The PANORAMA knowledge graph based prediction of outdoor air quality impact on health

Nareesa Karmali¹, Cedric Gil-Jardiné¹, Gavo Diallo¹

¹Team AHeaD, Bordeaux Population Health Inserm 1219, Univ. Bordeaux, F-33000, Bordeaux, France {first, last}@u-bordeaux.fr

Abstract

Current air quality surveillance ecosystems cannot 1 monitor air pollution levels across time and space 2 continuously, and lack the ability to integrate his-3 torical data from heterogeneous sources into arti-4 ficial intelligence algorithms capable of providing 5 more insights on the impact of severe air pollution. 6 Integrating in-situ sensor and satellite data has at-7 tracted recent attention, but the characteristics of 8 graph-based learning with heterogeneous air qual-9 ity data have not been explored in depth. 10

In this paper, we present the PANORAMA knowl-11 edge graph based prediction approach of the links 12 between pollutants exposure and health outcomes. 13 PANORAMA aims at integrating heterogeneous air 14 quality and health data in order to better predict and 15 explain their relationships. An Extract, Transform, 16 and Load process is used to design and incorpo-17 rate discrete data into the knowledge graph, and a 18 knowledge graph embedding model is trained and 19 tested to assess inferential capabilities of the graph. 20 A use case with a dataset related to the Gironde de-21 partment in France is showcased. The promising 22 preliminary results of a 31.376 MR, 0.312 MRR, 23 and 0.254, 0.335, and 0.524 Hits@ 1, 3, and 10 re-24 spectively on a TransE model suggest the potential 25 of leveraging semantic structuring and knowledge 26 graph technology for environmental public health.

27

1 Introduction 28

There has been a correlation between climate and air qual-29 ity and the incidence of disease over the past few decades. 30 Despite this, real-time information regarding air pollution is 31 limited due to the fragmentation of data sources. Accord-32 ing to the World Health Organization (WHO), air pollution 33 causes over 7 million premature deaths annually, the major-34 ity of which occur in low- and middle-income countries¹. By 35 2060, the number of deaths attributable to outdoor air pollu-36 tion is expected to double, yet many countries lack public air 37 quality information. Recent reports indicate that 141 coun-38 tries do not regularly monitor particulate matter, and only 23 39

countries have more than three monitors per million citizens [Holloway *et al.*, 2021].

40

41

Gaseous pollutants and particulate matter contribute to the 42 development and progression of disease[Kampa and Cas-43 tanas, 2008]. Globally, ground monitors are the most pop-44 ular source of air quality data, but their limited observation 45 range means they may not accurately reflect the air quality 46 where populations live and work [Ma et al., 2019]. Due 47 to its expansive spatial and geographic coverage and consis-48 tent data quality, satellite monitoring has gained popularity 49 in the past decade. Satellite data cannot directly assess par-50 ticulate matter concentration; aerosol optical depth data must 51 be post-processed into particulate matter concentrations. IoT 52 and other portable sensor options have emerged, enabling 53 low-cost, wearable technologies to monitor real-time air qual-54 ity. Integrating in-situ sensor and satellite data has become 55 increasingly crucial for providing comprehensive air quality 56 data for public health planning, policy development, and re-57 search. 58

There is gap regarding the automated prediction or fore-59 casting of disease caused by air quality factors while integrat-60 ing the various available data sources for air quality. This is 61 particularly the case in the context of Low and Middle Income 62 Countries (LMICs) where in-situ sensors are lacking [Raheja 63 et al., 2022]. Integrating information from multiple sources, 64 such as ground and satellite sensors, results in a multimodal 65 learning approach for air quality data analysis. Rather than 66 relying solely on a single modality for understanding the ef-67 fects of air pollutants, data integration enables a more com-68 prehensive knowledge base and thorough analysis. This area 69 has not been thoroughly explored. 70

In this paper, we described the PANORAMA approach 71 (Public heAlth iNdicatOrs generation from Air quality data) 72 for integrating heterogeneous air quality data and health re-73 lated outcomes into a triple-based knowledge graph [Somé et 74 al., 2019; Some et al., 2019; Hogan et al., 2022] in order to 75 enable predicting the relationship between outdoor air qual-76 ity and diseases and later explaining this relationship. We 77 report the design and development of the first version of the 78 PANORAMA knowledge graph, as well as the preliminary 79 results obtained by applying a missing links prediction pro-80 cess over the graph to detect evidenced association between 81 pollutants and health outcomes (diseases). In contrast to pre-82 vious studies, which employs traditional machine learning 83

¹https://www.who.int/westernpacific/health-topics/air-pollution

techniques for such a task, we use symbolic AI and knowl-edge graph-based techniques.

