

CovKG: A Covid-19 Knowledge Graph for Enabling Multidimensional Analytics on Covid-19 Epidemiological Data considering Spatiotemporal, Environmental, Health, and Socioeconomic Aspects

S. M. Shafkat Raihan^a, Rudra Pratap Deb Nath^{a,*}, Tonmoy Chandro Das^a, Torben Bach Pedersen^b and Debasish Ghose^c

^a *Department of Computer Science and Engineering, University of Chittagong, Chattogram, Bangladesh*
E-mails: shafkatraihan9001@gmail.com, rudra@cu.ac.bd, tonmoy.csecu@gmail.com

^b *Department of Computer Science, Aalborg University, Denmark*
E-mail: tbp@cs.aau.dk

^c *Kristiania University College, Norway*
E-mail: debasish.ghose@kristiania.no

Abstract. The Covid-19 pandemic is influenced by many environmental, health, and socioeconomic aspects such as air pollution, comorbidity, occupation, etc. Decision makers need better data on the mortality and morbidity of Covid-19 to efficiently withhold its spread. The majority of the data resources dedicated to Covid-19 focus on spatiotemporal aspects only. Furthermore, existing research often overlooks the integrated impact of combining multiple factors. In this study, we efficiently model and analyse Covid-19's epidemiological data from multiple dimensions, such as time, location, temperature, comorbidity, occupation, etc. Data warehousing technology is used to model and integrate data from disparate sources in a multidimensional format. Besides, to make the data interoperable and accessible, they are annotated, integrated, and published semantically using the Resource Description Framework (RDF) model in accordance with the Findability, Accessibility, Interoperability, and Reusability (FAIR) principles. To facilitate Online Analytical Processing (OLAP) compatibility, we annotate the Covid-19 knowledge graph—referred to as CovKG—with multidimensional semantics using QB and QB4OLAP vocabularies. CovKG is analyzed through an interactive analytical interface to observe the Covid-19 confirmed cases and deaths from thirteen aspects. Finally, the performance and quality of CovKG are assessed against prominent data stores modeling Covid-19 data. The ETL workflow typically takes around 42 minutes to load CovKG, which is connected to 10,951 external resources, has a size of about 5.3 GB, and consists of around 44 million RDF triples. When evaluated using competency queries, CovKG can answer 100% of the questions, whereas other prominent data stores can only provide the best answers for 39% of them.

Keywords: Knowledge Graph, Covid-19, Multidimensional Analysis, FAIR principles, Linked Data, Semantic Technology

* Corresponding author. E-mail: rudra@cu.ac.bd.

1. Introduction

The Covid-19 pandemic has affected nearly every country in the world. It is a contagious disease caused by the SARS-COV-2 virus, which exhibits a high mutation and transmission rate. Since the onset of the pandemic, numerous studies, datasets, and systems have been proposed and implemented worldwide to monitor the disease's epidemiology. Epidemiology is concerned with understanding the incidence, distribution, causes, and potential control of diseases in populations. In the case of Covid-19, this can be achieved by observing attributes such as daily confirmed and death cases, recoveries, critical cases, and more. The epidemiology of Covid-19 is influenced by various factors. While most systems and studies conducted on this issue primarily focus on the spatiotemporal aspect, it is important to note that environmental, health, and socioeconomic factors also play a significant role in influencing the spread of Covid-19.

Environmental factors include variables such as air pollution, humidity, temperature, precipitation, and wind speed, which can contribute to virus incubation and influence individual health. Health-related aspects involve factors like comorbidity and vaccine hesitancy. Comorbidity, which refers to the presence of another underlying disease, can impact the immune system and consequently affect the likelihood of contracting Covid-19. The socioeconomic aspects encompass elements such as occupation, ethnicity, urbanization, and so forth. For instance, occupations that entail leaving home, such as medical professions and law enforcement, carry a higher risk of exposure to the disease. Research focusing on these different aspects has often treated them in isolation rather than considering their integrated effects. Furthermore, the hierarchical structure of these factors is frequently overlooked. For example, most data repository systems designed to monitor the spatiotemporal spread of Covid-19 only track information at the country level and do not provide details at finer sub-national levels.

In this study, we utilize Business Intelligence (BI) [1] technologies to construct a robust data framework that empowers users to comprehensively address the previously mentioned problem. This approach sets itself apart from prior studies that focused solely on individual aspects and neglected to consider the hierarchical structure of the factors involved. BI encompasses a collection of disciplines and technological tools that provide intelligent support to decision-makers within organizations, enabling them to make efficient decisions regarding their business processes [2]. Therefore, global organizations like the World Health Organization (WHO) can utilize BI on Covid-19 epidemiological data, including confirmed cases and deaths, to analyze and mitigate the impact of Covid-19 and its transmission.

To achieve this, the data is procured from disparate sources and integrated into a Data Warehouse (DW) through an Extract-Transform-Load (ETL) workflow. In the ETL workflow, data are collected from disparate sources, transformed into an agreed upon format, and loaded in a data store that allows analysis, respectively [3]. Furthermore, to facilitate data analysis from multiple perspectives, the DW is designed in accordance with the Multidimensional (MD) model. This model offers an easily understandable framework in which data are organized within an n-dimensional space, commonly referred to as a data cube. This space is composed of dimensions (representing the cube's axes) and facts (representing the cells within the cube) [4]. Dimensions are ordered into hierarchies (composed of a number of levels) to explore and (dis)aggregate fact measures (i.e., numeric data) at various levels of detail [5]. For example, the `geographyHierarchy` hierarchy (*Admin2* → *Admin1* → *Country* → *Continent*) of the `Geography` dimension allows to (dis)aggregate the number of deaths at various administrative levels of detail.

We model the Covid-19 epidemiological data using fact constellation of data cuboids [6] as it helps managing and modelling in a sustainable manner. It also enables Online Analytical Processing (OLAP) functionality [7]. OLAP functionality provides quick and accurate results. It consists of a number of operations, such as roll up (where data is aggregated to a coarser granularity), drill down (where data is dis-aggregated to a finer granularity), drill out (where data is spread out along multiple cells), drill across (where data in two cubes is merged through one or more shared dimensions), slice (where the value of one dimension level is fixed and the analysis is done along the others), dice (where the value of one or more dimension levels is fixed to one or more levels and the analysis is done along the others) etc. [1] [8].

To contextualize and enable semantic integration on the data, the DW is implemented as a knowledge graph. To do this, the data cuboids are represented using the Resource Description Framework (RDF) [9] model. RDF is the W3C standard web data model designed for flexible data interchange on the web. The resulting knowledge

graph will truly prove to be beneficial if it is published using the Findability, Accessibility, Interoperability, and Reusability (FAIR) principles [10]. An efficient way to ensure this is to publish the data as Linked Data [11]. This will connect the integrated data to the vast ontology preserved in the Linked Open Data (LOD) cloud [12] which houses more than 1300 interlinked datasets. To achieve this, RDF, RDF Schema (RDFS) [9], and Web Ontology Language (OWL) [13] vocabularies are used in the data descriptions to apply various constraints to the data. To annotate data with MD semantics, Data Cube (QB) [14] and QB for OLAP (QB4OLAP) [15] vocabulary are used as well. Preparing the data in RDF format enables it to be queried using the RDF query language, SPARQL. To summarize, the unique contributions of this study are as follows:

- Producing a Covid-19 knowledge graph (CovKG)¹ with MD semantics from diverse sources. This involves collecting data from 20 different sources, defining a target model (a schema-level knowledge graph) with MD semantics by analyzing the source data, and integrating this data into the knowledge graph in a comprehensive and sustainable manner according to the semantics encoded in the target model.
- Linking data internally and externally with other knowledge graphs available in the LOD cloud.
- Developing an OLAP interactive interface to facilitate users creating OLAP queries using GUI components for deriving business critical knowledge.
- Conducting qualitative assessment using SPARQL queries and drawing statistical insights regarding the multiple dimensions of Covid-19 epidemiology.

The remainder of the paper is organized as follows. Section 2 defines various terms that are frequently referred throughout the study. Section 3 discusses the previous related work. Section 4 describes the datasets and methods used in the study to model the knowledge graph. Section 5 describes the knowledge graph generation process. Section 6 describes the features of CovKG. Section 7 describes the experimental evaluation. Finally, Section 8 provides the concluding remarks and suggestions for future work.

2. Preliminaries

In this section, we introduce the relevant terms that appear frequently throughout the paper.

2.1. Knowledge Graph

A knowledge graph (KG) is a semantic graph that manifests as interlinked network of real-world entities and visualizes the relationship between them. The KG comprises two elements: Terminology Box (TBox) and Assertion Box (ABox). The TBox defines the domain schema, while the ABox represents instances [16]. Formally, the TBox is defined as a 3-tuple: $TBox = (C, P, A^O)$, where C , P , and A^O represent the sets of concepts, properties, and terminological axioms. A concept is the blueprint of a group of instances sharing common properties. Properties establish relationships between instances of concepts or link instances of a concept to literals. Terminological axioms describe concepts, properties, and the interconnections and restrictions among them within the domain. The ABox assertions must conform to the definitions set by the TBox. In our context, the schema and instances of source datasets are called source TBoxes and ABoxes respectively. The TBox of the KG is formed by integrating and modelling the source data is the target TBox. It can have one or multiple target ABoxes. We refer to the KG consisting of the target TBox and ABoxes as the Covid-19 KG (CovKG).

In this paper, we express KG elements using the RDF model [9]. In RDF, real world entities are uniquely represented using internationalized resource identifiers (IRIs), and the description of the entities is expressed in the form of RDF triples which are three-part statements containing a subject, a predicate, and an object. For instance, in Figure 1, the subject `cdw:SpatioTemporalDataset` has an object `cdw:Admin1`. The relation between them is expressed through the predicate `cdw:hasAdmin1`. To express richer constraints on the KG, formal languages like RDFS and OWL are used in combination with RDF. For example, Figure 1 shows that in CovKG, `cdw:Admin1` is annotated as an `owl:Class` and is externally linked to Geonames KG using the `owl:sameAs` property. Given

¹Here, we use the terms “semantic data warehouse” and “knowledge graph” interchangeably.

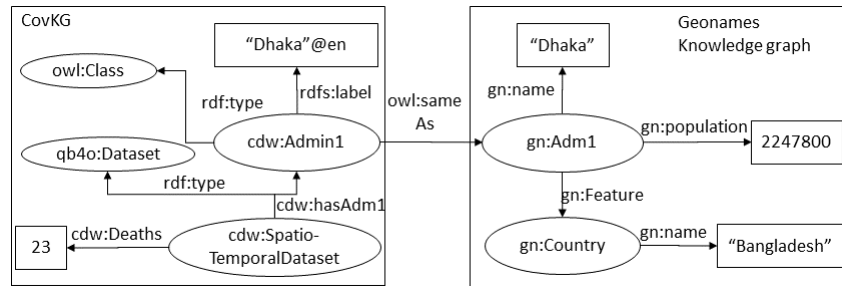


Fig. 1. A visual example of knowledge graphs. Here, `cdw:`, `qb4o:`, `gn:` `rdf:`, `owl:` represent their respective namespaces.

our emphasis on MD modeling, data needs to be annotated with MD semantics at both the schematic and instance levels. For this purpose, we employ the QB4OLAP vocabulary.

2.2. QB4OLAP

QB4OLAP is the extension of the RDF data cube vocabulary [14], which is the W3C standard for publishing statistical data as RDF. Despite being specialized for data cubes, QB does not facilitate OLAP queries to be conducted on the RDF data cubes. Thus QB4OLAP was designed as a vocabulary to represent OLAP cubes in RDF. It enables the implementation of OLAP operators as SPARQL queries directly on the RDF representation. In Figure 2, the schematic diagram of the QB4OLAP vocabulary is given.

