# The COVID-19 epidemiology and monitoring ontology

**Núria Queralt-Rosinach[1], Paul Schofield[2], Robert Hoehndorf[3], Claus Weiland[4], Erik Schultes[5], César Henrique Bernabé[1], and Marco Roos[1]**

**1** Leiden University Medical Center, The Netherlands **2** University of Cambridge. United Kingdom **3** King Abdullah University of Science and Technology, Saudi Arabia **4** Senckenberg Biodiversity and Climate Research Center, Germany **5** GO FAIR Foundation

corresponding author: n.queralt_rosinach@lumc.nl

## Motivation

One year ago, the novel COVID-19 infectious disease emerged and spread, causing high mortality and morbidity rates worldwide. In the OBO Foundry, there are more than one hundred ontologies to share and analyse large-scale datasets for biological and biomedical sciences. However, this pandemic revealed that we lack tools for an efficient and timely exchange of this epidemiological data which is necessary to assess the impact of disease outbreaks, the efficacy of mitigating interventions and to provide a rapid response (Editorial, 2021). Recently, several new COVID-19 ontologies have developed such as the Infectious Disease Ontology (IDO) ("IDO-COVID19 OWL ontology," n.d.) extension or Ontology of Coronavirus Infectious Disease (CIDO) (He, Yu, & Ong, 2020). Hence, our research question was to determine if there was a good representation of epidemiological quantitative concepts in OBO ontologies. Our objectives were to identify missing COVID-19 epidemiological terms and implement axiom patterns for extensions to existing ontologies or to build a new, logically well-formed, and accurate ontology in OBO. In this study we present our findings and contributions for the bio-ontologies community.

## Method

This work was conceived and mainly developed during open community hackathons COVID-19 Virtual BioHackathon 2020, BioHackathon Europe 2020, SWAT4HCLS Hackathon 2020. Our approach was based on first, extracting a list of relevant epidemiological terms through manual curation of recent COVID-19 epidemiological studies published in peer-reviewed journals, medRxiv and public health surveillance websites, and mapping them to existing OBO ontologies. Curation was focused on quantitative data and indicators. Second, developing a minimal ontological representation of COVID-19 epidemiological quantitative information. And third, to refine and evaluate the model with domain expert input.

Our formal modeling followed a rationale already used in other studies: 1) determine the domain and scope of the ontology; 2) ontology reuse and addressing poor ontological coverage of COVID-19 epidemiology; and 3) development of a conceptual model (Sánchez & Batet, 2011, p. @kouameasi2021). We extracted core domain knowledge concepts from existing sources(al., 2004, pp. @epibook2, @idehandbook2014). We re-used ontological terms and models as much as possible using ontology search engines such as the Ontology Lookup Service (OLS) and the Open Biological and Biomedical Ontology (OBO) Foundry. To build an interoperable biomedical ontology, we decided to build an OBO ontology and use the OWL 2, a DL-based formalism and semantic web standard for knowledge representation, to enable data sharing and formal reasoning. We used knowledge-engineering best practices following the OBO principles and modularization guidelines (Rector, n.d.) to achieve a logically well-formed model. Finally, we based our decisions on building a FAIR (findable, accessible, interoperable, reusable)

(Wilkinson et al., 2016) resource for health data and research following recent recommendations published by international data standard organizations (Group, n.d., p. @ehden2020). More information on the method, the list of sources used for curation and extracted terms, and the developed OWL ontology are open and publicly available for reproducibility and community re-use on https://github.com/NuriaQueralt/covid19-epidemiology-ontology.

## Results

We provide a formal ontological model for COVID-19 epidemiology and monitoring (graphical and OWL representations are in our GitHub repository). With the rise of new variants of the virus that may challenge vaccine efficacy, a compatible logical model for quantities that enables researchers to represent and share machine-readable patient monitoring and epidemiology surveillance data for rapid analysis, modeling and response is an urgent need. In this work we re-used the SIO design pattern for measurements, a model already applied to patient health data for rare diseases in the European Joint Programme on Rare Diseases (EJP RD), clinical research data in the Leiden University Medical Center (LUMC) (al., n.d.) and the measurements schema in the new GA4GH Phenopackets release (GA4GH, n.d.). The taxonomic structure is extended from IDO, a core ontology for infectious diseases. For domain concepts we re-used General Formal Ontology (GFO) (Herre et al., 2006) to formalize timelines concepts using the 'chronoid' and the GFO-based 'mortality' model approach (F, F, R, Z, & D., 2012). To link patient-population is an RDA COVID-19 recommendation on data sharing, thus we checked common data models such as (OMOP) and re-used the relationship used in Phenopackets based on *composition* semantics.

