on Wikidata.

جامعة صفاقس
University of Sfax

# Analyzing CORD-19 research publications

Task 2

- Processing CORD-19 research publications using Semantic Scholar API and PubMed Central Entrez API and corresponding Python Libraries.

# Analyzing CORD-19 research publications

Task 2

- Annotating CORD-19 research publications with named entities from semantic databases such as Wikidata.
  - Eliminating stop words
  - Extracting n-grams
  - Finding n-grams in knowledge graphs using APIs
- The annotation process can be contextualized and restricted to the items included in a given class of a knowledge graph (Diseases, drugs, etc…).

جامعة صفاقس
University of Sfax

Multivac: Annotation tool for COVID-19 scholarly publications

# Analyzing CORD-19 research publications

Task 2

- Using Topic Modelling (Latent Dirichlet Allocation) of CORD-19 scholarly publications to retrieve the main concepts related to COVID-19 pandemic.
- Redundancy of returned concepts can be solved using word similarity metrics such as semantic similarity measures and word embeddings.

# Analyzing CORD-19 research publications

Task 2

- Annotating CORD-19 research publications with relations from semantic databases such as Wikidata.

- Various methods:
  - Human screening and annotation of CORD-19 using a tool such as https://brat.nlplab.org/.
  - Using a benchmark of semantic relations of the same types (particularly drug interactions and drug-disease relations) with ML techniques (CNN, RNN and LSTM) or word embeddings (Word2Vec, BERT and ELMo) to annotate and retrieve semantic relations from CORD-19.
  - Inferring semantic relations for topic modelling of CORD-19.

# Analyzing CORD-19 research publications

Task 2

- Two data models of relation annotation:
  - SAT+R: *Subject*, *Action*, *Target* triplets and additional *Relations*
    - Limited number of relation types linking annotated entities
    - Supported relations are generic (Subject, Target, Has property…)
    - Semantic links between concepts (*Green*) are extracted from analyzed text and annotated as action entities (*Red*).

# Analyzing CORD-19 research publications

Task 2

- Two data models of relation annotation:
  - Subject-Object relation annotations
    - Every type of biomedical information is represented by a relation type
    - Only identified concepts are annotated in the analyzed text
    - Actions are represented as links between concepts and are not consequently annotated as entities.

جامعة صفاقس
University of Sfax

# Analyzing CORD-19 research publications

Task 2

- Using Bibliometric-Enhanced Information Retrieval to enhance the extraction of biomedical relations
  - MeSH Keywords
    - PubMed Records include MeSH Keywords describing the output of the corresponding research publication. MeSH Keywords involve a Heading describing a biomedical concept and a qualifier specifying the studied pattern in the Heading/Qualifier form.
    - A combination of two MeSH Keywords can be used to infer a given biomedical relation.
    - Example: "Hepatitis C/therapy" and "Sofosbuvir/therapeutic use" in https://pubmed.ncbi.nlm.nih.gov/32526210/ can be used to determine that Sofosbuvir is used as a drug for Hepatitis C.
  - Publication types
    - Section titles in literature reviews can provide an idea about the types of available relations in each section of the reviews.
    - Example: "Drug interactions" section in a literature review about "Azithromycin" involves semantic relations about significant drug interactions of Azithromycin.

جامعة صفاقس
University of Sfax

# Adding COVID-19-related items and statements in KG

Task 3

- Human creation of COVID-19-related items and statements
  - Item creation: https://www.wikidata.org/wiki/Special:NewItem
  - Statement creation: https://www.wikidata.org/wiki/<item>

جامعة صفاقس
University of Sfax

# Adding COVID-19-related items and statements in KG

Task 3

- Wikidata API (https://www.wikidata.org/w/api.php)
  - Mass extraction of COVID-19 multidisciplinary information
  - Mass adjustment and modification of COVID-19 multidisciplinary information

## MediaWiki API help

This is an auto-generated MediaWiki API documentation page.

Documentation and examples: https://www.mediawiki.org/wiki/Special:MyLanguage/API:Main_page

## Main module

[Documentation · FAQ · Mailing list · API Announcements · Bugs & requests]

**Status:** The MediaWiki API is a mature and stable interface that is actively supported and improved. While we try to avoid it, we may occasionally need to make breaking changes; subscribe to the mediawiki-api-announce mailing list for notice of updates.

- Source: MediaWiki
- License: GPL-2.0-or-later

**Erroneous requests:** When erroneous requests are sent to the API, an HTTP header will be sent with the key "MediaWiki-API-Error" and then both the value of the header and the error code sent back will be set to the same value. For more information see API: Errors and warnings.

**Testing:** For ease of testing API requests, see Special:ApiSandbox.

**Parameters:**

    action:        Which action to perform.