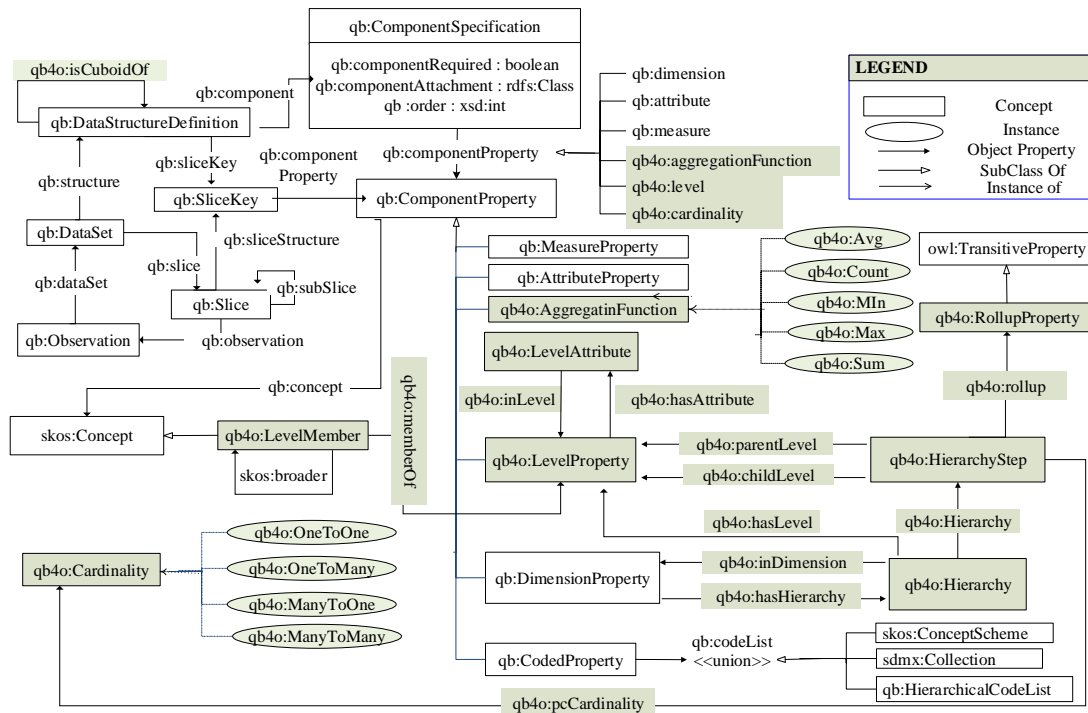


Fig. 2. The QB4OLAP vocabulary (reproduced from [6]).

In the figure, the prefix `qb:` represents the terms defined in the QB vocabulary, QB4OLAP terms are denoted with the prefix `qb4o:` and displayed with a gray background. RDF classes, properties, and class instances are

represented by capitalized terms, noncapitalized terms, and capitalized italic terms respectively. An arrow from class A to class B, with the label *rel* points out that *rel* is an RDF property of domain A and range B. White triangle arrows represent subclass, or subproperty relationships whereas black diamonds represent `rdf:type` relationships (instances). It can be seen that the vocabulary presents various constructs to represent the dimension (`qb:DimensionProperty`), its levels (`qb4o:LevelProperty`), level members (`qb4o:LevelMember`), level attributes (`qb4o:LevelAttributes`), level's rollup properties (`qb4o:RollupProperty`), dimension hierarchies (`qb4o:Hierarchy`), and hierarchy steps (`qb4o:HierarchyStep`). QB4OLAP offers the `qb:Dataset` to delineate an observation dataset. The structure of this dataset is outlined using `qb:DataSetStructureDefinition`. This structure can take the form of either a cube, if specified in dimensions and measures, or a cuboid, if outlined in levels of dimensions and measures. In Section 5, we illustrate the utilization of various QB4OLAP components for annotating CovKG with MD semantics, both in the TBox (Listing 1) and ABox (Listing 4 and 5).

3. Related Work

In this section, we conduct a comprehensive examination of prior relevant research related to the current study's topic. Through the analysis of prior studies, datasets, and systems focused on Covid-19, we categorize them into two specific groups:

1. Relevant research papers within the domain of our interest.
2. Prominent public data repositories dedicated to reporting on Covid-19.

Table 1 presents a summary of the comparative analysis of research papers and prominent data repositories regarding various features associated with COVID-19 data. The table lists the sources of previous research or repositories, indicating their involvement with confirmed and death cases, the core technologies utilized, usage of knowledge graphs, adherence to FAIR principles, compatibility with external datasets, provision of query interfaces or dashboards, availability of downloadable data, ability to conduct visual data analysis, covered aspects, and consideration of multiple dimensions. Core technologies serve as indicators of whether: 1) their processes involve either discovery techniques or surveys, 2) they employ data mining or pattern mining to uncover hidden knowledge, 3) they utilize DW/OLAP technology for descriptive analysis, 4) RDF technology is used for semantic annotation, and 5) Natural Language Processing (NLP) is employed for processing scientific open data.

The table reveals that the majority of the research papers focus on analyzing confirmed cases, with only [21] considering both death and confirmed counts. Most of these studies collect data from secondary sources and utilize data or pattern mining techniques to unveil hidden insights. However, [20] and [21] employ Data Warehousing technology to enable OLAP-like analysis. [22] retrieves data from scientific literature using NLP techniques and apply KG methods to analyze drug-drug interactions. Similarly, [23] also utilizes KGs, covering various aspects, although they do not enable MD analysis. Both [22] and [23] adhere to FAIR principles and provide downloadable data. While most of the research papers offer query interfaces, [18] and [19] do not. However, none of the studies provide a dashboard to facilitate end-users for data analysis.

Since the beginning of the Covid-19 pandemic, various government and non-government organizations worldwide have made considerable efforts to make real-time Covid-19 data available to the public. These steps have provided researchers with abundant data to conduct essential research necessary to tackle this global catastrophe. Moreover, they have made the data publicly available through online platforms, allowing users to search, view, analyze, and download data related to Covid-19. Articles [24] - [29] are some of the most prominent data repositories dedicated to monitoring Covid-19 epidemiology. Some of them focus on global Covid-19 scenario, such as Worldometer [24], while there are those which target a specific region, such as Bangladesh Dynamic Dashboard for Covid-19 [27]. Note that the study in [26] records number of Covid-19 waves instead of death and confirmed counts.

The studies and repositories mentioned are undoubtedly valuable, providing prolific data for analyzing Covid-19. However, as outlined in Table 1 they are not devoid of shortcomings. For instance, studies [17]- [20] utilized individual-level data. Although individual-level data offers a large number of influencing parameters, such parameters often have a significant amount of missing data due to undisclosed personal information. In the case of modeling

Table 1:
Overview of related works.

Category	Reference	Death info	Confirmed cases info	Core Technologies					KG?	FAIR?	Compatible with external dataset?	Query interface/ Dash board	Enable visual Data Analysis?	Data Downlo- adable?	Covered aspects	Covered multiple dimen- sions?
				Data collection	Pattern/ data mining	DW / OLAP	RDF	NLP								
Research Papers	[17]	X	✓	DD	X	X	X	X	X	X	✓	X	X	SP	X	
	[18]	X	✓	DD	X	X	X	X	X	X	X	X	X	T	X	
	[19]	X	✓	DD	✓	X	X	X	X	X	X	X	X	SE	X	
	[20]	X	✓	DD	✓	✓	X	X	X	X	✓	X	X	SE	X	
	[21]	✓	✓	DD	X	✓	X	X	X	X	✓	X	X	E	✓	
	[22]	X	✓	DD	✓	X	X	✓	✓	✓	✓	X	X	H	X	
	[23]	X	✓	DD	X	X	X	✓	✓	✓	✓	X	✓	SP, SE, H	X	
	[24]	✓	✓	S	X	X	X	X	X	X	✓	✓	✓	SP	X	
	[25]	✓	✓	S	X	X	X	X	X	X	✓	✓	✓	SP	X	
Prominent Data Repositories	[26]	X	✓	S	X	X	X	X	X	X	✓	✓	✓	SP	X	
	[27]	✓	✓	S	X	✓	X	X	X	X	✓	✓	✓	SP	X	
	[28]	✓	✓	S	X	X	X	X	X	X	✓	✓	✓	SE, SP	X	
	[29]	✓	✓	S	X	X	X	X	X	X	✓	✓	✓	SP	X	
CovKG		✓	✓	DD	✓	✓	✓	✓	✓	✓	✓	✓	SP, T,F, H, SE	✓		

*[DD-Data Discovery, S-Survey] † [SP-Spatial, T-Temporal, E-Environment, H-Health, SE-Socioeconomic]

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

1 occupation data, they did not employ any recognized occupation classification model. In contrast, our study aggregated population-level data—specifically, the number of confirmed cases and deaths—rather than individual cases. 2 Such data are widely available and have the least amount of missing values. We also included occupation as one of 3 the dimensions for analyzing Covid-19 and followed the ISCO-08 system to model it, which is an internationally 4 renowned occupation classification system [30]. 5

6 In [31], a theoretical framework for modeling COVID-19 data was presented, but its implementation was not 7 realized. The authors in [23] proposed an existing ontology using Wikidata to serve as a knowledge base for COVID- 8 19. However, it is general-purpose ontology where Covid-19 information is just one facet. Our study designed 9 and implemented a novel ontology (TBox of CovKG) and CovKG, dedicated and specialized for Covid-19 data 10 modeling. 11

12 The data warehouse modeled in [21], named COVID-WAREHOUSE, specifically focuses on data from Italy. 13 Our ontology covers twenty-two countries, with Italy being one of them. It also analyzes Covid-19 from thirteen 14 perspectives, including weather and air pollution aspects. Furthermore, authors in [22] contributed to the study of 15 relation between comorbidities and Covid-19, but they did not consider structured data. Our ontology is based on 16 a general-purpose, structured dataset of confirmed cases and deaths, providing the flexibility to conduct interdis- 17 ciplinary research with ease. The prominent data stores share the common limitation that data is not available in 18 a semantic format. Our CovKG is a semantic data warehouse created using the RDF model. Users can link its 19 components to external KGs. This makes data sharing and new knowledge discovery relatively easier. 20

21 In summary, our study addresses these research gaps through the application of multidimensional semantic data 22 integration, forming, and analyzing CovKG from multiple perspectives at fine granularities, and aims to integrate it 23 as part of linked open data. 24

25 4. Methodology and Knowledge Graph Modelling 26

27 The entire methodology to generate CovKG by modeling Covid-19 epidemiological data in a multidimensional 28 format is illustrated in Figure 3. Initially, the data is collected from various data sources, which are obtained from the 29 Web. Most of the sources present their data in CSV, XLS(X), and JSON formats. After collecting the raw data, we 30 semantically design the TBox of CovKG (data warehouse) with MD semantics using a demand-driven approach. 31 Then, the CovKG is populated using an ETL pipeline. In the ETL pipeline, the relevant data is extracted from 32 the sources, transformed semantically according to the semantics encoded in the target schema, and finally, the 33 transformed data is loaded into a data warehouse in the form of a semantic graph. Using an OLAP interface, the 34 OLAP operability of the CovKG is assessed. Finally, using SPARQL queries, qualitative assessment and statistical 35 analysis are performed. 36

37 In this section, we outline data sources and design the target TBox for CovKG. In the next section, we explore 38 the generation of CovKG'. 39

40 4.1. Description of Data Sources 41

42 We use different searching techniques ([32], [33], [34]) for discovering Covid-19 related datasets [35]. Dur- 43 ing the search process, we highlight the aspects of data and selectively focus on datasets that capture information 44 from multiple dimensions. For instance, ‘Covid-19 daily province level confirmed case dataset by occupation’ is an 45 example of a self-explanatory query. We extract the data either directly from the source sites or by utilizing Applica- 46 tion Programming Interfaces (APIs). The datasets are available as either CSV, Spreadsheets, or JSON files. The data 47 mainly consist of Covid-19 confirmed cases and death counts. The data are collected with a focus on spatiotemporal, 48 environmental, health, and socioeconomic aspects. An overview of the data sources is provided in Table 2. 49

50 **Spatiotemporal dataset:** Spatiotemporal data are collected from 11 sources spanning over 18 countries. Addi- 51 tionally, Turkey and Romania’s data obtained for occupation dataset are also used. Daily data for Austria, Italy, 52 Lithuania, Poland, Slovakia, and Sweden are available at state, province, county, voivodeship, region, and country 53 level respectively. Croatia, Denmark and Finland’s data are available at country, municipality and region level re- 54 spectively. Greece, Ireland and Netherland’s data are available at prefecture, county and municipality level. Indian 55

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51Table 2:
Overview of the data sources.