86 2 Related Work

In this section, we highlight some relevant approaches on air
quality monitoring and the prediction of the impact of poor air
quality on health outcomes. They are categorized in two categories: traditional Machine Learning (ML) based approaches
and graph based approaches.

92 2.1 Traditional ML approaches

Statistical and machine learning prediction methods have re-93 ceived increasing attention in recent years for the advantages 94 they provide in describing complex linear and non-linear rela-95 tionships [Wen et al., 2019]. Prediction methods involving air 96 quality have been categorized into two main classes: physi-97 cal and data-driven methods. Physical-based approaches have 98 been widely used for air quality modelling in the environ-99 mental community, and includes strategies such as chemical 100 transport models, operational street canyon models, Gaus-101 sian plume models, and reduced-scale models. An advantage 102 of physical models is their ability to provide insights about 103 physical and chemical processes that govern atmospheric 104 pollutants and thus explain the variability in the chemical 105 and dynamical mechanisms of pollutants [Zhao et al., 2019; 106 Zhou *et al.*, 2019]. However, physical methods often require 107 a great deal of empirical assumptions and expert-based ad-108 justments to adapt to different regions and contexts. Multiple 109 linear regression has been used quite extensively for under-110 standing the relationships between air pollution observations, 111 health, and other meteorological variables. This technique 112 was built upon when creating land-use regression, which ex-113 tends the linear regression model to make use of land char-114 acteristics, road use, traffic, and other physical environmental 115 variables as predictors. This spatial interpolation modelling 116 method involves the inclusion of Euclidean buffers around 117 pollution monitoring stations, while conventional linear re-118 gression only reflects the coefficients within a study area 119 [Beckerman et al., 2013; Shaddick et al., 2018]. With this 120 method, land use information can be used to interpolate how 121 pollution may be distributed at locations that do not have di-122 rect monitors [Beckerman et al., 2013]. However, land-use 123 regression has been criticized for its inflexibility and limita-124 tions when dealing with time-series or temporal data, pre-125 dictor interaction, or collinearity [Beckerman et al., 2013; 126 Kerckhoffs et al., 2019]. Following the use of multiple linear 127 regression and Poisson regression for time-series data, gen-128 eralized additive models (GAM) became popular in the late 129 1990s for dealing with complex and non-linear time-series 130 data both widely and in the air quality field [Ma et al., 2020]. 131 The GAM is more flexible for seasonality and time trends, as 132 well as nonlinear relationships, and defines the model using a 133 non-parametric smoothing function [Khojasteh et al., 2021]. 134

The study in [Khojasteh *et al.*, 2021] used the Dickey-Fuller test to determine the long-term effects of air pollution on respiratory morbidity and mortality and developed nonlinear autoregressive and artificial neural network models to accurately predict it. They found that NO2, SO2, O3, and PM10 did not affect respiratory morbidity and mortality in the sensitivity analysis.

142

179

2.2 Graphs based approaches

There have been a variety of use cases of graph learning in 143 healthcare. Patient records for instance, particularly elec-144 tronic health records (EHRs), have been commonly repre-145 sented as networks and used to personalize patient predic-146 tions for precision medicine [Li et al., 2022]. However, the 147 field of graph learning in environmental public health has not 148 been thoroughly investigated, and the characteristics of graph 149 representation of heterogeneous air quality sensor data are 150 not fully understood. Due to the complex topographic struc-151 ture and diverse types of interactions between nodes, graph 152 learning can be challenging to develop and implement [Li et 153 al., 2022]. [Ferrer-Cid et al., 2021] described the topology 154 of an air pollution monitoring network using graph learn-155 ing techniques. The primary nodes of air pollution moni-156 toring networks are ground sensors, but the limited coverage 157 of these sensors demonstrates the critical need for modeling 158 techniques that can increase spatial resolution and utilize the 159 knowledge of neighboring sensors to reconstruct the signal 160 in regions without sensor data [Ferrer-Cid et al., 2021]. The 161 study in [Lin et al., 2020] integrates theories of spatial statis-162 tics into neural networks for the prediction of geographically-163 based spatiotemporal phenomena on a fine spatial scale. Us-164 ing observed air quality measurements from a network of 165 low-cost sensors as well as contextual data describing envi-166 ronmental characteristics, the objective is to estimate air qual-167 ity values at locations lacking sensors across a fine spatial 168 grid. 169