        **abusefiltercheckmatch:** Check to see if an AbuseFilter matches a set of variables, an edit, or a logged AbuseFilter event.

        **abusefilterchecksyntax:** Check syntax of an AbuseFilter filter.

جامعة صفاقس
University of Sfax

# Adding COVID-19-related items and statements in KG

Task 3

- QuickStatements (https://quickstatements.toolforge.org)
  - Upload batches of semantic relations to Wikidata
  - Can be programmatically used

# Adding COVID-19-related items and statements in KG

Task 3

- Wikidata Integrator (https://pypi.org/project/wikidataintegrator/)
  - A Python Library to analyze, add and adjust Wikidata statements

جامعة صفاقس
University of Sfax

# Adding COVID-19-related items and statements in KG

Task 3

- To apply a bot on Wikidata
  - Create a Python code
  - Publish it in a GitHub repository
  - Apply for a bot flag
  - Run the bot on a server

### Wikidata:Requests for permissions/Bot/RefB (WikiCred)

< Wikidata:Requests for permissions | Bot

The following discussion is closed. **Please do not modify it.** Subsequent comments should be made in a new section. A summary of the conclusions reached follows.

✓ Approved--Ymblanter (talk) 10:10, 2 September 2020 (UTC)

**RefB (WikiCred)**   [ edit ]

RefB (WikiCred) (talk · contribs · new items · SUL · Block log · User rights log · User rights · xtools)
**Operator:** Csisc (talk · contribs · logs)
**Task/s:** This bot will add reference support to biomedical statements in Wikidata.
**Code:** https://github.com/Data-Engineering-and-Semantics/refb/

**Function details:**

- This bot identifies unsupported biomedical relations on Wikidata using a SPARQL query.
- To find references supporting the extracted Wikidata statements, all that should be done is to use the PubMed Central search engine (based on Biopython and NCBI Entrez API) to find publications where the subject and the object of each statement co-occur. The algorithm will return the PMC ID of the reference for each assessed Wikidata statement and the sentence proving it within the full text of the reference.
- All we need to do is to convert PMC IDs into Wikidata IDs using Wikidata Hub, and then add the obtained references to Wikidata using the QuickStatements API.
- The source code of this bot is build using Python 3.5.
- Further details about the bot can be found here.

--Csisc (talk) 14:31, 29 July 2020 (UTC)

# Adding COVID-19-related items and statements in KG

Task 3

- Wikidata items are aligned to several open datasets and knowledge graphs particularly in the context of Linked Open Data Cloud
  - Other open knowledge graphs are automatically extracting COVID-19 information from CORD-19
  - Wikidata can integrate these information if the licenses of these open knowledge bases waive all the legal barriers (CC0 or Public Domain)



Legend:
- User Generated
- Life Science
- Publications
- Cross Domain
- Government
- Geography
- Linguistics
- Media
- Social Networking
- Other

# Updating the data model for COVID-19

Task 4

- Many methods:
  - Inferring new COVID-19 related classes from Topic Modelling of CORD-19 and mass import them to Wikidata using QuickStatements tool.
  - Deriving the data model of classes by analyzing the characteristics of COVID-19 knowledge in Wikidata using full screening and SPARQL queries.
  - Human updates of data models and ShEx validation schemas.



EntitySchema  Discussion                                    Read  View history   Search Wikidata

### COVID-19 dashboards, search engines and datasets (E205)

| language code | label | description | aliases | edit |
|---|---|---|---|---|
| en | COVID-19 dashboards, search engines and datasets | Entity schema of COVID-19 dashboards, search engines and datasets | | edit |

check entities against this Schema   edit

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

#Reference: https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Data_models/COVID-19_apps