Aspect	Dataset	Sources	Covered countries	# of instances
Spatiotemporal	Spatiotemporal	Humanitarian data exchange portal [36], [37]	Afghanistan	2,316,677
		GitHub https://shorturl.at/LQT78 , https://covid19.go.id/	Indonesia	
		Dynamic Covid-19 dashboard for Bangladesh [27]	Bangladesh	
		European centre for disease prevention and control [38]	Austria, Italy, Lithuania, Poland, Slovakia, Sweden	
		Naqvi et al. [39]	Croatia, Denmark, Finland	
		GitHub [40]	Greece	
		Geohive open data repository [41]	Ireland	
		NL COVID-19 geohub repository [42]	Netherlands	
		Kaggle [43]	India	
		South Africa provincial breakdown dashboard [44]	South Africa	
Environment	Weather	Haratian et al. [45]	United States	1,337,073
		World weather online [46]	Countries listed under the Spatiotemporal aspect plus Turkey and Romania	
Health	Air pollution	South Africa air quality information system (SAAQIS) [47], Central control room for air quality management (CPCBCCR) [48]	South Africa, India	157,662
		Vaccine hesitancy	Denmark, Germany, Netherlands	
		Comorbidity	United States	
		Ethnicity	United States	
Socioeconomic	Place of death	Centers for disease control and prevention (CDC) [50]	United States	8,394
		Centers for disease control and prevention (CDC) [50]	United States	
		Sari et al. [51], Windsor-Sheldard and Rabiya [52], Hamecan et al. [53]	United Kingdom, Romania, Turkey	
Socioeconomic	Occupations	Haratian et al. [45]	United States	1,572,241
		Urbanicity	United States	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

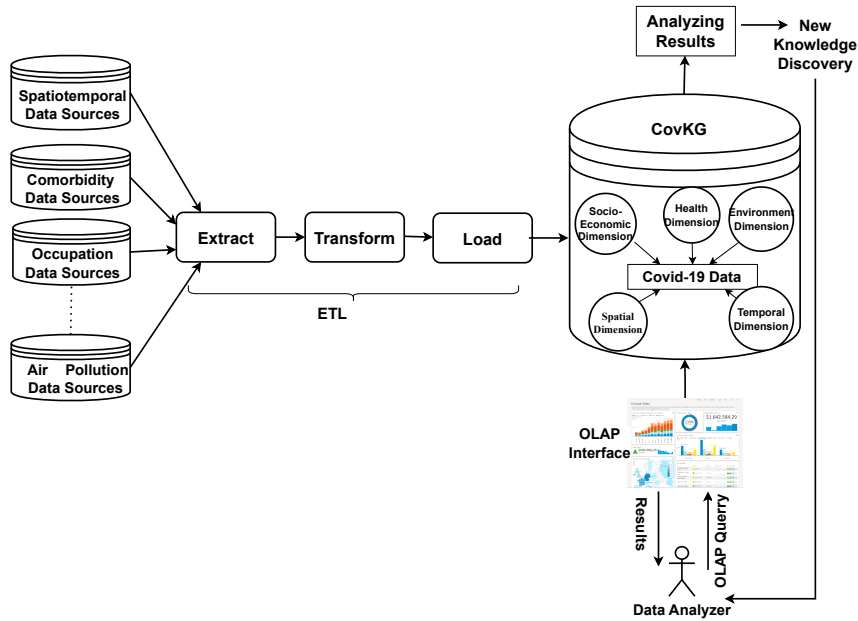


Fig. 3. Incorporation of the ETL workflow to form CovKG from disparate data sources, followed by conducting analysis on CovKG.

data are available at states and unions level whereas South Africa's data are at provincial level. Finally, country level data for USA are collected.

Weather dataset: Historic weather data for various geographic locations specified by latitude and longitude are available. The collected data contains information such as temperature (degrees Celsius), humidity (%), precipitation (millimeters), and wind speed (kilometers per hour).

Air pollution dataset: We collect the daily state-level air pollution data for South Africa and India. The pollutants considered are ground-level Ozone (O_3), particulates ($PM_{2.5}$ and PM_{10}), Sulfur dioxide (SO_2), Carbon monoxide (CO), and Nitrogen dioxide (NO_2). These pollutants are commonly measured in monitoring air pollution, as per [54].

Vaccine hesitancy dataset: Questionnaire data are collected, which document and classify hesitant behavior of the subjects regarding vaccination.

Comorbidity dataset: Monthly death counts at the state level, grouped by comorbidity in the U.S.A. are available.

Ethnicity dataset: USA death counts of various races at the state level on a monthly basis are available.

Place of death dataset: Monthly death counts at the state level based on place of death in the U.S.A. are available.

Occupation dataset: Occupation datasets are collected in light of ISCO-08 standard for occupation classification. Data on deaths among various occupations due to Covid-19 in England and Wales are collected. Both Romania and Turkey's data are available at the provincial level.

Urbanicity Dataset: Urbanicity indicates how rural or urban a region is. The urbanicity dataset is constructed from the spatiotemporal data collected for the USA by mapping counties to levels of urbanicity using the urban-rural classification scheme for US counties [55].

4.2. Modeling the Target TBox of CovKG

In this section, we design the MD schema of CovKG to integrate the sources defined in Section 4.1, as depicted in Figure 4. Since we integrate data from various dimensions and disparate sources, not all data points will align, after integration. For instance, count of deaths of a specific ethnicity at a certain place on a certain time may be available from one data source. Same may be available in case of deaths of patients of a certain comorbidity from another data source. However, to place them in a single cube, we need to know the overlap between the counts i.e.,

the exact count of deaths in that ethnicity who had that specific comorbidity. This is what we mean, when we say data points will not align. Furthermore, data concerning the same dimensions collected from different sources may be available at different hierarchical levels. To address this issue, we employ data cuboids [6]. These cuboids are subsets of data cubes and are represented with respect to one or more dimension levels, in contrast to data cubes, which are represented with respect to all dimensions. This approach allows data cuboids to facilitate separation of concerns when analyzing the data.

In Figure 4, the green cube shapes represent the data cuboids, while the blue rectangles represent the dimensions. An exclusive relationship (\otimes) between two levels indicates that the cuboid can contain data from either of the levels. Measures of the cuboids are *Total Confirmed Cases* and *Total Deaths*. Among the 13 dimensions, `cdw:GeographyDim` and `cdw:TimeDim` constitute the spatiotemporal perspective; `cdw:TemperatureDim`, `cdw:HumidityDim`, `cdw:WindDim`, `cdw:PrecipitationDim`, and, `cdw:AirPollutionDim` constitute the environmental aspect; `cdw:VaccineHesitencyDim` and `cdw:ComorbidityDim` constitute the health perspective; `cdw:EthnicityDim`, `cdw:PlaceofDeathDim`, `cdw:OccupationDim`, and `cdw:UrbanicityDim` constitute the socioeconomic perspective.

The `cdw:TimeDim` dimension has the `cdw:Calendar` hierarchy, which consists of `cdw:Day`, `cdw:Month` and `cdw:Year` levels, from the finest to coarsest level. This hierarchy allows the user to see the temporal evolution of the confirmed cases as well as deaths. The `cdw:GeographyDim` dimension has the `cdw:Geography` hierarchy containing `cdw:Admin2` as the finest level, which represents the second administrative level as per Geonames [56]. The other higher levels are `cdw:Admin1`, `cdw:Country`, and `cdw:Continent`.

The `cdw:TemperatureDim` contains the hierarchy `cdw:Temperature`. It houses two levels (`cdw:ThermanlSubtype`, `cdw:Thermanltype`) based on the classification model of [57]. Each of the `cdw:HumidityDim`, `cdw:WindDim` and `cdw:PrecipitationDim` dimensions contains only one level. The `cdw:AirPollutionDim` dimension represents the daily air pollution through various pollutants. This dimension has the `cdw:AirPollution` hierarchy, which contains the `cdw:PollutionLevels` and `cdw:Pollutants` levels. The `cdw:PollutionLevels` level represents the various levels of pollution determined according to the levels depicted in [54]. Since Covid-19 is a respiratory disease, it is intuitive that air pollution may have relationships with its epidemiological behavior [58–61].

The dimension `cdw:VaccineHesitencyDim` represents the hesitancy to accept vaccines. Its `cdw:Hesitancies` hierarchy has two levels: `cdw:HesitancyScore` and `cdw:VaccineAvailabilityYear`. It is based on the questionnaire used in the research done in [49]. The vaccination intent was determined based on the survey respondents' responses to the question of whether they will get vaccinated if a Covid-19 vaccine becomes available to them in 2021. The respondent can respond with scores 1 (Strongly Agree), 2 (Agree), 3 (Neutral), 4 (Disagree), and 5 (Strong Disagree). The level `cdw:VaccineAvailabilityYear` groups the responses based on year. We have selected vaccine hesitancy as a dimension as it represents people's tendency not to take vaccine, hence assist in the further propagation of the pandemic. In 2019, WHO listed vaccine hesitancy as one of the top ten threats to global health [62].

The dimension `cdw:ComorbidityDim` represents affliction with other diseases alongside Covid-19. Under its `cdw:Diseases` hierarchy, there are two levels: `cdw:Disease` and `cdw:DiseaseType`. Diseases are categorized into three disease types: respiratory diseases, circulatory diseases, and other diseases. Various research has shown that the presence of comorbidities such as diabetes, respiratory diseases, cardiovascular diseases etc., can influence Covid-19's impact on the patient's body [63–65]. The `cdw:EthnicityDim` dimension represents the races of the cases in the hierarchy `cdw:Ethnicities`. This hierarchy holds two levels - `cdw:Race` and `cdw:RaceType`. The level `cdw:RaceType` consists of two instances - Hispanic and Non-Hispanic. Although further classification for Hispanic is not available, Non-Hispanic is divided into five groups, namely, white, black, American Indian / Alaskan native, Asian, and native Hawaiian / Pacific Islander. Multiple researches have shown the existence of ethnic disparity in Covid-19 cases [66–70]. Place of death is an important indicator of the availability of healthcare facilities in a geographical area of interest. The dimension `cdw:PlaceOfDeathDim` indicates the place where the Covid-19 patients died. Under the `cdw:PlaceofDeath` hierarchy, it has two levels: `cdw:DeathPlaceSubtype` and `cdw:DeathPlaceType`. `cdw:DeathPlaceType` has four types: Healthcare, Homelike, Other, and Unknown. The finer level `cdw:DeathPlaceSubtype` has three instances both for Healthcare and Homelike, while Other and *Unknown* are not further divided. The three instances in

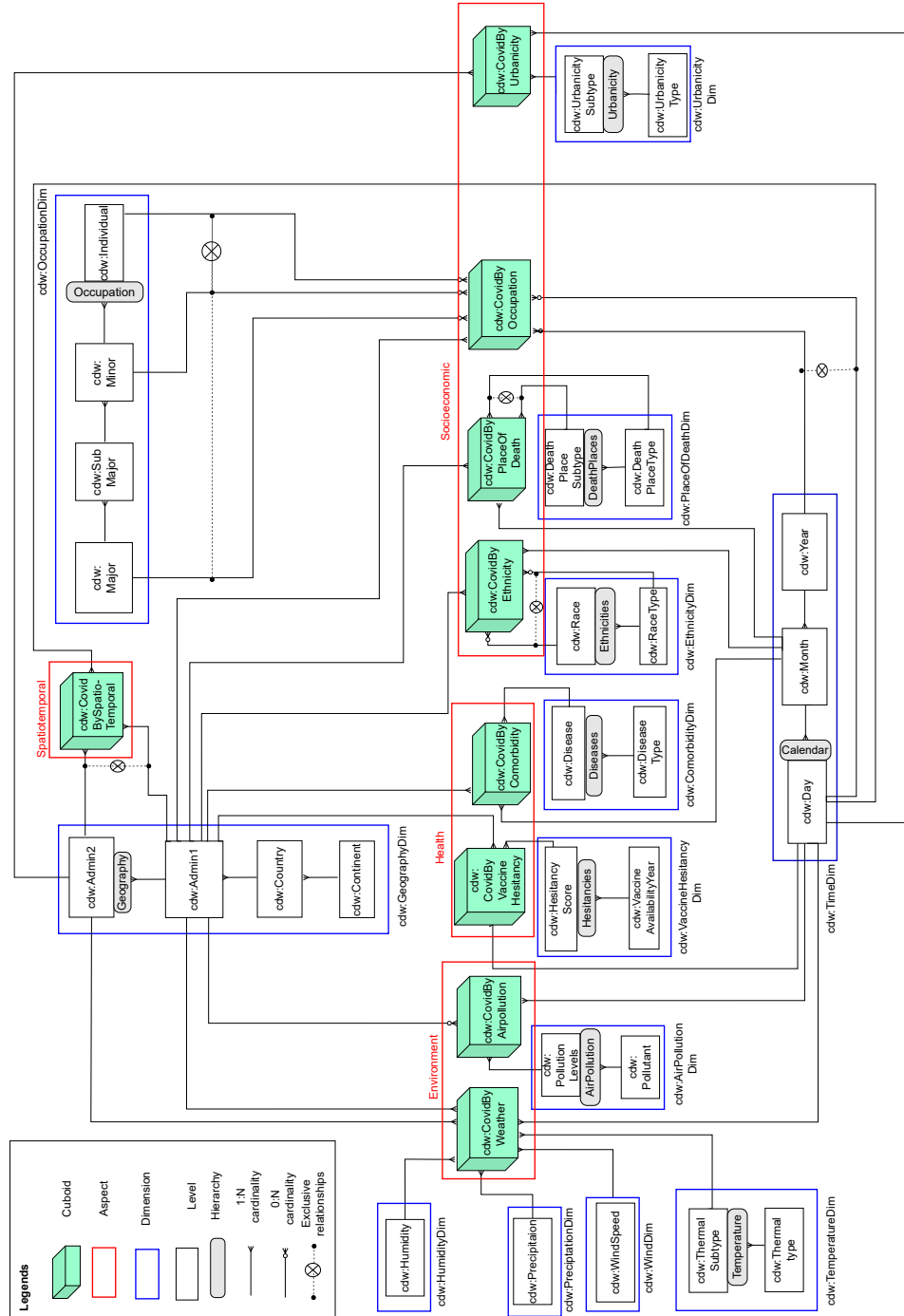


Fig. 4. Multidimensional schema of Covid data warehouse.

cdw:DeathPlaceSubtype for Healthcare are: Inpatient, Outpatient, and Dead on arrival. The three instances of Homelike are: Decedent’s home, Nursing home, and Hospice facility.