The current study uses a knowledge graph, a heterogeneous 170 graph that captures knowledge from different datasets, to in-171 tegrate a variety of air quality and health outcome data [Li et 172 al., 2022]. Knowledge graphs are networks of data expressed 173 through nodes and edges corresponding to important entities 174 of information and the relations that link them together, re-175 spectively. Using the principles of Linked Data, knowledge 176 graphs provide a tool for combining, exploring, and travers-177 ing data in context with other connected resources. 178

3 The PANORAMA approach

Our objective is to develop a versatile method that can inte-180 grate diverse data types and modalities to predict the effects 181 of exposure to poor air quality. To achieve this, we aim to rely 182 on readily available Open Data, especially in LMICs, where 183 reliable statistics and monitoring tools could be scarce. Ini-184 tially, we integrated satellite data on air quality monitoring, 185 ground sensor data, geographical subdivision data, and hos-186 pital admissions data to validate the approach's fundamental 187 principle. Our primary focus was on the Gironde department 188 in the southwest of France, where we were able to retrieve the 189 necessary multimodal data to apply the approach. By leverag-190 ing these data types, we aim to develop a generic method that 191 can be used in the future to predict the effects of air pollution 192 exposure. 193

Figure 2 exhibits the overall workflow of the PANORAMA 194 approach, which consists of the steps described in the following sections. 196

197 **3.1 Design of the knowledge model**

A model was designed based on nodes of interest to be in-198 tegrated regarding air quality data monitoring and the rela-199 tionships between these nodes. Therefore, the Knowledge 200 graph construction begins with modelling the nodes and rela-201 tions that are set to be included in the graph. The knowledge 202 graph is a graph-based information represented in the form of 203 triplets, to represent entities and their relationships [Li et al., 204 2020] that are a set of symbols. Nodes were selected based on 205 discrete air quality data available for the Gironde department. 206 The links and their corresponding labels were more partic-207 ular, and thus chosen depending on the semantics available 208 within each discrete dataset and the intersecting information 209 between the datasets. Custom names for relations were as-210 signed to links between the nodes. The high level view of the 211 graph is depicted in figure 1. 212

213 **3.2 Knowledge sources selection**

The resulting conceptual model in the previous step was used 214 to determine the information to be retrieved from the knowl-215 216 edge sources. Several knowledge sources were assessed for location, reliability and frequency of measurements, range 217 of pollutants, and availability of health outcome data. Such 218 evaluation of knowledge sources was necessary to understand 219 how fragmented, unavailable, and unreliable data may af-220 fect the ability to integrate and produce downstream findings. 221 While this is an expected challenge considering the data of in-222 terest, it is essential that the PANORAMA knowledge graph 223 is developed and tested on consistent and valid data to ensure 224 its applicability in settings where data may be more disparate, 225 difficult to obtain, or unreliable. Two air quality data sources 226 and one health outcome data source were selected for their 227 frequency of measurements, spatial and temporal coverage, 228 229 reliability, and ability to align location.

230 Copernicus Atmospheric Monitoring Service (CAMS)

The Copernicus Atmospheric Monitoring Service (CAMS) is 231 the European Union's Earth Observation Programme, mon-232 233 itoring and forecasting European air quality and pollutants through satellite and non-satellite observations². CAMS pro-234 vides consistent and quality-controlled information through 235 daily forecasts of air composition and forecasts of global 236 long-range transport of pollutants. CAMS focuses on five 237 main areas: Air quality and atmospheric composition, ozone 238 layer and ultraviolet radiation, emissions and surface fluxes, 239 solar radiation, and climate forcing. Free and openly ac-240 cessible, CAMS provides data from two groups of missions: 241 (i) the Sentinels, which have been developed specifically for 242 CAMS and provided dedicated observations, and (ii) the Con-243 tributing Missions operated by National, European, or Inter-211 national organizations and siphon data to CAMS. For satel-245 246 lite data obtained from CAMS, which includes measurements 247 with a date, value, the pollutant that is measured, and the geographical coordinates of the measurement, the following cus-248 tom links were established: hasValue, hasDate, hasPollutant, 249 fromLocation. 250