start = @<app>

<app>  EXTRA wdt:P31  {
    wdt:P31 [ wd:Q90790055 wd:Q91136116 wd:Q91137337 ]; # instance of a COVID-19 dashboard, search engine or dataset
    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
}
```

# Validating existing COVID-19 information

Task 5

- Property constraints and statements
  - Define structural constraints for the definition of Wikidata statements
  - Identify the links between Wikidata properties

# Validating existing COVID-19 information

Task 5

| Constraint type | Description |
|---|---|
| single value constraint | Constraint used to specify that this property generally contains a single value per item |
| format constraint | Constraint used to specify that the value for this property has to correspond to a given pattern |
| 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 |
| distinct values constraint | Constraint used to specify that the value for this property is likely to be different from all other items |
| Commons link constraint | Constraint used to specify that the value must link to an existing Wikimedia Commons page |
| 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 |
| mandatory qualifier constraint | Constraint used to specify that the listed qualifier has to be used |
| symmetric constraint | Constraint used to specify that the referenced entity should also link back to this entity |
| used as qualifier constraint | Constraint used to specify that a property must only be used as a qualifier |
| value requires statement constraint | Constraint used to specify that the referenced item should have a statement with a given property |
| relation of type constraint | relation establishing dependency between types/metalevels of its members |
| allowed qualifiers constraint | Constraint used to specify that only the listed qualifiers should be used. Novalue disallows any qualifier |
| value type constraint | Constraint used to specify that the referenced item should be a subclass or instance of a given type |
| allowed units constraint | Constraint used to specify that only listed units may be used |
| multi-value constraint | Constraint used to specify that a property generally contains more than one value per item |
| one-of constraint | Constraint used to specify that the value for this property has to be one of a given set of items |
| range constraint | Constraint used to specify that the value must be between two given values |

جامعة صفاقس
University of Sfax

# Validating existing COVID-19 information

Task 5

- Data Models and ShEx
  - Specify the required statements for the definition of a Wikidata item
  - Available at https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Data_models



Generic properties [edit]

| Title | ID | Data type | Description | Examples | Inverse |
|---|---|---|---|---|---|
| publisher | P123 | Item | publisher: organization or person responsible for publishing books, periodicals, games or software | Smittestopp <publisher> Norwegian Institute of Public Health | - |
| creator | P170 | Item | creator and author: maker of this creative work or other object (where no more specific property exists). Paintings with unknown painters, use "anonymous" (Q4233718) as value. series ordinal (P1545) can be added to identify the order of an individual in the list of the creators. | Systematic Platform for Essential and Epidemiological Data analysis of COVID-19 <creator> Houcemeddine Turki | - |
| title | P1476 | Monolingual text | original title and title: published title of a work, such as a newspaper article, a literary work, a website, or a performance work | Smittestopp <title> Smittestopp (Norwegian Bokmål) | - |
| developer | P178 | Item | video game developer and software developer: organisation or person that developed the item | Smittestopp <developer> Simula Research Laboratory | - |
| use | P366 | Item | use: main use of the subject (includes current and former usage) | Smittestopp <use> contact tracing | - |
| country of origin | P495 | Item | country of origin: country of origin of this item (creative work, food, phrase, product, etc.) | Smittestopp <country of origin> Norway | - |
| language of work or name | P407 | Item | language: language associated with this creative work (such as books, shows, songs, or websites) or a name (for persons use "native language" (P103) and "languages spoken, written or signed" (P1412)) | Smittestopp <language of work or name> Norwegian | - |
| official website | P856 | URL | official website and home page: URL of the official homepage of an item (current or former) [if the homepage changes, add an additional statement with preferred rank. Do not remove the former URL] | Coronavirus Australia <official website> https://www.health.gov.au/resources/apps-and-tools/coronavirus-australia-app | - |
| publication date | P577 | Point in time | date of publication: date or point in time when a work was first published or released | Coronavirus Australia <publication date> 29 March 2020 | - |
| copyright license | P275 | Item | license: license under which this copyrighted work is released | Systematic Platform for Essential and Epidemiological Data analysis of COVID-19 <copyright license> CC0 | - |
| on focus list of Wikimedia project | P5008 | Item | WikiProject focus list: property to indicate that an item is of particular interest for a Wikimedia project. This property does not add notability. Items should not be created with this property if they are not notable for Wikidata. See also P6104, P972, P2354. | Smittestopp <on focus list of Wikimedia project> WikiProject COVID-19 | - |

COVID-19 dashboards, search engines and datasets (E205)

| language code | label | description | aliases | edit |
|---|---|---|---|---|
| en | COVID-19 dashboards, search engines and datasets | Entity schema of COVID-19 dashboards, search engines and datasets | | ✎edit |

check entities against this Schema | ✎edit

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

#Reference: https://www.wikidata.org/wiki/Wikidata:WikiProject_COVID-19/Data_models/COVID-19_apps

start = @<app>

<app>  EXTRA wdt:P31  {
    wdt:P31 [ wd:Q90790055 wd:Q91136116 wd:Q91137337 ]; # instance of a COVID-19 dashboard, search engine or dataset
    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
}
```

جامعة صفاقس
University of Sfax

جامعة صفاقس
University of Sfax

# Validating existing COVID-19 information

Task 5

- Logical constraints implemented in SPARQL to validate relational statements

| Constraint | Description |
|---|---|
| Defining the scheme of a Wikidata property | |
| T1 | Identify common use cases of $P$: ($C_S$,$C_O$) pairs |
| T2 | Identify inverse properties of $P$ corresponding to each common use case: ($C_S$, $R^{-1}$,$C_O$) statements |
| Identifying the deficiencies of the scheme | |
| T3 | For each returned $P^{-1}$, identify $P(S,O)$ relations supported by references and corresponding to the most common ($C_S$, $P^{-1}$, $C_O$) statement but not available in Wikidata |
| T4 | Identify $P(S,O)$ relations not corresponding to the most common scheme of $P$ |