Occupations are represented by the cdw:OccupationDim dimension, whose cdw:Occupation hierarchy has four levels as per ISCO-08 [71]. Among the four levels, cdw:Individual represents individual occupations and cdw:Minor represents the minor categories which group the individual professions. Similarly, sub-major occupations categorize minor ones, which in turn are aggregated into major occupations. Occupation is an important socioeconomic factor to take into consideration when analyzing Covid-19’s epidemiology because occupational distribution can describe people’s interaction behavior and hence provide a precursor to the transmission path of the disease.

The dimension cdw:UrbanicityDim represents the urban-rural classification of a geographical area as per the urban-rural classification scheme for counties [72]. The cdw:UrbanicitySubtype level is composed of six instances: Large central metro, large fringe metro, medium metro, small metro, micropolitan, and noncore. Among them, the first four belong to the metropolitan (urban) category and the last two belong to the noncore (rural) category. Noncores and metropolitans constitute the level cdw:UrbanicityType.

5. Generation of CovKG

In this section, we will elaborate on how the task of generating CovKG is accomplished using an ETL pipeline.

In the *Extraction phase*, we extract data from different sources as described in Section 4.1 and cleanse and format those data to conform with the target TBox. This cleansing and formatting tasks include extraction of micro data, elimination of aggregated data, conversion from other formats to CSV, adding unique ids, filtering out irrelevant and noisy data, and so on. Then, the *Transformation phase* semantically transforms the extracted data according to the semantics encoded in the target TBox and generates CovKG, which is in turn loaded into the triple store Openlink Virtuoso [73] in the *Load phase*.

The *Transformation process* unfolds in four steps: 1) The *Target TBox generation* step implements the target TBox defined in Section 4.2. 2) The *Source TBox generation* step generates TBoxes from the data sources. 3) The *SourceToTarget mappings generation* step establishes mappings between the source and target TBoxes. 4) Then, the *Target ABox generation* process generates assertions consistent with the target TBox. Below, we provide a detailed description of each of these steps.

5.1. Target TBox generation

The purpose of this step is to represent the target TBox as outlined in Section 4.2 using the constructs provided by RDFS, OWL, and QB4OLAP (as defined in Section 2) alongside the RDF model. Users have the flexibility to create a TBox with MD semantics either manually or by utilizing tools such as Protege [74], WebVOWL [75], or PoolParty [76]. Listing 1 demonstrates a part of the target TBox annotated with QB4OLAP constructs. In this listing, cdw: represents the namespace of the data cube, which is <https://bike-csecu.com/datasets/covid/cdw#>. The cubic structure of the dataset cdw:SpatioTemporalDataset is defined by cdw:SpatioTemporalCuboid (lines 5-12), which contains the cdw:Admin2 and cdw:Admin1 levels of the cdw:GeographyDim dimension, and cdw:Day of the cdw:TimeDim dimension. Both cdw:Admin2 and cdw:Admin1 are included in this cuboid, as some countries have data available only at the first administrative level, while others have data at the second administrative level.

The cdw:TimeDim dimension contains the cdw:calendarHierarchy hierarchy which is composed of cdw:Day, cdw:Month, and cdw:Year levels (lines 14-21). A hierarchy step of cdw:calendarHierarchy (lines 23-26) represents that cdw:Day is related to its parent level cdw:Month through the rollup property defined by cdw:inMonth (lines 23-26). The level cdw:Day is defined as an instance of the qb4o:LevelProperty class, and it contains a set of attributes (lines 30-32). cdw:dayID is a level attribute (lines 35-37) and cdw:inMonth is defined as a roll-up relation (lines 39-41). The measures cdw:Confirmed and cdw:Deaths are defined using the qb:MeasureProperty class. They are defined as decimals to facilitate aggregation functions like average, which may return floating point values.

```

1 1 #DATASETS
2 2 cdw:SpatioTemporalDataset a qb:DataSet;
3 3   qb:structure cdw:SpatioTemporalCuboid.
4 4 #CUBOIDS
5 5 cdw:SpatioTemporalCuboid a qb:DataStructureDefinition;
6 6   dct:conformsTo <http://purl.org/qb4olap/cubes>;
7 7   qb4o:isCuboidOf cdw:COVID_DW;
8 8   qb:component [ qb:measure cdw:Confirmed; qb4o:aggregateFunction qb4o:sum];
9 9   qb:component [ qb:measure cdw:Deaths; qb4o:aggregateFunction qb4o:sum];
10 10  qb:component [ qb4o:level cdw:Admin1; qb4o:cardinality qb4o:OneToMany];
11 11  qb:component [ qb4o:level cdw:Admin2; qb4o:cardinality qb4o:OneToMany];
12 12  qb:component [ qb4o:level cdw:Day; qb4o:cardinality qb4o:OneToMany].
13 13 #DIMENSIONS
14 14 cdw:TimeDim a qb:DimensionProperty;
15 15   rdfs:label "Time Dimension"@en;
16 16   qb4o:hasHierarchy cdw:CalendarHierarchy.
17 17 #HIERARCHIES
18 18 cdw:CalendarHierarchy a qb4o:Hierarchy;
19 19   rdfs:label "Calendar Hierarchy"@en;
20 20   qb4o:inDimension cdw:TimeDim;
21 21   qb4o:hasLevel cdw:Day, cdw:Month, cdw:Year.
22 22 #HIERARCHY STEPS
23 23 _:hsl6 a qb4o:HierarchyStep;
24 24   qb4o:inHierarchy cdw:CalendarHierarchy;
25 25   qb4o:childLevel cdw:Day;
26 26   qb4o:parentLevel cdw:Month;
27 27   qb4o:pcCardinality qb4o:OneToMany;
28 28   qb4o:rollup cdw:inMonth.
29 29 #LEVELS
30 30 cdw:Day a qb4o:LevelProperty;
31 31   rdfs:label "Day"@en;
32 32   qb4o:hasAttribute cdw:dayID, cdw:dayName, cdw:inMonth;
33 33   rdfs:range cdw:Day.
34 34 #ATTRIBUTES
35 35 cdw:dayID a qb4o:LevelAttribute;
36 36   rdfs:label "Day ID"@en;
37 37   rdfs:range xsd:string.
38 38 #ROLLUP RELATIONSHIPS
39 39 cdw:inMonth a qb4o:LevelAttribute, qb4o:RollupProperty;
40 40   rdfs:label "Rollup property to roll up from day to month"@en;
41 41   rdfs:range onto:Month.
42 42 #MEASURES
43 43 cdw:Confirmed a qb:MeasureProperty;
44 44   rdfs:label "Total Confirmed Cases"@en;
45 45   rdfs:range xsd:decimal.
46 46 cdw:Deaths a qb:MeasureProperty;
47 47   rdfs:label "Total Deaths"@en;
48 48   rdfs:range xsd:decimal.

```

Listing 1: Target TBox defining spatiotemporal cuboid and the time dimension. Prefixes are omitted due to space constraints.

5.2. Source TBox generation

After creating the target TBox of CovKG, the next task is to populate CovKG from the available data sources. In order to accomplish this, we must establish mappings between the target and source constructs at the TBox level. Therefore, it is essential to derive TBoxes from the existing sources and augment them with OWL and RDFS

constructs. Various vocabularies/tools, such as R2RML mapping [77], Direct mapping [78], and NonSemanticToTBoxDeriver [3], can be used for this purpose. In this study, we utilize NonSemanticToTBoxDeriver. This tool is employed to extract conceptual information embedded in non-semantic structured data and transform it into a semantic form. Listing 2 displays the generated source TBox from the spatiotemporal dataset. Here, the table name is used as an OWL class and the attributes are considered as OWL datatype properties. The `onto:` namespace is used to create semantic counterparts of the various elements of the source data.

```

1 onto:SpatiotemporalFact a owl:Class.
2 onto:adm2ID a owl:DatatypeProperty;
3   rdfs:domain onto:SpatiotemporalFact;
4   rdfs:range xsd:string.
5 onto:Death a owl:DatatypeProperty;
6   rdfs:domain onto:SpatiotemporalFact;
7   rdfs:range xsd:decimal.
8 onto:adm1ID a owl:DatatypeProperty;
9   rdfs:domain onto:SpatiotemporalFact;
10  rdfs:range xsd:string .
11 onto:Confirmed a owl:DatatypeProperty;
12  rdfs:domain onto:SpatiotemporalFact;
13  rdfs:range xsd:decimal.
14 onto:dayID a owl:DatatypeProperty;
15  rdfs:domain onto:SpatiotemporalFact;
16  rdfs:range xsd:string.

```

Listing 2: Source TBox for the fact table SpatiotemporalFact.

5.3. Source-to-Target Mappings generation

Source data can be heterogeneous in nature. To handle this situation, sources should be mapped to the target at the TBox level. The communication between the source and target is materialized in the form of intermediate mapping definitions that assist complex data flows between the sources and target. As source TBoxes are generated for both dimension tables and fact tables, the mappings are to be generated for both as well. Source-to-target mappings of dimension tables are used later to produce the semantic assertions of the level instances. Similarly, source-to-target mappings of fact tables are used later to create semantic assertions of the fact tables. Listing 3 shows mapping definitions between different constructs of `onto:SpatiotemporalFact` and `cdw:spatioTemporalDataset`. The mapping file are annotated with Source-to-Target Mapping (S2TMAP) vocabulary [79]: an OWL-based mapping vocabulary.

```

1 #Dataset mapping
2 cdw:spatiotemporalfacts_COVID_Schema a map:Dataset ;
3   map:source  '/SpatiotemporalFacts';# source location
4   map:target  '/COVID_Schema'.# target location
5 #Concept mapping
6 cdw:SptempFact_SpTempDataset a map:ConceptMapper;
7   map:sourceConcept  onto:SpatiotemporalFact;
8   map:targetConcept  cdw:spatioTemporalDataset.
9   map:dataset  cdw:spatiotemporalfacts_COVID_Schema;
10  map:iriValue  "CONCAT(onto:adm1ID,CONCAT(_,CONCAT(onto:adm2ID,CONCAT(_,onto:dayID))))";
11  map:iriValueType  map:Expression;
12  map:matchedInstances  "All";
13  map:relation  skos:exact;
14 #Property mapping
15 map:PropMap_Confirmed_Confirmed a map:PropertyMapper ;
16   map:ConceptMapper  cdw:SptempFact_SpTempDataset;
17   map:sourceProperty  onto:Confirmed;
18   map:sourcePropertyType  map:SourceProperty;

```



```

1 19     map:targetProperty      cdw:Confirmed.
2 20 map:PropertyMapper_01_dayID_day a  map:PropertyMapper;
3 21     map:ConceptMapper      cdw:SptempFact_SpTempDataset;
4 22     map:sourceProperty     onto:dayID;
5 23     map:sourcePropertyType map:SourceProperty;
6 24     map:targetProperty     cdw:Day.

```

Listing 3: Source-to-Target mapping file of `onto:SpatioTemporalFact` and `cdw:spatioTemporalDataset`.