We filled the gap for epidemiological surveillance terms in OBO adding 100 new terms. From an initial set of 138 manually extracted terms, only 38 are covered by bio-ontologies, 21% (30 terms) IDO ("Infectious Disease Ontology OWL ontology," n.d.) and 24% Statistical Methods Ontology (STATO) (33 terms) ("STATO: the statistical methods ontology OWL ontology," n.d.) (although including fallbacks this percentage could increase to 50%) and the rest by epidemiological-related ontologies such as Apollo Structure Vocabulary (APOLLO_SV) ("Apollo structured vocabulary owl ontology," n.d.) and Genomic Epidemiology Ontology (GenEpiO) ("Genomic Epidemiology Ontology OWL ontology," n.d.). We noticed that the Epidemiology Ontology (EPO) ("Epidemiology Ontology OBO Foundry Website," n.d.) is not maintained since its publication and has been deprecated from OBO Foundry and IDO is working towards epidemiological enrichment (Babcock, Beverley, Cowell, & Smith, 2021). While interoperability within the OBO landscape is fostered by adopting the BFO backbone structure, the link with GFO can lead to incompatible temporal regions due to logical inconsistency ("Toward semantic interoperability with linked foundationalontologies in ROMULUS," n.d.). Another issue that may be improved is the current absence of axioms and definition patterns that relate epidemiology (i.e., observations of a population) to clinical ontologies (i.e., observations on an individual) and allow reasoning for discovery. The re-use of the Entity-Quality(EQ) model (CJ et al., 2011) or the adaptation of the Resources, events, agents (REA) model (Mabee PM, 2020) will be evaluated. In the future, we will evaluate our ontology with domain experts and logical competency questions (Almeida Falbo, 2014). Moreover, we expect to use this model in FAIR-based projects such as Trusted World of Corona (TWOC) (funder, n.d.) to publish epidemiological claims as nanopublications for trust (Groth, Gibson, & Velterop, 2010). We aim at FAIR reasoning and analytics of person-level real world observations over epidemiological surveillance information (Sherimon, 2020). Therefore, checking common data models such as Phenopackets or OHDSI standards was done to enable the development of applications to discover patterns with ontology-guided machine learning algorithms and translational research.

## Conclusion

In the context of an infectious disease outbreak it is imperative to have these data as FAIR as possible to facilitate rapid analysis and support timely evidence-based decision making and trust. To enable the community to provide machine-readable epidemiological quantitative data and make it easier to share, we contributed with the development of an ontological representation, which was built based on ontology engineering best-practices such as reuse and ontology formalization through upper-level ontologies (i.e., GFO, SIO).

## Jupyter notebooks, GitHub repositories and data repositories

- Github repo. covid-19-epidemiology-ontology, https://github.com/NuriaQueralt/covid19-epidemiology-ontology, CC0

## Paper data sources

- Paper.md and paper.bib location:

https://github.com/NuriaQueralt/covid19-epidemiology-ontology

## Acknowledgements

## References

al., F. M. N. et. (2004). *Vigilancia epidemiológica*. McGraw-Hill Interamericana.

al., N. Q.-R. et. (n.d.). FAIR data management to access patient data. https://repository.publisso.de/resource/frl%3A6424232.

Almeida Falbo, R. de. (2014). SABiO: Systematic approach for building ontologies. *FOIS*.

Apollo structured vocabulary owl ontology. (n.d.).

Babcock, S., Beverley, J., Cowell, L. G., & Smith, B. (2021). The infectious disease ontology in the age of covid-19. *Preprint*. doi:https://doi.org/10.31219/osf.io/az6u5

Bochove, K. van, Vos, E., Moinat, M., Sandijk, S. van, Korthout, T., & Mohtashani, P. (n.d.). EHDEN - d4.5 - roadmap for interoperability solutions. https://zenodo.org/record/4474373#.YI66yq4p5hH.