- Logical constraints implemented in SPARQL to validate statistical statements

| Constraint | Description |
|---|---|
| Validating qualifiers of COVID-19 epidemiological statements | |
| V1 | Verify $Z$ as a date > November 01, 2019 |
| V2 | Verify $Q$ as any subclass of (P279*) of medical diagnosis (Q177719) |
| 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 $c_Z \leq c_{Z+1}$) |
| 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 |

# Validating existing COVID-19 information

Task 5

- Scholarly databases like CORD-19 should not only be used to extract COVID-19 information.

- Scholarly databases can be searched to find interesting references for unsupported statements in Wikidata

Extract unreferenced Wikidata statements

Identify the most relevant PubMed Central publications

Find the supporting sentence for claims

Align PMC ID with Wikidata ID of each reference

Add obtained references to Wikidata

# Results and Discussion

Current situation of COVID-19 semantic information in Wikidata

# Factors for the growth of COVID-19 multilingual coverage in Wikidata
Positive correlations between the language support for COVID-19 and a significant number of factors

| | Medical Wikipedia articles https://w.wiki/Z6a | | Medical Wikidata labels https://w.wiki/Z6h | | Wikipedia and Wikidata users https://w.wiki/Z6W | |
|---|---|---|---|---|---|---|
| Rank | Language | Number of medical articles | Language | Number of labels | Language | Number of users |
| 1 | English | 16670 | English | 65986 | English | 9600 |
| 2 | German | 8911 | French | 37053 | French | 2580 |
| 3 | Arabic | 8596 | German | 22432 | German | 2490 |
| 4 | French | 7258 | Spanish | 21505 | Spanish | 2330 |
| 5 | Spanish | 6979 | Arabic | 18581 | Russian | 1790 |
| 6 | Italian | 6498 | Italian | 18074 | Italian | 1430 |
| 7 | **Polish** | 6071 | Japanese | 17992 | Chinese | 1120 |
| 8 | Portuguese | 5652 | **Dutch** | 17985 | Japanese | 1090 |
| 9 | Russian | 5564 | Chinese | 17462 | Portuguese | 979 |
| 10 | Japanese | 4651 | Russian | 17165 | Arabic | 688 |

| | COVID Wikidata content https://w.wiki/ZSq | | COVID Wikipedia pages https://w.wiki/ZSt | | COVID Wikipedia edits https://covid-data.wmflabs.org/perProjectNoHumans | | COVID-19 pandemic Wikipedia pageviews https://w.wiki/ZTG | |
|---|---|---|---|---|---|---|---|---|
| Rank | Language | Number of labels | Language | Number of articles | Language | Number of edits | Language | Average daily pageviews |
| 1 | English | 1429 | English | 561 | English | 250306 | English | 52872 |
| 2 | **Dutch** | 785 | Arabic | 517 | German | 126359 | Russian | 41246 |
| 3 | Arabic | 623 | German | 431 | French | 42029 | Spanish | 37722 |
| 4 | **Catalan** | 579 | Portuguese | 427 | Chinese | 41545 | Chinese | 27598 |
| 5 | German | 561 | **Korean** | 408 | Spanish | 30869 | German | 20707 |
| 6 | French | 517 | Chinese | 396 | Arabic | 19963 | **Italian** | 8490 |
| 7 | Japanese | 503 | **Vietnamese** | 392 | Russian | 18719 | French | 7959 |
| 8 | Chinese | 483 | French | 379 | Japanese | 11508 | Portuguese | 7648 |
| 9 | Portuguese | 463 | Spanish | 370 | **Ukrainian** | 10599 | Japanese | 5227 |
| 10 | Spanish | 433 | **Indonesian** | 363 | **Hebrew** | 10386 | Arabic | 4300 |

# External databases aligned to Wikidata items

Scholarly research publications and clinical trials

| Wikidata ID | Wikidata Property | Count |
|---|---|---|
| P356 | DOI | 45101 |
| P698 | PubMed ID | 42294 |
| P6179 | Dimensions Publication ID | 16944 |
| P932 | PMCID | 12590 |

Diseases and clinical signs

| Wikidata ID | Wikidata Property | Diseases count | Symptoms count |
|---|---|---|---|
| P672 | MeSH tree code | 40 | 12 |
| P2892 | UMLS CUI | 38 | 11 |
| P494 | ICD-10 | 32 | 8 |
| P4229 | ICD-10-CM | 32 | 1 |
| P3827 | JSTOR topic ID | 32 | 10 |