In S2TMAP, a property-level mapping is nested within a concept-level mapping, which is further encapsulated within a mapping dataset. A mapping dataset is defined as an instance of `map:Dataset`, which captures the references of the source and target TBoxes (lines 2-4). A concept-mapping outlines the correspondence between a source and a target concept (lines 6-13). The source and target concepts are defined using the `map:sourceConcept` and `map:targetConcept` properties. The linkage between a concept-mapping and its mapping dataset is established through the `map:dataset` property. The properties `map:iriValue` and `map:iriValueType` signify that values of `onto:admID`, `onto:adm2ID`, and `onto:dayID` are concatenated to generate unique IRIs for the observations of `cdw:SpatioTemporalDataset`. The "All" value of `map:matchedInstances` indicates that all source instances are mapped.

A source and target property are mapped using the `map:PropertyMapper` (lines 15-19). The connection between a property-mapping and its corresponding concept-mapping is established via `map:conceptMapper`. The specification of the target property within the property-mapping is done using `map:targetProperty`, and this target property can be associated with either a source property or an expression. In this specific instance, the target property `cdw:Confirmed` is mapped to the source property `onto:Confirmed`.

5.4. Target ABox generation

Using the target TBox, along with the source datasets (extracted and cleansed ones) and source-to-target mapping definitions as inputs, this *Target ABox generation* process generates the target ABoxes from the source datasets based on the semantics specified in the target TBox. In QB4OLAP, dimensional data is physically stored in levels, where each level member is identified by a unique IRI and is semantically linked with its relevant level attributes and roll-up properties.

```

34 1 day:1 a qb4o:LevelMember;
35 2   cdw:dayID "1";
36 3   cdw:dayName "2020-01-01";
37 4   cdw:inMonth month:1;
38 5   qb4o:memberOf cdw:Day;
39 6   owl:sameAs wd:Q57396575.

```

Listing 4: An level member of the `cdw:Day` level in the target ABox.

```

44 1 stempd:_1224_52 a qb:Observation;
45 2   cdw:Confirmed "8096";
46 3   cdw:Deaths "270";
47 4   cdw:Admin1 _:b;
48 5   cdw:Admin2 adm2:1224;
49 6   cdw:Day day:52;
50 7   qb:dataset cdw:SpatioTemporalDataset.

```

Listing 5: An observation of the `cdw:SpatioTemporalDataset` dataset in the target ABox.

Table 3

A quantitative overview of the dimensions present in CovKG.

Aspect	Dimension	# of instances	# of attributes
Spatiotemporal	cdw:GeographyDim	4,287	11
	cdw:TimeDim	1,135	8
Sub Total		5,422	19
Environment	cdw:TemperatureDim	15	11
	cdw:HumidityDim	4	4
	cdw:WindDim	4	4
	cdw:PrecipitationDim	3	4
	cdw:AirPollutionDim	42	9
Sub Total		68	32
Health	cdw:VaccineHesitencyDim	18	5
	cdw:ComorbidityDim	24	5
Sub Total		42	10
Socioeconomic	cdw:EthnicityDim	7	5
	cdw:PlaceofDeathDim	10	5
	cdw:OccupationDim	610	11
	cdw:UrbanicityDim	8	5
Sub Total		635	26
Grand Total		6,167	87

In Listing 4, a level member of the `cdw:Day` level, with an IRI value of `day:1`, is implemented using the `qb4o:LevelMember` class. This level member has the attributes `cdw:dayID` with a value of "1", `cdw:dayName` with a value of "2020-01-01", and `cdw:inMonth` with a value of `month:1`. The `cdw:inMonth` represents the relationship of the `cdw:Day` level with its parent level `cdw:Month`. Using the property `qb4o:memberOf`, the relationship of the level member to its containing level is defined. Moreover, the `owl:sameAs` property is used to link to an external data resource. In this case, the resource is the Wikidata entity `wd:Q57396575`, which is the Wikidata entry for January 1, 2020. QB4OLAP employs an observation (an instance of `qb:Observation`) to depict a fact (line 1 in Listing 5). An observation is identified by a distinct IRI and is semantically enriched by combining multiple members from different levels, integrating values for various measure properties. Listing 5 depicts that the dataset of the observation is `cdw:SpatioTemporalDataset` and its cuboid structure consists of `cdw:Admin1`, `cdw:Admin2` and `cdw:Day` levels and two measures `cdw:Confirmed` and `cdw:Deaths`.

6. Description of CovKG

In this section, we describe CovKG from the perspective of both dimensions and facts. Additionally, we provide an overview of the embedded links in CovKG leading to external datasets.

6.1. Dimension and fact overview

Table 3 provides an overview of the dimensions of CovKG. It shows the aspect the dimension represents, the total number of level instances, the number of level attributes. In total, CovKG has 28 levels, 87 level attributes, and 6,167 level members. Nine separate Turtle ABox files are created for nine cuboids. These files are then concatenated with the respective level ABoxes of the dimensions used in the fact tables. Table 4 sheds light on the size measurements

Table 4
Overview of size metrics of the data cuboids.

Semantic cuboid	Source ABox raw size (in MB)	Source ABox processed size (in MB)	# of observations (in million)	# of RDF triples (in million)	Target ABox size (in MB)
Spatiotemporal	838.9	377.4	2.32	16.25	1,719.5
Weather	102.3	26.9	1.34	14.75	1,646.1
Air Pollution	12.8	3.1	0.16	1.14	121.7
Vaccine hesitancy	179.7	0.029	0.002	0.051	499.7
Comorbidity	65.6	2.6	0.032	0.26	27.9
Ethnicity	1.4	0.106	0.0073	0.095	9.6
Place of death	2.7	0.121	0.0084	0.1	10.9
Occupations	4.6	0.035	0.0019	0.06	6.1
Urbanicity	705.4	29.5	1.57	11.04	1,185.0
Total	1,913.4	439.79	5.44	43.75	5,226.5

of CovKG. The raw source ABoxes are available in CSV format, containing a lot of noise data such as irrelevant information not useful for the purpose of this study, negative values, and blank values. The raw source ABoxes are processed and cleaned, after which their sizes reduce significantly, as can be seen in the table.

Another reason for this reduction in size was the utilization of dimension tables. Dimension information in the fact tables, such as names and labels, or floating-point sensor data, are replaced by integer indices pointing to the relevant dimension tables. This information is instead placed in the dimension table, which the fact table can refer to when needed. Additionally, the arrangement of dimension tables in hierarchical levels allows for the performance of OLAP operations on these fact tables. Furthermore, floating-point data of air pollution and weather dimensions are replaced by categorical data, which also contributes to the reduction in size after processing.

The number of RDF triples in Table 4 is significantly larger than the number of observations. This is due to the fact that for each observation, multiple RDF triples are generated. For instance, if an observation has five attributes, containing three dimension attributes and two measure attributes, then there will be seven RDF triples representing that observation in the fact table.

Concatenating the respective dimension level ABox files to the fact ABox files also increase the number of triples. The target ABox sizes are relatively large. As can be seen in Table 4, the spatiotemporal, weather, and urbanicity datasets have gigabyte-scale sizes. This is primarily because the spatiotemporal data contains finer spatial data in the form of second-level administrative unit and first-level administrative unit data for twenty-one countries. Spatiotemporal data was one of the most readily available types of data collected in this study. On the temporal side, day-level data is available, contributing to the increase in size. Additionally, the inclusion of daily confirmed and death data for 3143 U.S. counties significantly contributed to the overall data size.

Urbanicity data was also based on the USA's spatiotemporal data. The weather data was available for all countries at the same fine level as the spatiotemporal data and had more dimension fields than the spatiotemporal data. However, USA's temperature, precipitation, and humidity data were not available, which is why the weather dataset is still smaller in size than the spatiotemporal dataset.

However, the core reason behind the exponential size of the RDF data is that Turtle requires more text characters to represent data than tabular data such as CSV. However, this size tradeoff is reasonable, as the dataset achieves the ability to infer new knowledge based on the available information. Moreover, it gains the capability to be linked to larger knowledge networks to mine further insights, utilizing techniques such as federated query.

6.2. Linking CovKG to external datasets

Linked open datasets contain references to similar elements across other external datasets, allowing for the sharing and reusability of previous knowledge. This also helps in avoiding the inclusion of redundant data,

Table 5

Number of links to external datasets among level instances and the programmatic time taken to link.

Level	Number of links	Processing time (sec)
cdw:Day	1,096	1.712
cdw:Month	36	0.751
cdw:Year	3	0.726
cdw:Individual	872	0.972
cdw:Minor	260	0.496
cdw:Submajor	86	0.482
cdw:Major	20	0.42
cdw:Continent	14	2.292
cdw:Country	44	2.333
cdw:Admin1	856	2.504
cdw:Admin2	7,664	6.414
Total	10,951	19.102

contributing to maintaining scalability. The linking can be done to the concepts in the target TBox as well as level instances in the target ABox. The ABox links can then be referenced by the cuboids. CovKG is linked to a total of four reputable external KGs. This is achieved using the OWL property `owl:sameAs`, as demonstrated in Listing 4 and Listing 6. In the target TBox, all the levels of `cdw:GeographyDim`, `cdw:TimeDim`, `cdw:AirPollutionDim`, `cdw:PrecipitationDim`, `cdw:WindDim`, `cdw:HumidityDim`, as well as the level `cdw:Race` of `cdw:EthnicityDim`, are linked to Wikidata [80] and DBpedia [81] KG. Moreover, levels under `cdw:GeographyDim` are linked to the Geonames KG [56]. Geonames is renowned for collecting and presenting geographical information at highly fine levels in semantic form. Demonstrations of some TBox links are shown in Listing 6 under the comment ‘Level Concepts in target TBox’ (lines 2-7).

In the ABoxes, level instances of the `cdw:GeographyDim` dimension are linked to Wikidata and Geonames. Level instances of the `cdw:TimeDim` dimension’s levels are linked to Wikidata. The level instances of the `cdw:OccupationDim` dimension are linked to Wikidata and the European Skills, Competences, qualifications and Occupations (ESCO) ontology [82]. The ESCO ontology makes the ISCO-08 occupation taxonomy available in semantic form. Demonstrations of some ABox links are shown in Listing 6 under the comment ‘Level instances in target ABox’ (lines 9-17).

This linking process is implemented using the RDFlib Python library [83]. Initially, the IRIs of the level members to external knowledge graphs, such as Wikidata [80], Geonames [56], and ESCO [82], are collected through SPARQL queries on the Wikidata SPARQL endpoint. The query results are imported as CSV files and undergo pre-processing to eliminate duplicates and irrelevant data. Finally, RDFlib is employed to link these IRIs and the ABox triples of the corresponding dimension levels or members. The number of links to the level instances as well as their respective programmatic linking times are reported in Table 5. In summary, CovKG is linked to 10,951 external resources and the total linking time is 19 seconds.

```

1 #Levels Concepts in target TBox
2 cdw:Admin1 a qb4o:LevelProperty;
3     owl:sameAs wiki:Q10864048, geonames:A.ADM1, dbpedia:First-
4     level_administrative_division.
5 cdw:Day a qb4o:LevelProperty;
6     owl:sameAs wiki:Q573, dbpedia:Day.
7 cdw:Humidity a qb4o:LevelProperty;
8     owl:sameAs wiki:Q180600, dbpedia:Air_humidity.
9 #Levels instances in target ABox
10 adm1:1 a qb:LevelMember;
11     cdw:adm1Name "Badakhshan";
12     owl:sameAs wiki:Q165376, geoname:1147745.

```

```

1 12 indiOcc:0110 a qb:LevelMember;
2 13     cdw:individualOccupationName  "Commissioned Armed Forces Officers";
3 14     owl:sameAs wiki:Q108305412, esco:C0110.
4 15 day:1 a qb:LevelMember;
5 16     cdw:dayName  "2020-01-01";
6 17     owl:sameAs wiki:Q57396575.

```

Listing 6: Examples of links to external datasets at conceptual level.

6.3. Availability

The dump files of CovKG can be found at <http://bike-csecu.com/datasets/covid>, and CovKG was stored in the OpenLink Virtuoso Triplestore. Users can remotely access CovKG through the SPARQL endpoint at <http://bike-csecu.com:8890/sparql> and write their own SPARQL queries based on their requirements to obtain answers. To verify the correctness and conduct comparative analysis of CovKG, we developed a set of competency questions (refer to Table 7 and Table 9). All these competency questions are translated into equivalent SPARQL queries to retrieve the answers from CovKG. The set of competency and correctness queries can be accessed, posed to the repository, and answered through a user interface available at <https://bike-csecu.com/datasets/covid/query>. We also provide an interactive OLAP interface, available at <https://github.com/bi-setl/SETL>, allowing users to create their OLAP queries using GUI components and retrieve the answer by posing the query to the related graphs. The OLAP interface is described in Section 7.2.1.