Atmo Nouvelle-Aquitaine pollutants data

Atmo Nouvelle-Aquitaine is the regional air observatory in 252 the Nouvelle-Aquitaine region of France³. According to the 253 law on air and rational use of energy, Atmo has several goals, 254 including the monitoring of air 24 hours a day and forecast-255 ing the air quality in the region, predicting pollution episodes 256 and altering the authorities. As of December 2020, Atmo 257 Nouvelle-Aquitaine has a fixed sensor network of 44 pollu-258 tion measurement stations and over 100 analyzers, which op-259 eration continuously days a week, 24 hours a day. For data 260 obtained from Atmo service, which includes measurements 261 with a unit, date, value, and name, the pollutant that is mea-262 sured, and the location of the sensor, the following custom 263 links were established: hasValue, hasDate, hasUnit, hasPol-264 lutant, fromLocation. 265

SAMU Urgences de France

SAMU Urgences de France (SUDF) is the French organiza-267 tion of prehospital emergency medicine and ensure appropri-268 ate medical attention, from which emergency calls receive the 269 most appropriate response as soon as possible according to 270 the urgency and severity. From the emergency call, a variety 271 of services are available depending on the nature of the call. 272 These options include medical advice, a private ambulance, a 273 general practitioner, etc. There is one SAMU per French de-274 partment, approximately one for every 500,000 inhabitants, 275 and nearly 100 in total. SAMU data are organized at the level 276 of entries. Specifically, entry-level data are collected for each 277 contact with SAMU and include the entry identification num-278 ber, date, main complaint, INSEE area code, age and sex of 279 the individual, and possible diagnostic information including 280 an International Classification of Diseases (ICD) code. For 281 data obtained from SAMU Urgences de France, which in-282 cludes entries with an identification number, date, main com-283 plaint, sex, age, possible diagnosis, and the location of the en-284 try, the following custom links were established: hasReason, 285 hasDate, isSex, isAge, hasCondition, hasICD, hasDiagnosis, 286 fromLocation. 287

The National Institute of Statistics and Economic Studies 288 The National Institute of Statistics and Economic Studies 289 (INSEE) collects, analyses and disseminates information on 290 the French economy and society⁴. Created in 1946, INSEE 291 is a Directorate-General of the Ministries for the Economy 292 and for Finances and is the official body for the design, pro-293 duction and dissemination of official statistics. Pertaining to 294 the present study, INSEE describes and analyzes regions and 295 territories in France according to codes. Within France, IN-296 SEE provides codes for each region, department, arrondisse-297 ment, canton, and municipality, and thus provides granular 298 geographic and regional information across the country. For 299 specificity, INSEE provides codes, descriptors, and the rela-300 tionship between geographical levels. All data provided by 301

²Copernicus Services — Copernicus. Accessed April 25, 2023. https://www.copernicus.eu/en/copernicus- services

³Atmo Nouvelle-Aquitaine, l'observatoire régional de l'air. Atmo Nouvelle-Aquitaine. Published December 5, 2016. Accessed April 25, 2023. https://www.atmonouvelleaquitaine.org/article/atmo-nouvelle-aquitainelobservatoire-regional-de-lair

⁴Getting to know INSEE — Insee. Accessed April 25, 2023. https://www.insee.fr/en/information/2381925



Figure 1: High level view of the conceptual model of the knowledge graph. The link to be predicted is the one between the health outcomes node (Emergency Help) and the Pollutants node.

INSEE is open-access and free. The department of Gironde,
which is of interest in this study, is associated with code 33,
and the codes for each arrondissement, canton, and municipality in this department begin with this code and add subsequent digits as the geographic level of interest becomes more
granular.

308 International Classification of Diseases, Version 10

The WHO International Classification of Diseases (ICD) Version 10 Vocabulary was obtained from BioPortal, an online biomedical ontology and terminology repository [Noy *et al.*, 2009]. Medical and clinical terminology are described in the ICD as codes, and these codes are used as the main foundation for most disease and health records and statistics in medical care, research, and public health.