# External databases aligned to Wikidata items

Humans and sovereign states

| Wikidata ID | Wikidata Property | Sovereign states | Humans |
|---|---|---|---|
| P214 | VIAF ID | 159 | 654 |
| P7859 | WorldCat Identities ID | 146 | 548 |
| P244 | Library of Congress authority ID | 125 | 458 |
| P213 | ISNI | 100 | 443 |
| P646 | Freebase ID | 124 | 379 |
| P2002 | Twitter username | 16 | 353 |

Other items

| Wikidata Class | Wikidata ID | Wikidata Property | Count |
|---|---|---|---|
| drug [Q11173] | P6689 | MassBank accession ID | 44 |
| drug [Q11173] | P4964 | SPLASH | 31 |
| protein [Q8054] | P638 | PDB structure ID | 31 |
| film [Q11424] | P345 | IMDb ID | 25 |

جامعة صفاقس
University of Sfax

Proteins encoded by SARS-CoV-2 genes (note that some genes encode multiple proteins) and the currently known human protein interaction partners (live data: https://w.wiki/beR).

SARS-CoV-2 interactions with the human proteome

جامعة صفاقس
University of Sfax

A

B

| disease | diseaseLabel | symptom_count | symptoms |
|---|---|---|---|
| Q wd:Q21396183 | arsenic pentoxide exposure | 12 | headache // abdominal pain // brain diseases // respiratory failure // cough // dyspnea // nausea // fever // anorexia // diarrhea // delirium // conjunctivitis |
| Q wd:Q706845 | Lassa fever | 10 | headache // fatigue // cough // abdominal pain // nausea // brain diseases // fever // myalgia // diarrhea // conjunctivitis |
| Q wd:Q21173341 | cadmium dust exposure | 9 | headache // anemia // cough // nausea // dyspnea // chills // anosmia // myalgia // diarrhea |
| Q wd:Q21173343 | cadmium oxide exposure | 9 | headache // anemia // cough // nausea // dyspnea // chills // anosmia // myalgia // diarrhea |
| Q wd:Q51993 | Ebola hemorrhagic fever | 8 | headache // nausea // dyspnea // fever // myalgia // diarrhea // conjunctivitis // abdominal pain |
| Q wd:Q21167939 | benzene exposure | 8 | headache // fatigue // abdominal pain // nausea // dyspnea // respiratory failure // anorexia // diarrhea |

A) Currently listed symptoms of COVID-19, with qualifiers indicating their frequency. (live data: https://w.wiki/N8f).

B) Other medical conditions sorted by the number of shared symptoms with COVID-19. (live data: https://w.wiki/bqV; adapted from https://scholia.toolforge.org/disease/Q84263196)

**Symptoms of COVID-19 and similar conditions**

A) Correlation between the current number of cases and mortality rates in every country, calculated from numeric summary data for each region. Countries coloured randomly (live data: https://w.wiki/bf$).

B) Age distribution of notable persons who have died of COVID-19 (blue), compared to the death age distribution for people who were born after 1901 (green), calculated from individual dates of birth and death (live data: https://w.wiki/be7 and https://w.wiki/but).

Summary epidemiological data on the COVID-19 pandemic

جامعة صفاقس
University of Sfax

A) Common words and word combinations (ngrams) in the titles of publications (live data: https://w.wiki/cFu).

B) Co-occurrence of topics in publications with one of the COVID-related items as a topic, with ribbon widths proportional to the number of publications sharing those topics (log scale). Topics coloured by group as determined by louvain clustering, topics shared in fewer than 5 publications omitted (interactive version: https://csisc.github.io/WikidataCOVID19SPARQL/Fig8B.html; live data: https://w.wiki/bww).

COVID-19 publication topics

A

| organization | organizationLabel | bankruptcyDate | countryLabel | inception | industries | parents | subsiduaries |
|---|---|---|---|---|---|---|---|
| Q wd:Q2208025 | STA Travel | 20 August 2020 | Germany | 1 January 1979 | tourism industry | DKSH | |
| Q wd:Q7606770 | Stein Mart | 12 August 2020 | | 1 January 1902 | retail | | |
| Q wd:Q5206569 | DW Sports Fitness | 3 August 2020 | United Kingdom | 1 January 2009 | retail | | |
| Q wd:Q2749082 | Lord & Taylor | 2 August 2020 | United States of America | 1 January 1826 | retail | | |
| Q wd:Q64059182 | Le Tote | 2 August 2020 | | 1 January 2012 | clothing, sharing economy | | |
| Q wd:Q3305660 | Tailored Brands | 2 August 2020 | United States of America | 1 January 1973 | retail | | Men's Wearhouse |
| Q wd:Q15109854 | California Pizza Kitchen | 30 July 2020 | United States of America | 1 January 1985 | hospitality industry | | |

B

C

A) Tabular output of SPARQL query
B) Bankruptcies per month
C) ratios of different industries associated with bankrupt companies. (live data: https://w.wiki/cG6).