7. Experimental Evaluation

In this section, we discuss experiments conducted on CovKG to evaluate its performance. First the ETL performance is measured by the ETL runtime. After that, we make the qualitative assessment of CovKG. Finally, we present some interesting analytical findings.

7.1. ETL performance overview

Here, the ETL time performance is discussed. The machine on which the ETL is run is a computer of processor Intel(R) Core(TM) i5-8400 CPU. Processor speed was 2.81 GHz with 8 GB RAM. The operating system is Windows 10 Pro 64-bit system.

Table 6 shows the ETL time performance (measured in seconds) in outputting CovKG. We do not record the time for the extraction phase as it depends on users' expertise, internet speed, API performance. Rather, the steps which can be calculated fairly are recorded. It can be seen that larger fact ABoxes such as the spatiotemporal, urbanicity, and weather take longer to pass through the ETL process. The greatest proportion is taken up at RDF loading time, where the facts are loaded to the triple store. We use Openlink Virtuoso [73] as triple store because of generating fast query results, simple interface, and code correction ability.

The source TBox generation time is shown here. The target TBox generation's time was not provided here because it was a manual process where the authors had to carefully design the structure. The entire ETL process takes 2547.271 seconds, around 42 minutes. The majority of the time is spent on RDF loading, primarily due to loading RDF graphs into the triple store. Among the cuboids, the weather cuboid takes the longest time, as it contains the most attributes.

7.2. Qualitative analysis

The previous subsection focused on the quantitative performance evaluation of CovKG. In this section, we assess the quality of CovKG in terms of its business analytical capabilities, compare its performance with other repositories, and evaluate its correctness.

Table 6
ETL program time taken (seconds) by the ETL process for each cuboid.

Semantic cuboid	TBox generation	Source to target mapping	ABox generation	RDF loading	Total (per cuboid)
Spatiotemporal	6.63	1.29	123	227.67	358.59
Weather	10.65	1.23	150	960	1,121.88
Air Pollution	7.68	1.63	8	39.97	57.28
Vaccine Hesitancy	7.48	1.71	1.3	5.62	16.11
Comorbidity	7.71	1.19	5	10.58	24.48
Ethnicity	7.56	1.28	1	5.55	15.39
Place of death	7.53	1.34	2	5.22	16.09
Occupations	8.54	1.30	1.73	1.801	13.371
Urbanicity	7.5	1.22	77	838.36	924.08
Total (per phase)	71.28	12.19	369.03	2,094.771	Grand Total =2,547.271

7.2.1. Enabling business analytics

After generating CovKG, we evaluate its business analytical capabilities. This assessment focuses on whether CovKG has become OLAP-compatible, and to do so, we provide $SETL_{BI}$ as discussed in Section 6.3. $SETL_{BI}$'s *OLAP Layer* enables OLAP analysis over CovKG annotated with MD semantics. Therefore, if CovKG can be loaded into the *OLAP Layer*, OLAP operations can be conducted, and results can be generated, then it confirms that CovKG is ready for business analytics.

In Figure 5, we demonstrate how CovKG is enabled for business analytics using the *OLAP Layer* of $SETL_{BI}$. For brevity, we illustrate the business analytics enabling of only one of the nine fact ABoxes, which is the Vaccine Hesitancy ABox. Figure 5a shows that the Vaccine Hesitancy cuboid is successfully loaded into the tool. To load the CovKG, a user needs to click the *Load File* button (marked by a red rectangle), which prompts to select the target TBox and ABox files. After selecting the files, the tool loads them if the ABox file is OLAP compatible. The upper left part of Figure 5a shows that the target TBox and ABox files have loaded successfully (marked by the yellow rectangle). After loading them, the user selects the Vaccine Hesitancy cuboid from the drop-down list using the arrow icon (marked by the orange square). Next, (s)he clicks the *Extract Cube* button (marked by the orange rectangle) to extract the cuboid's structure. The *Visualization* panel on the left displays dimensions, hierarchies, levels, and measures in a tree view. Users can expand the tree view by clicking on small dots, as shown by the red arrow. Levels, measures, and aggregation functions are selected from this panel.

The *Filtering* panel in the middle shows the available attributes when a level is selected. Relevant dimension ABox triples are also shown if they are part of the cuboid's ABox file. For instance, when `cdw:Country` level (marked by the blue arrow) is selected, attributes like `cdw:countryName` are displayed (as shown by the purple arrow) for the user to choose. Additionally, one can use checkboxes to select instances for slice and dice operations (the three pink arrows). The rightmost *Summary* panel shows the selected level, their associated attributes, and level instances. For the example in the figure, `qb4o:avg` function is selected for the measure `cdw:Confirmed` (marked by the light pink arrow). Figure 5b demonstrates the result of a slice OLAP operation with averaging as the aggregate function, displaying data only for Denmark, Germany, and the Netherlands along the geography dimension.

7.2.2. Comparative analysis of functionality in relation to leading data repositories

We compare the functionality of CovKG with that of other data repositories by first formulating a set of competency queries. These competency queries are designed to assess the ability of the datasets to answer questions with multiple aspects [84]. The questions are listed in Table 7. This evaluation is conducted in comparison to the responses of well-known repositories: Worldometer [24], WHO Dashboard [25], World Bank Dashboard [26], Dynamic Dashboard for Bangladesh [27], and CDC COVID Data Tracker Dashboard [29]. The responses are summarized in Table 8.

Table 7

The competency queries designed for qualitative analysis of CovKG.

Query no.	Competency query statement
Q1	What is the name and population of the level x (example:district) with max confirmed cases on date y(example Januray 1, 2021) in location z (example:Bangladesh)?
Q2	What is the number of deaths among occupation x (example: Engineers) in location y (example: Romania) in the year z (example : 2020) at hot temperature?
Q3	What disease comorbidity has the highest number of deaths during month x (example:February) of year y(example:2020)?
Q4	Which has more infections of Covid-19? Urban(metropolitan) or Rural((non-metropolitan)?
Q5	Which countries have the strongest vaccine hesitancy to vaccines available after one year?
Q6	How many deaths occurred in homelike environments in month x (example:January) of year y (example:2021)?
Q7	What kind of thermal weather has most number of confirmed and/or death on date x (example:January 1, 2021) in location y (example: Feni,Chittagong, Bangladesh)?
Q8	What kind of humidity has most number of confirmed and/or death on date x (example:January 1, 2021) in location y (example: Feni,Chittagong, Bangladesh)?
Q9	What kind of precipitation has most number of confirmed and/or death on date x (example:January 1, 2021) in location y (example: Feni,Chittagong, Bangladesh)?
Q10	What kind of windspeed has most number of confirmed and/or death on date x (example:January 1, 2021) in location y (example: Feni,Chittagong, Bangladesh)?
Q11	How many patients died of Covid-19 in Asia in region of hazardous air pollution in 2020?
Q12	What race has the highest number of deaths in 2021?
Q13	What is the total number of confirmed cases in medical professions?

Table 8

Assessing the comparative functionality of CovKG against prominent data sources by asking each the thirteen competency questions, to which they answer in yes, partially or no.

Competency query no.	CovKG	Bangladesh dashboard[27]	CDC dashboard [29]	WHO dashboard [25]	World Bank dashboard [26]	Worldo meter[24]
Q1	Yes	Partially	Partially	No	No	No
Q2	Yes	No	No	No	No	No
Q3	Yes	No	Yes	No	No	No
Q4	Yes	No	Yes	No	Yes	No
Q5	Yes	No	Partially	No	Partially	No
Q6	Yes	No	Yes	No	No	No
Q7	Yes	No	No	No	No	No
Q8	Yes	No	No	No	No	No
Q9	Yes	No	No	No	No	No
Q10	Yes	No	No	No	No	No
Q11	Yes	No	No	No	No	No
Q12	Yes	No	Yes	No	No	No
Q13	Yes	No	No	No	No	No

7.2.3. Correctness

In our study, data from various sources undergo a semantic ETL process, resulting in CovKG. It is pivotal to ensure the correctness of the ETL process. While the concept of correctness is extensive and beyond the scope of this work, we conducted a partial assessment for CovKG by devising queries for which we already knew the answers. These queries fall into two categories: 1) Common knowledge: Globally recognized information. 2) Special knowledge: Information available from specific sources. Table 9 presents the assessment of CovKG's correctness, and the query details can be found at <https://bike-csecu.com/datasets/covid/query>.

We collected the correct answers from reliable sources. For instance, the correct answers regarding occupa-

Table 9
Assessment of correctness of the ETL process in generating CovKG.

Correctness query	Type	Correct answer	CovKG's answer
How many continents are there?	Common knowledge	7	7
Which year among 2020-2022 is a leap year?	Common knowledge	2020	2020
How many occupations are there under ISCO-08 submajor group?	Common knowledge	43	43
How many nonhispanic white people died of Covid-19 in New Mexico, USA in the month of January, 2021?	Special knowledge	180	180
How many confirmed cases were reported in Ireland's Mayo county in July 12, 2020?	Special knowledge	1,505	1,505
How many confirmed and death cases were reported in Netherland's Buren county on September 21, 2021?	Special knowledge	Confirmed: 2,807 Death : 20	Confirmed: 2,807 Death : 20

tion group numbers were obtained from the International Labour Organization's ISCO-08 classification page (<https://ilostat.ilo.org/resources/concepts-and-definitions/classification-occupation/>). Information on the number of confirmed COVID-19 cases among non-Hispanic whites in the USA and confirmation data for Ireland's Mayo County were sourced from the CDC's COVID Data Tracker [29] and Geohive OpenData repository [41], respectively. Details on the confirmation and death statistics for the municipality of Buren in the Netherlands were gathered from the NL COVID-19 Geo Hub repository [42]. Upon reviewing Table 9, we noted that CovKG provided correct answers for all queries.

7.3. Analytical findings

We analyze CovKG using statistical methods to extract insights from the multidimensional data it represents. Below, we briefly discuss some examples of new insights gained from socioeconomic, health, and environmental aspects. Figure 6 provides insights into the *socioeconomic aspect* of COVID-19, focusing on the occupation factor. The figure illustrates the number of confirmed and death cases among the ten major ISCO-08 occupation classes. Professionals have the highest number of confirmed cases, while services and sales workers have the highest number of deaths. This finding is intriguing because the professionals category includes medical professions such as doctors and nurses. In contrast, services and sales workers encompass professions like travel attendants, transport conductors, travel guides, waiters, cleaning and housekeeping supervisors in offices, hotels, and other establishments, as well as undertakers and embalmers. These occupations involve regular public contact. The relatively high number of confirmed cases but lower death rate among professionals suggests a level of health awareness, even in roles that require frequent public interaction.

Health-related insights can be gained by examining comorbidity, as illustrated in Figure 7. Among various comorbidities, influenza, pneumonia, and respiratory failure stand out with higher numbers of COVID-19 deaths. This correlation is noteworthy because COVID-19, influenza, and pneumonia are all respiratory diseases, suggesting that damage to the respiratory system by one of these conditions may facilitate the spread of the others. This finding

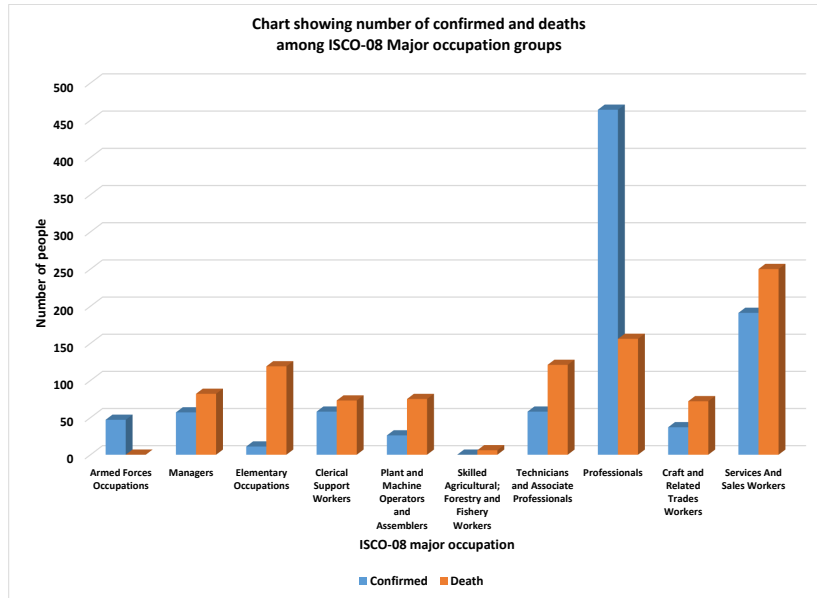


Fig. 6. Total number of deaths and confirmed cases in 2020 among ISCO-08 major occupations.