316 **3.3** Extraction and transformation of knowledge

Information was extracted from the knowledge sources, 317 cleaned, and made available as triples in Turtle format with 318 PANORAMA identifiers and the name of each relationship 319 between entities. CAMS data was delivered in netCDF for-320 mat for all global geographical coordinates and required the 321 initial step of conversion to tabular Comma-Separated Value 322 (CSV) format prior to cleaning. Post-conversion, data was 323 restricted to the geographical coordinates of Gironde, France 324 according to Google coordinates. Data obtained from Atmo 325 and CAMS were subject to the same preparation: (i) Dates 326 between 01-01-2017 and 31-12-2019; (ii) Only pollutants of 327 interest included: PM2.5, PM10, CO, NO2, O3, SO2. Data 328 obtained from SUDF was similarly prepared: (i) Dates be-329 tween 01-01-2017 and 31-12-2019; (ii) Only entries which 330 include a possible diagnosis (ICD code). Then, post-cleaning 331 was necessary in order to normalize the entities to ensure 332 that links could be created between the disparate knowl-333 edge sources. Indeed, in air quality and health data, dif-334 ferent terms or descriptions exist for the same entity. En-335 tity normalization is required to map the terms, ensuring 336 that linking entities are represented analogously between re-337 sources. In the PANORAMA knowledge graph, all nodes 338 must be harmonized according to their location (which acts 339 as a proxy) to create links between the selected knowledge 340 sources. Entity normalization was achieved using INSEE ge-341 ographical data, which describes different levels of regions 342 in Gironde as codes. SUDF data was previously mapped 343 to INSEE geographical codes prior to data being obtained. 344 For instance, the Atmo Sensor Location << Bordeaux -345 Bastide >> was normalized to the INSEE Commune << 346 Bordeaux >> with the INSEE Code "33063.0", while 347

<< Ambes 2 >> was normalized to "Ambes" with IN- 348SEE Code << 33528.0 >>. CAMS data described location 349 through geographical coordinates to the tenth place. Location was mapped to INSEE codes accordingly. For instance, 351 CAMS Coordinate << 44.8"N, -0, 6"W >> was mapped 352 to the INSEE Commune << Bordeaux >> with the code 353 << 33063.0 >>. 354

3.4 Loading into the triplestore

RDF statements are comprised of a three-part structure link-356 ing resources, which can be identified with a URI, together 357 depicting two entities and the relation between them. This se-358 mantic data is represented as subject-predicate-object triples 359 for implementation into a knowledge graph. A set of triples 360 makes an RDF graph, and in a graph, nodes represent enti-361 ties and edges represents relationships between entities. RDF 362 triples are well suited for representing highly interconnected 363 data. The set of PANORAMA triples are loaded and intercon-364 nected in a knowledge graph database. In the present study, 365 11 custom relations are defined in this study: hasValue, has-366 Date, hasUnit, hasPollutant, fromLocation, hasReason, isSex, 367 isAge, hasCondition, hasICD, hasDiagnosis. 368

Triples were extracted from unstructured tabular data us-369 ing RDFLib 6.1.1, a Python package for working with RDF. 370 For each data source, a tabular data frame was loaded and 371 parsed with RDFLib. A namespace of 'http://panorama.org/' 372 was set. A namespace is a mapping that connects URIs to 373 a set prefix. For instance, a URI concept was created for 374 each measurement accordingly for Atmo Nouvelle-Aquitaine 375 and CAMS data: measureConcept = URIRef("http:376 *//panorama.org/Measure*"). Finally, all the triples have 377 been loaded into a semantic graph database, Blazegraph⁵. 378 The latter is an open source Triplestore and graph database 379 for storing and querying linked data and support RDF and 380 SPARQL Protocol and RDF Query Language (SPARQL). 381

3.5 Predicting the links between pollutants and health outcomes

As already mentioned, we turn the identification of links be-384 tween pollutants and health outcomes to a knowledge graph 385 completion problem. Generally, knowledge graphs are in-386 complete sets with many missing facts. Thus, Knowledge 387 Graph Completion aims to fill in missing triples to support 388 graph completeness and improve prediction and the perfor-389 mance of other downstream applications, such as link predic-390 tion [Guan et al., 2018]. Knowledge graph embeddings are 391

382 383

⁵https://blazegraph.com/



Figure 2: Overall workflow of PANORAMA. Nodes are selected based on interest for integration and connected with edges to form the conceptual knowledge graph. After loading in a triplestore, the embedding model is developed and link prediction is conducted to predict the possible association between pollutants and health outcomes.

supervised learning models that learn vector representations 392 of nodes and edges of labeled, directed multigraphs. The aim 393 is to map and project entities and relations into a low- di-394 mensional and continuous vector space, on which other oper-395 ations can be conducted. In our case, a Knowledge Graph 396 Embedding Model (KGEM) was trained and evaluated on 397 the PANORAMA knowledge graph using several model algo-398 rithms and a set of hyperparameters (see figure 4. A training 399 knowledge graph is subjected to the architecture of a model 400 with a couple layers: (i) lookup layer (simply assigns an em-401 bedding to each node and edge type of the graph); (ii) scoring 402 layer, interacting with a loss function; (iii) negatives gener-403 ation process; (iv) optimization (to train embeddings, which 404 are used by the above scoring function to predict new links, 405 the downstream task). 406

In this study, we have tested three different embeddings models, respectively Translating Embeddings (TransE)[Bordes *et al.*, 2013], RotatE[Huang *et al.*, 2021] and ComplEx[Trouillon *et al.*, 2016] which is a semantic matching model that serves as an alternative to translational distance models. RotatE is another translational distance model inspired by TransE, but instead switches to complex vector



Figure 3: Examples of triples extracted in the study. A measure node, acquired from tabular Atmo or CAMS data, is assigned to the http://panorama.org/ namespace.