Bankrupt publicly listed businesses due to the COVID-19 pandemic

جامعة صفاقس
University of Sfax

# Validation using logical constraints

Relational statements



Statistical statements

| cases | deaths | recoveries | tests | hospitalizations | Overall |
|---|---|---|---|---|---|
| 2856 | 2467 | 189 | 9 | 10 | 5496 |

**An infrastructure for knowledge graph validation based on interactions between consistency rules, property statements and RDF validation languages**

جامعة صفاقس
University of Sfax

# Limitations

- Several aspects of COVID-19 information can be more represented in other knowledge graphs. The inclusion of these COVID-19 information is blocked by the CC0 License of Wikidata (e.g. ORKG)

- Several types of information are still not supported by Wikidata (e.g. Structured outcomes of COVID-19 scholarly publications)

- Several aspects of COVID-19 information are more considered in textual resources such as Wikipedia than in Wikidata and other knowledge graphs (e.g. https://en.wikipedia.org/wiki/COVID-19_pandemic_on_cruise_ships)

# Applications

COVID-19 KG-driven applications

# Multilingual topic modelling of social interactions

*In progress*

# Knowledge-Based Systems

Education and Social Recommendation

# Conclusion

Take-Home Messages

# Conclusion

- Wikidata and other collaborative multidisciplinary knowledge graphs can create more efficiently a semantic database for COVID-19.

- Despite their slight limitations, Collaborative multidisciplinary knowledge graphs like Wikidata can return interesting findings about COVID-19 due to the integration of COVID-19 multidisciplinary information with non-COVID-19 information.

- Due to its interesting coverage of COVID-19, Wikidata can be used for a variety of applications using semantic web tools.

جامعة صفاقس
University of Sfax

# To cite the work

- **Main Work:**
  - Turki, H., Shafee, T., Hadj Taieb, M. A., Ben Aouicha, M., Vrandečić, D., Das, D., & Hamdi, H. (2019). Wikidata: A large-scale collaborative ontological medical database. *Journal of biomedical informatics*, 99, 103292. doi:10.1016/j.jbi.2019.103292.
  - Turki, H., Hadj Taieb, M. A., Shafee, T., Lubiana, T., Jemielniak, D., Ben Aouicha, M., Labra Gayo, J. E., Banat, M., Das, D., & Mietchen, D. (2020). Representing COVID-19 information in collaborative knowledge graphs: a study of Wikidata. *Zenodo*. doi:10.5281/zenodo.4028482.
  - Turki, H., Jemielniak, D., Hadj Taieb, M. A., Labra Gayo, J. E., Ben Aouicha, M., Banat, M., Shafee, T., Prud'Hommeaux, E., Lubiana, T., Das, D., & Mietchen, D. (2020). Using logical constraints to validate information in collaborative knowledge graphs: a study of COVID-19 on Wikidata. *Zenodo*. doi:10.5281/zenodo.4008358.
  - Waagmeester, A., Willighagen, E. L., Su, A. I., Kutmon, M., Gayo, J. E. L., Fernández-Álvarez, D., ... & Koehorst, J. J. (2020). A protocol for adding knowledge to Wikidata, a case report. *BioRxiv*. doi:10.1101/2020.04.05.026336.
  - Waagmeester, A., Stupp, G., Burgstaller-Muehlbacher, S., Good, B. M., Griffith, M., Griffith, O. L., ... & Keating, S. M. (2020). Science Forum: Wikidata as a knowledge graph for the life sciences. *ELife*, 9, e52614. doi:10.7554/eLife.52614.

- **Applications:**
  - Xianxian, H. (2019). Knowledge Graph (KG) for Recommendation System. *Medium*. https://medium.com/@hxianxian/knowledge-graph-kg-for-recommendation-system-8fe2c6cd354.