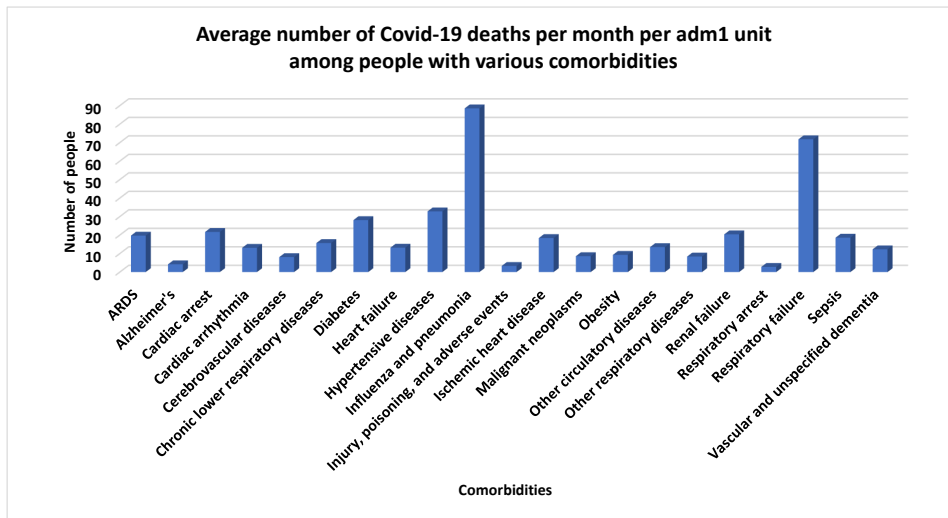


Fig. 7. Number of Covid-19 deaths per month per Admin1 unit among people with various comorbidities.

aligns with previous research on the correlation between influenza and COVID-19, as shown in [85]. Following respiratory diseases, the next highest number of deaths is observed among patients with hypertensive diseases, which are associated with high blood pressure.

Several research studies have concluded that nitrogen dioxide pollution is positively correlated with the transmission and mortality of COVID-19 [86], [61], [87]. These studies were conducted in China and the USA. In our CovKG, we integrate daily subnational air pollution data for South Africa and India. The *environmental aspect* can be observed through the total and average statistics of confirmed and death cases in relation to levels of nitrogen dioxide in South Africa and India, as depicted in Figure 8..

It can be seen in Figure 8 that both the average and total of both deaths and confirmed cases are high for unhealthy levels of nitrogen dioxide. Hazardous and very unhealthy levels exist that are above unhealthy. Yet, unhealthy falls

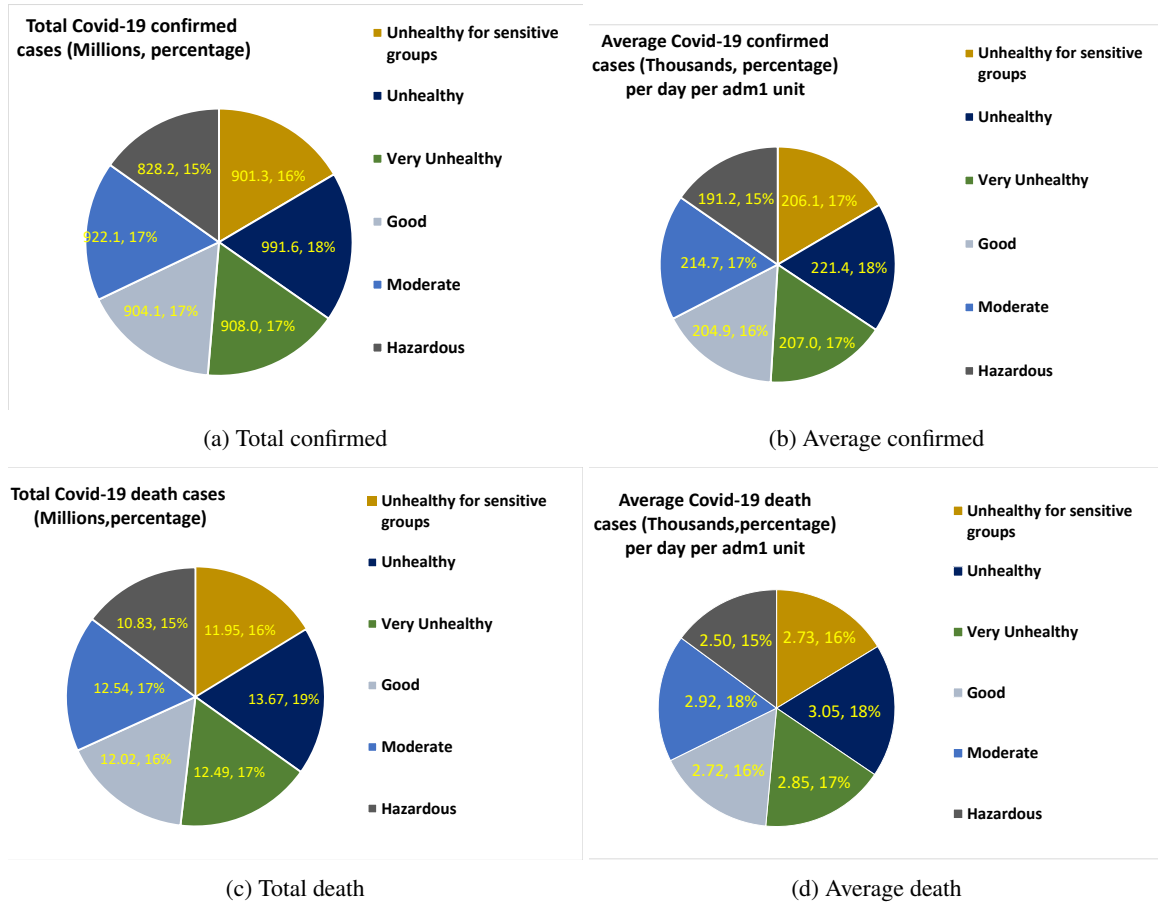


Fig. 8. Covid-19 situation with respect to nitrogen dioxide pollution in India and South Africa shown in millions and percentage for total and thousands and percentage for average per day per Admin1 unit.

in the excessive side of the nitrogen dioxide spectrum. This observation underscores the evident positive correlation with COVID-19’s epidemiology.

8. Conclusion and Future Work

In this study, we generated a multidimensional and semantically annotated Covid-19 knowledge graph titled CovKG that integrates data on Covid-19 epidemiology from disparate sources and facilitates analysis from spatiotemporal, socioeconomic, health, and environmental perspectives. To our knowledge, no previous research generated a multidimensional knowledge graph dedicated to Covid-19 to such an extent as this study. CovKG allows OLAP operations and SPARQL queries to draw new insights from available data. The ETL workflow typically takes around 42 minutes to load CovKG, which is connected to 10,951 external resources, has a size of about 5.3 GB, and consists of about 44 million RDF triples. Moreover, due to being structured as per linked data standards, it is published as per FAIR principles, which is highly essential in cases of global phenomena like Covid-19. The qualitative assessment shows that CovKG is OLAP-compatible, can answer different aspect queries, and yields correct results when compared to other repositories. CovKG was also explored using statistical analysis to get insights into the Covid-19 situation.

CovKG faces limitations due to data sparsity stemming from the unavailability of more comprehensive data sources. The data that was used to create the model were procured from free sources. If paid sources could be

accessed, the data model could be richer. Also, cross-cuboid relationships among some cuboids were not explored. Hence, in future, we will enrich the data model even more by accessing paid sources. We will also apply data mining algorithms and knowledge graph exploration techniques to enable robust cross-cuboid analysis.

References

- [1] C.S. Jensen, T.B. Pedersen and C. Thomsen, Multidimensional databases and data warehousing, *Synthesis Lectures on Data Management* **2**(1) (2010), 1–111.
- [2] S. Negash, Business intelligence, *Communications of the association for information systems* **13**(1) (2004), 15.
- [3] R.P.D. Nath, K. Hose, T.B. Pedersen and O. Romero, SETL: A programmable semantic extract-transform-load framework for semantic data warehouses, *Information Systems* **68** (2017), 17–43.
- [4] R.P. Nath, *Aspects of Semantic ETL*, Aalborg Universitetsforlag, 2020.
- [5] R. Kimball and M. Ross, *The data warehouse toolkit: the complete guide to dimensional modeling*, John Wiley & Sons, 2011.
- [6] L. Etcheverry, S.S. Gomez and A. Vaisman, Modeling and querying data cubes on the semantic web, *arXiv preprint arXiv:1512.06080* (2015).
- [7] S. Chaudhuri and U. Dayal, An overview of data warehousing and OLAP technology, *ACM Sigmod record* **26**(1) (1997), 65–74.
- [8] W.H. Inmon, *Building the data warehouse*, John wiley & sons, 2005.
- [9] B. McBride, The resource description framework (RDF) and its vocabulary description language RDFS, in: *Handbook on ontologies*, Springer, 2004, pp. 51–65.
- [10] FAIR Principles - GO FAIR, (Accessed on 02/03/2024).
- [11] L. Yu, Linked open data, in: *A Developer's Guide to the Semantic Web*, Springer, 2011, pp. 409–466.
- [12] The Linked Open Data Cloud, (Accessed on 01/03/2024).
- [13] D.L. McGuinness, F. Van Harmelen et al., OWL web ontology language overview, *W3C recommendation* **10**(10) (2004), 2004.
- [14] J. Tension, R. Cyganiak and D. Reynolds, The rdf data cube vocabulary, Technical Report, Technical report, W3C Working Draft 05 April, 2012. <http://www.w3.org/TR...>, 2012.
- [15] L. Etcheverry and A.A. Vaisman, QB4OLAP: a new vocabulary for OLAP cubes on the semantic web, in: *Proceedings of the Third International Conference on Consuming Linked Data*, Vol. 905, CEUR-WS. org, 2012, pp. 27–38.
- [16] F. Baader, *The description logic handbook: Theory, implementation and applications*, Cambridge university press, 2003.
- [17] S. Shang, C.K. Leung, Y. Chen and A.G. Pazdor, Spatial Data Science of COVID-19 Data, in: *2020 IEEE 22nd International Conference on High Performance Computing and Communications*, IEEE, 2020, pp. 1370–1375.
- [18] Y. Chen, C.K. Leung, S. Shang and Q. Wen, Temporal data analytics on COVID-19 data with ubiquitous computing, in: *2020 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCLOUD/SocialCom/SustainCom)*, IEEE, 2020, pp. 958–965.
- [19] C.K. Leung, Y. Chen, S. Shang and D. Deng, Big data science on COVID-19 data, in: *2020 IEEE 14th International Conference on Big Data Science and Engineering (BigDataSE)*, IEEE, 2020, pp. 14–21.
- [20] C.K. Leung, Y. Chen, C.S. Hoi, S. Shang and A. Cuzzocrea, Machine learning and OLAP on big COVID-19 data, in: *2020 IEEE International Conference on Big Data (Big Data)*, IEEE, 2020, pp. 5118–5127.
- [21] G. Agapito, C. Zucco and M. Cannataro, COVID-warehouse: A data warehouse of Italian COVID-19, pollution, and climate data, *International Journal of Environmental Research and Public Health* **17**(15) (2020), 5596.
- [22] A. Sakor, S. Jozashoori, E. Niazmand, A. Rivas, K. Bougiatiotis, F. Aisopos, E. Iglesias, P.D. Rohde, T. Padiya, A. Krithara et al., Knowledge4COVID-19: A semantic-based approach for constructing a COVID-19 related knowledge graph from various sources and analyzing treatments' toxicities, *Journal of Web Semantics* **75** (2023), 100760.
- [23] H. Turki, M.A. Hadj Taieb, T. Shafee, T. Lubiana, D. Jemielniak, M.B. Aouicha, J.E. Labra Gayo, E.A. Youngstrom, M. Banat, D. Das et al., Representing COVID-19 information in collaborative knowledge graphs: the case of Wikidata, *Semantic Web* (2022), 1–32.
- [24] Worldometer, Coronavirus Death Toll and Trends—Worldometer (2020). <https://www.worldometers.info/coronavirus/>.
- [25] WHO Coronavirus (COVID-19) Dashboard, (Accessed on 01/17/2023). <https://covid19.who.int/>.
- [26] World Bank, COVID-19 Household Monitoring Dashboard, (Accessed on 01/22/2023). <https://www.worldbank.org/en/data/interactive/2020/11/11/covid-19-high-frequency-monitoring-dashboard>.
- [27] COVID, DGHS, Dynamic Dashboard for Bangladesh. 2021 [visited: 2021 Mar 25], 19. <http://dashboard.dghs.gov.bd/webportal/pages/covid19.php>.
- [28] Data - COVID-19 - Eurostat, (Accessed on 01/17/2023).
- [29] Centers for Disease Control and Prevention., COVID Data Tracker., Atlanta, GA: US Department of Health and Human Services, CDC; 2023, January 22. (Accessed on 01/17/2023). <https://covid.cdc.gov/covid-data-tracker/#datatracker-home>.
- [30] Classifying the Standard Occupational Classification 2020 (SOC 2020) to the International Standard Classification of Occupations (ISCO-08) - Office for National Statistics, (Accessed on 11/04/2022). <https://t.ly/OhWyx>.
- [31] O. Duda, V. Pasichnyk, N. Kunanets, R. Antonii and O. Masiuk, Multidimensional Representation of COVID-19 Data Using OLAP Information Technology, in: *2020 IEEE 15th International Conference on Computer Sciences and Information Technologies (CSIT)*, Vol. 2, IEEE, 2020, pp. 277–280.