CJ, M., M, B., TZ, B., J, D., A, I., MA, H., DP, H., et al. (2011). Cross-product extensions of the gene ontology. *J Biomed Inform*, *44*(1), 80–6. doi:doi:10.1016/j.jbi.2010.02.002

Editorial. (2021). How epidemiology has shaped the covid pandemic. *Nature*, (589), 491–492. doi:https://doi.org/10.1038/d41586-021-00183-z

Epidemiology Ontology OBO Foundry Website. (n.d.). http://www.obofoundry.org/ontology/epo.html.

F, S., F, F., R, F., Z, M., & D., S. (2012). Towards an ontological representation of morbidity and mortality in description logics. *J Biomed Semantics*, *Suppl 2(Suppl 2)*(3), S7. doi:doi:10.1186/2041-1480-3-S2-S7

funder, H.-H. (n.d.). The trusted world of corona (twoc). Retrieved from /url%7Bhttps://www.health-holland.com/project/2020/trusted-world-of-corona%7D

GA4GH. (n.d.). Phenopackets v2. https://github.com/phenopackets/phenopacket-schema/issues/261.

Genomic Epidemiology Ontology OWL ontology. (n.d.). http://purl.obolibrary.org/obo/genepio/releases/2020-08-09/genepio.owl.

Groth, P., Gibson, A., & Velterop, J. (2010). The anatomy of a nanopublication'. *Information Services & Use*, *30*(1-2), 51–56. doi:10.3233/ISU-2010-0613

Group, R. C.-1. W. (n.d.). RDA covid-19 recommendations and guidelines on data sharing. https://zenodo.org/record/3932953#.YI7Dba4p5hH.

He, Y., Yu, H., & Ong, E. et a. (2020). CIDO, a community-based ontology for coronavirus disease knowledge and data integration, sharing, and analysis. *Sci Data*, *7*(181). doi:https://doi.org/10.1038/s41597-020-0523-6

Herre, H., Heller, B., Burek, P., Hoehndorf, R., Loebe, F., & Michalek, H. (2006). *General formal ontology (gfo) - a foundational ontology integrating objects and processes [version 1.0]* (WorkingPaper). Unknown Publisher; Unknown Publisher.

IDO-COVID19 OWL ontology. (n.d.). http://purl.obolibrary.org/obo/2020-21-07/ido-covid-19.owl/.

Infectious Disease Ontology OWL ontology. (n.d.). http://purl.obolibrary.org/obo/ido/2017-11-03/ido.owl.

Kouamé, K.-M., & Mcheick, H. (2021). Ontological approach for early detection of suspected covid-19 among copd patients. *Appl. Syst. Innov.*, (4), 21. doi:https://doi.org/10.3390/asi4010021

Mabee PM, D. W., Balhoff JP. (2020). A logical model of homology for comparative biology. *Syst Biol.*, *1*(69), 345–362. doi:doi:10.1093/sysbio/syz067

MacMahon, B., & Trichopoulos, D. (n.d.). Harvard Medical School of Public Health, Boston, Massachussetts: Marbán SL.

Rector, A. L. (n.d.). Modularisation of domain ontologies implemented in description logics and related formalisms including owl. http://www.cs.man.ac.uk/~rector/papers/rector-modularisation-kcap-2003-distrib.pdf.

S, S.-B., R, R., & M, K. (2014). Infectious disease epidemiology.y. *Handbook of Epidemiolog*, 2041–2119. doi:doi:10.1007/978-0-387-09834-0_34

Sánchez, D., & Batet, M. (2011). Semantic similarity estimation in the biomedical domain: An ontology-based information-theoretic perspective. *J. Biomed. Inform*, (44), 749–759.

Sherimon, V. et al. (2020). Covid-19 ontology engineering-knowledge modeling of severe acute respiratory syndrome coronavirus 2 (sars-cov-2). *(IJACSA) International Journal of Advanced Computer Science and Applications*, *11*(11).

STATO: the statistical methods ontology OWL ontology. (n.d.). http://purl.obolibrary.org/obo/stato.owl.

Toward semantic interoperability with linked foundationalontologies in ROMULUS. (n.d.). https://researchspace.csir.co.za/dspace/bitstream/handle/10204/7042/Khan_2013.pdf?sequence=1&isAllowed=y.

Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., et al. (2016). The FAIR guiding principles for scientific data management and stewardship. *Scientific data*, *3*.