# References

- Vrandečić, D., & Krötzsch, M. (2014). Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, *57*(10), 78-85. doi: 10.1145/2629489.
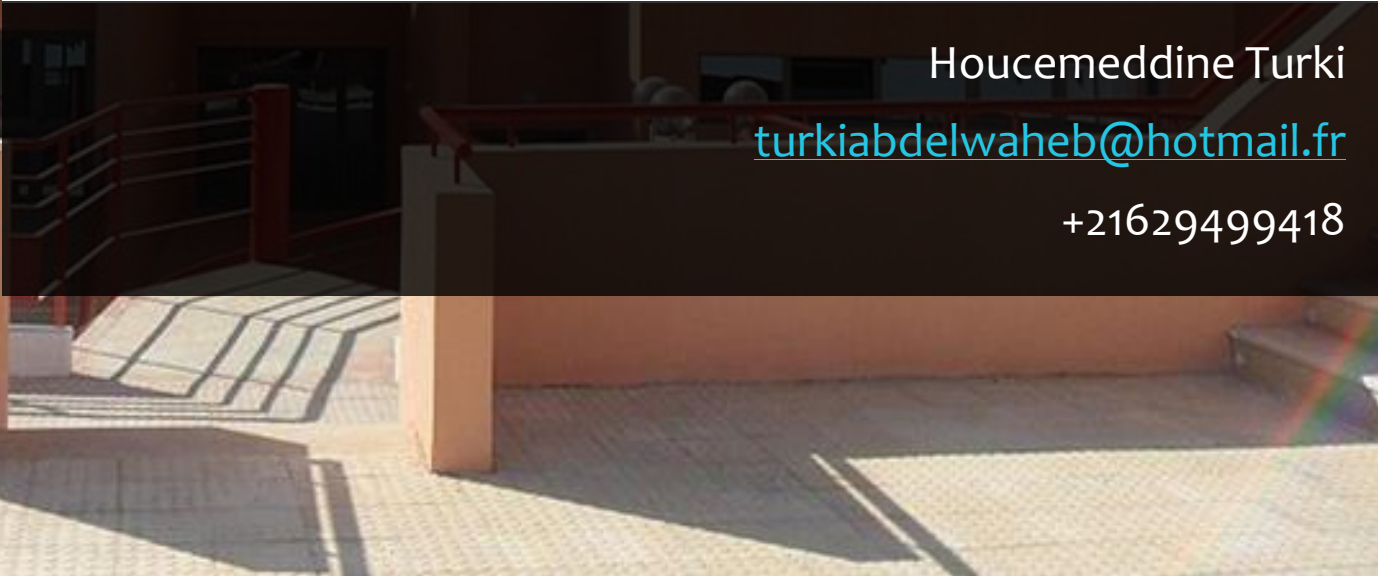
- Wimmics Research Team (2020). Covid-on-the-Web dataset. *Zenodo*. doi:10.5281/zenodo.3833752.

- Michel, F., Gandon, F., Ah-Kane, V., Bobasheva, A., Cabrio, E., Corby, O., … & Simon, M. (2020, November). Covid-on-the-Web: Knowledge Graph and Services to Advance COVID-19 Research. In *International Semantic Web Conference*.

- CNRS (2020). Multivac Platform. *GitHub*. https://github.com/multivacplatform.

- Colavizza, G., Costas, R., Traag, V. A., Van Eck, N. J., Van Leeuwen, T., & Waltman, L. (2020). A scientometric overview of CORD-19. *BioRxiv*. doi:10.1101/2020.04.20.046144.

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

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

- Peng, Y., Yan, S., & Lu, Z. (2019, August). Transfer Learning in Biomedical Natural Language Processing: An Evaluation of BERT and ELMo on Ten Benchmarking Datasets. In *Proceedings of the 18th BioNLP Workshop and Shared Task* (pp. 58-65).

- Piad-Morffis, A., Estevez-Velarde, S., Estevanell-Valladares, E. L., Gutiérrez, Y., Montoyo, A., Muñoz, R., & Almeida-Cruz, Y. (2020). Knowledge Discovery in COVID-19 Research Literature. *OpenReview*. https://openreview.net/forum?id=CWfGhEFOTKU

- Wu, Y., Liu, M., Zheng, W. J., Zhao, Z., & Xu, H. (2012). Ranking gene-drug relationships in biomedical literature using latent dirichlet allocation. In *Biocomputing 2012* (pp. 422-433). doi:10.1142/9789814366496_0041.

- Piad-Morffis, A., Gutiérrez, Y., & Muñoz, R. (2019). A corpus to support ehealth knowledge discovery technologies. *Journal of biomedical informatics*, *94*, 103172. doi:10.1016/j.jbi.2019.103172.

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

# Credits

- https://commons.wikimedia.org/wiki/File:Du_cot%C3%A9_du_bab_diwan.jpg

- https://creativecommons.org/licenses/by/4.0/

- http://www.webdo.tn/2018/10/16/luniversite-de-sfax-1ere-en-tunisie-801eme-dans-le-monde/

- https://commons.wikimedia.org/wiki/File:A_Rainbow_Of_Books_-_Flickr_-_Dawn_Endico.jpg

- Colavizza, G., Costas, R., Traag, V. A., Van Eck, N. J., Van Leeuwen, T., & Waltman, L. (2020). A scientometric overview of CORD-19. *BioRxiv*. doi:10.1101/2020.04.20.046144.