- [32] R. Baumgartner, W. Gatterbauer and G. Gottlob, Web Data Extraction System., *Encyclopedia of database systems, Second Edition* **1** (2018).
- [33] K. Chakrabarti, S. Chaudhuri, Z. Chen, K. Ganjam, Y. He and W. Redmond, Data services leveraging Bing's data assets., *IEEE Data Eng. Bull.* **39**(3) (2016), 15–28.
- [34] N. Dalvi, R. Kumar and M. Soliman, Automatic wrappers for large scale web extraction, *arXiv preprint arXiv:1103.2406* (2011).
- [35] Y. Roh, G. Heo and S.E. Whang, A survey on data collection for machine learning: a big data-ai integration perspective, *IEEE Transactions on Knowledge and Data Engineering* **33**(4) (2019), 1328–1347.
- [36] Welcome - Humanitarian Data Exchange, (Accessed on 08/29/2022).
- [37] Afghanistan: Coronavirus(COVID-19) Subnational Cases - Humanitarian Data Exchange, (Accessed on 08/29/2022).
- [38] Data on the weekly subnational 14-day notification rate of new COVID-19 cases, (Accessed on 08/29/2022).
- [39] A. Naqvi, COVID-19 European regional tracker, *Scientific Data* **8**(1) (2021), 1–14.
- [40] covid19-data-greece/data/greece/regional at master · Covid-19-Response-Greece/covid19-data-greece, (Accessed on 01/24/2023).
- [41] GeoHive Open Data, (Accessed on 01/24/2023).
- [42] NL COVID-19 Hub, (Accessed on 01/24/2023).
- [43] COVID-19 in India | Kaggle, (Accessed on 08/29/2022).
- [44] South Africa Provincial Breakdown | Covid-19 South Africa, (Accessed on 01/24/2023).
- [45] A. Haratian, H. Fazelinia, Z. Maleki, P. Ramazi, H. Wang, M.A. Lewis, R. Greiner and D. Wishart, Dataset of COVID-19 outbreak and potential predictive features in the USA, *Data in Brief* **38** (2021), 107360.
- [46] W.W. Online, Worldweatheronline.com (2016), (Accessed on 08/29/2022).
- [47] P. Gwaze and S.H. Mashele, South African Air Quality Information System (SAAQIS) mobile application tool: Bringing real time state of air quality to South Africans, *Clean Air Journal* **28**(1) (2018), 3–3.
- [48] CCR, (Accessed on 01/24/2023).
- [49] H. Khan, M.E. Dabla-Norris, F. Lima and A. Sollaci, *Who doesn't want to be vaccinated? Determinants of vaccine hesitancy during COVID-19*, International Monetary Fund, 2021.
- [50] Conditions Contributing to COVID-19 Deaths, by State and Age, Provisional 2020-2022 | Data | Centers for Disease Control and Prevention, (Accessed on 09/04/2022).
- [51] E. Sari, G. Kağan, B.Ş. Karakuş and Ö. Özdemir, Dataset on social and psychological effects of COVID-19 pandemic in Turkey, *Scientific Data* **9**(1) (2022), 1–7.
- [52] B. Windsor-Shellard and R. Nasir, Coronavirus (COVID-19) related deaths by occupation, England and Wales: deaths registered between 9 March and 28 December (2021).
- [53] M.-G. Hâncean, J. Lerner, M. Perc, I. Oană, D.-A. Bunaciu, A.A. Stoica and M.-C. Ghiță, Occupations and their impact on the spreading of COVID-19 in urban communities, *Scientific reports* **12**(1) (2022), 1–12.
- [54] M.C.S. Gaddekar, Air Quality Index (AQI) Basics, *Journal homepage: www.ijrpr.com ISSN 2582* (2022), 7421.
- [55] Data Access - Urban Rural Classification Scheme for Counties, (Accessed on 01/24/2023).
- [56] G. GeoNames, The GeoNames geographical database, 2004.
- [57] K. Piotrowicz, D. Ciaranek, A. Wypych, A. Razi and J. Mika, Local weather classifications for environmental applications, *Aerul și Apa. Componente ale Mediului= Air and Water. Components of the Environment* **2013** (2013).
- [58] M. Travaglio, Y. Yu, R. Popovic, L. Selley, N.S. Leal and L.M. Martins, Links between air pollution and COVID-19 in England, *Environmental pollution* **268** (2021), 115859.
- [59] E.B. Brandt, A.F. Beck and T.B. Mersha, Air pollution, racial disparities, and COVID-19 mortality, *Journal of Allergy and Clinical Immunology* **146**(1) (2020), 61–63.
- [60] X. Wu, R.C. Nethery, M.B. Sabath, D. Braun and F. Dominici, Air pollution and COVID-19 mortality in the United States: Strengths and limitations of an ecological regression analysis, *Science advances* **6**(45) (2020), eabd4049.
- [61] D. Liang, L. Shi, J. Zhao, P. Liu, J.A. Sarnat, S. Gao, J. Schwartz, Y. Liu, S.T. Ebel, N. Scovronick et al., Urban air pollution may enhance COVID-19 case-fatality and mortality rates in the United States, *The Innovation* **1**(3) (2020), 100047.
- [62] Ten threats to global health in 2019, (Accessed on 09/04/2022).
- [63] A. Sanyaolu, C. Okorie, A. Marinkovic, R. Patidar, K. Younis, P. Desai, Z. Hosein, I. Padda, J. Mangat and M. Altaf, Comorbidity and its impact on patients with COVID-19, *SN comprehensive clinical medicine* **2**(8) (2020), 1069–1076.
- [64] B. Wang, R. Li, Z. Lu and Y. Huang, Does comorbidity increase the risk of patients with COVID-19: evidence from meta-analysis, *Aging (albania NY)* **12**(7) (2020), 6049.
- [65] W.-j. Guan, W.-h. Liang, Y. Zhao, H.-r. Liang, Z.-s. Chen, Y.-m. Li, X.-q. Liu, R.-c. Chen, C.-l. Tang, T. Wang et al., Comorbidity and its impact on 1590 patients with COVID-19 in China: a nationwide analysis, *European Respiratory Journal* **55**(5) (2020).
- [66] D. Pan, S. Sze, J.S. Minhas, M.N. Bangash, N. Pareek, P. Dival, C.M. Williams, M.R. Oggioni, I.B. Squire, L.B. Nellums et al., The impact of ethnicity on clinical outcomes in COVID-19: a systematic review, *EclinicalMedicine* **23** (2020), 100404.
- [67] P. Patel, L. Hiam, A. Sowemimo, D. Devakumar and M. McKee, Ethnicity and covid-19, Vol. 369, British Medical Journal Publishing Group, 2020.
- [68] P. Chin-Hong, K.M. Alexander, N. Haynes and M.A. Albert, Pulling at the heart: COVID-19, race/ethnicity and ongoing disparities, *Nature Reviews Cardiology* **17**(9) (2020), 533–535.
- [69] M.J. Townsend, T.K. Kyle and F.C. Stanford, Outcomes of COVID-19: disparities in obesity and by ethnicity/race, Vol. 44, Nature Publishing Group, 2020, pp. 1807–1809.
- [70] H. Ali, A. Alshukry, S.K. Marafie, M. AlRukhayes, Y. Ali, M.B. Abbas, A. Al-Taweel, Y. Bukhamseen, M.H. Dashti, A.A. Al-Shammari et al., Outcomes of COVID-19: Disparities by ethnicity, *Infection, Genetics and Evolution* **87** (2021), 104639.

- [71] I.L. Office, *International Standard Classification of Occupations 2008 (ISCO-08): Structure, group definitions and correspondence tables*, International Labour Office, 2012.
- [72] D.D. Ingram and S.J. Franco, *2013 NCHS urban-rural classification scheme for counties*, Vol. 2014, US Department of Health and Human Services, Centers for Disease Control and . . . , 2014.
- [73] O. Erling, Virtuoso, a Hybrid RDBMS/Graph Column Store., *IEEE Data Eng. Bull.* **35**(1) (2012), 3–8.
- [74] M.A. Musen, The protégé project: a look back and a look forward, *AI matters* **1**(4) (2015), 4–12.
- [75] S. Lohmann, V. Link, E. Marbach and S. Negru, WebVOWL: Web-based visualization of ontologies, in: *Knowledge Engineering and Knowledge Management: EKAW 2014 Satellite Events, VISUAL, EKMI, and ARCOE-Logic, Linköping, Sweden, November 24-28, 2014. Revised Selected Papers. 19*, Springer, 2015, pp. 154–158.
- [76] T. Knap, P. Hanečák, J. Klímek, C. Mader, M. Nečaský, B. Van Nuffelen and P. Škoda, UnifiedViews: an ETL tool for RDF data management, *Semantic Web* **9**(5) (2018), 661–676.
- [77] W.W.W. Consortium et al., R2RML: RDB to RDF mapping language (2012).
- [78] M. Arenas, A. Bertails, E. Prud’hommeaux, J. Sequeda et al., A direct mapping of relational data to RDF, *W3C recommendation* **27** (2012), 1–11.
- [79] R.P. Deb Nath, O. Romero, T.B. Pedersen and K. Hose, High-level ETL for semantic data warehouses, *Semantic Web* **13**(1) (2022), 85–132.
- [80] D. Vrandečić and M. Krötzsch, Wikidata: a free collaborative knowledgebase, *Communications of the ACM* **57**(10) (2014), 78–85.
- [81] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak and Z. Ives, Dbpedia: A nucleus for a web of open data, in: *international semantic web conference*, Springer, 2007, pp. 722–735.
- [82] J. De Smedt, M. le Vrang and A. Papantoniou, ESCO: Towards a Semantic Web for the European Labor Market., in: *Ldow@ www*, 2015.
- [83] D. Krech, RdfLib: A python library for working with rdf, *Online* <https://github.com/RDFLib/rdfLib> (2006).
- [84] A. Ghose, K. Hose, M. Lissandrini and B.P. Weidema, An open source dataset and ontology for product footprinting, in: *The Semantic Web: ESWC 2019 Satellite Events: ESWC 2019 Satellite Events, Portorož, Slovenia, June 2–6, 2019, Revised Selected Papers 16*, Springer, 2019, pp. 75–79.
- [85] B. Alosaimi, A. Naeem, M.E. Hamed, H.S. Alkadi, T. Alanazi, S.S. Al Rehily, A.Z. Almutairi and A. Zafar, Influenza co-infection associated with severity and mortality in COVID-19 patients, *Virology journal* **18**(1) (2021), 1–9.
- [86] Y. Yao, J. Pan, Z. Liu, X. Meng, W. Wang, H. Kan and W. Wang, Ambient nitrogen dioxide pollution and spreadability of COVID-19 in Chinese cities, *Ecotoxicology and environmental safety* **208** (2021), 111421.
- [87] J. Lipsitt, A.M. Chan-Golston, J. Liu, J. Su, Y. Zhu and M. Jerrett, Spatial analysis of COVID-19 and traffic-related air pollution in Los Angeles, *Environment International* **153** (2021), 106531.