- https://commons.wikimedia.org/wiki/Category:COVID-19_Study_of_Wikidata

- Xianxian, H. (2019). Knowledge Graph (KG) for Recommendation System. *Medium*. https://medium.com/@hxianxian/knowledge-graph-kg-for-recommendation-system-8fe2c6cd354.

- https://commons.wikimedia.org/wiki/File:Int%C3%A9rieur_2_du_centre_de_recherche,_Technopole_de_Sfax.jpg

# Thank You

Houcemeddine Turki

turkiabdelwaheb@hotmail.fr

+21629499418

- Appendix A: Algorithm for the development of CORD-19 NEKG

- Appendix B: RDF data format for Wikidata

- Appendix C: Wikidata prefixes

- Appendix D: Useful links

- Appendix E: Hidden Tools - FM3S

- Appendix F: Hidden Tools - SNOWL Model

# Appendices

Useful information

جامعة صفاقس
University of Sfax

# Appendix A: Algorithm for the development of CORD-19 NEKG

https://www.inria.fr/fr/covid-web

جامعة صفاقس
University of Sfax

# Appendix B: RDF data format for Wikidata

# Appendix C: Wikidata prefixes

```
PREFIX wd: <http://www.wikidata.org/entity/>
PREFIX wds: <http://www.wikidata.org/entity/statement/>
PREFIX wdv: <http://www.wikidata.org/value/>
PREFIX wdt: <http://www.wikidata.org/prop/direct/>
PREFIX wikibase: <http://wikiba.se/ontology#>
PREFIX p: <http://www.wikidata.org/prop/>
PREFIX ps: <http://www.wikidata.org/prop/statement/>
PREFIX pq: <http://www.wikidata.org/prop/qualifier/>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX bd: <http://www.bigdata.com/rdf#>
```

# Appendix D: Useful links

| Page | URL |
|------|-----|
| Statistics | https://www.wikidata.org/wiki/Special:Statistics |
| Database Reports | https://www.wikidata.org/wiki/Wikidata:Database_reports |
| Database download | https://www.wikidata.org/wiki/Wikidata:Database_download |
| User access levels | https://www.wikidata.org/wiki/Wikidata:User_access_levels |
| Wikidata Tour | https://www.wikidata.org/wiki/Wikidata:Tours |
| Wikidata Tools | https://www.wikidata.org/wiki/Wikidata:Tools |
| Wikidata Hub | https://hub.toolforge.org/ |
| Scholia tool (Scholarly information) | https://scholia.toolforge.org/ |
| SPEED tool (Epidemiological information) | https://speed.ieee.tn/ |
| SARS-CoV-2-Queries (Genomics) | https://egonw.github.io/SARS-CoV-2-Queries/ |
| Bot flag for RefB (Bot adding references to biomedical statements) | https://www.wikidata.org/wiki/Wikidata:Requests_for_permissions/Bot/RefB_(WikiCred) |
| Source code for RefB | https://github.com/Data-Engineering-and-Semantics/refb/ |

# Appendix E: Hidden Tools - FM3S

A measure of sentence-level semantic similarity

Springer Link

International Conference on Hybrid Artificial Intelligence Systems

HAIS 2015: Hybrid Artificial Intelligent Systems pp 515-529 | Cite as

## FM3S: Features-Based Measure of Sentences Semantic Similarity

Authors | Authors and affiliations

Mohamed Ali Hadj Taieb ✉, Mohamed Ben Aouicha, Yosra Bourouis

جامعة صفاقس
University of Sfax

# Appendix F: Hidden Tools - SNOWL Model

An ontology for the alignment between the namespaces and entity types in social media



SpringerLink

Search  Log in

Regular Paper | Published: 27 July 2020

## SNOWL model: social networks unification-based semantic data integration

Hiba Sebei ✉, Mohamed Ali Hadj Taieb & Mohamed Ben Aouicha

*Knowledge and Information Systems* (2020) | Cite this article

**66** Accesses | **2** Altmetric | Metrics

Download PDF ⬇

Sections  Figures  References

Abstract
Introduction
Related works
Ontology model of social networks data
SNOWL ontology evaluation
Ontology deployment
Conclusion and Future work
Notes
References
Author information
Additional information
Rights and permissions

## Abstract

Integrating social networks data in the process of promoting business and marketing applications is widely addressed by several researchers. However, regarding the isolation between social network platforms managing such data has become a challenging task facing data scientist. In this respect, the present paper is designed to put forward a special semantic data integration approach, whereby a unified presentation and access to social networks data can be maintained. To this end, the novel SNOWL (Social Network OWL) ontology aims to

جامعة صفاقس
University of Sfax