Figure 4: A KGEM comprises several layers: a lookup layer, a scoring layer, a negatives generation process, and optimization. Embedding of nodes and links of the original knowledge graph allows for downstream tasks such as link prediction.

space C to define relations as rotational from head to tail of a 414 triple. 415

416

417

4 **Results**

4.1 The PANORAMA Knowledge graph

The graph is depicted with all conceptual nodes and edges in 418 figure 6. It contains 1,922,998 triples, consisting of 443,332 419 entities and 13 relation types (Table 1). The graph con-420 tains entities from CAMS, Atmo, SUDF, ICD10, and INSEE. 421 There are two pre-defined relations in the graph acquired 422 from RDF Vocabulary: rdf : type, which is an instance of 423 rdf : *Property* that is used to declare that a resource is an 424 instance of a class, and rdfs : *label*, which is another instance 425 of *rdf* : *Property* that provides a label or readable version 426 of a resource's name. An example of Turtle format rendering 427 for a pollutant measure is depicted in figure 5. 428

There are a total of 2,145 unique ICD10 codes included in 429 the graph, identified from SUDF entries and 6 pollutants have 430 been included: PM2.5, PM10, NO_2 , O_3 , SO_2 and CO. The 431 different measurements were acquired from both Atmo and 432 CAMS data, but CO was only available in CAMS data. There 433 are a total of 6 locations included in the knowledge graph. 434 They are assigned to the following INSEE communes: "Le 435 Temple", "Ambes", "Merignac", "Bordeaux", "Bassens", and 436 "Talence". These communes are assigned to INSEE URIS 437

Triples	E	R	Training	Testing
1,922,998	443,332	13	1,634,548	288,450

Table 1: The number of triples, entities, and relations along with the training and testing split of triples for the first PANORAMA knowledge graph.

Figure 5: Triples from the CAMS dataset. Relations, in green, are assigned to nodes according to the measurement of interest



Figure 6: The conceptual knowledge graph obtained after the integration of the different knowledge sources. The green banners indicate the knowledge sources. Nodes are indicated with yellow squares and Literal nodes are indicated with bolded text. Links are indicated with directed black arrows and are labelled according to the relation.

identified from the INSEE RDF space according to the map-ping exemplified earlier.

440 4.2 Predicting the link between pollutants and diseases

Link prediction is the main downstream task used to evaluate 442 knowledge graph embedding models. There are two main 443 types of link prediction that can be conducted: (i) Trans-444 ductive link prediction, in which all links that are predicted 445 are relations already seen in the training graph; (ii) Inductive 446 link prediction, in which the goal is to predict links between 447 unseen entities in the training graph. For the PANORAMA 448 knowledge graph, transductive link prediction has been used 449 to predict the relation of mayHaveParticipant, which is al-450 ready seen in the training graph between ICD10 and Pollutant 451 entities. In this work, the PyKEEN framework⁶ was used to 452 perform the analysis over the knowledge graph. PyKEEN is 453 a Python package for reproducible, facile knowledge graph 454 embeddings. 455

There exist 2,150 triples in the dataset including the may-456 HaveParticipant link between ICD10 and Pollutant entities, 457 which serves as the future link of interest in the KGEM. For a 458 proof of concept, the KGEM was trained on the following hy-459 perparameters: model algorithms TransE, ComplEx, RotatE, 460 epoch size 2, 5, embedding size 50, batch size 128. Simi-461 lar to previous knowledge graph work, test triples are ranked 462 against all candidate triples not appearing in the training set. 463 The knowledge graph dataset was divided into an 85-15 split 464 for training and testing, respectively. This evaluation strategy 465 employs a learning to rank problems, in which we can see 466

how well each positive triple in a test set ranks against syn-467 thetic negatives, built under the same procedure used in the 468 training. This uses information retrieval metrics such as Mean 469 Rank (MR), Mean Reciprocal Rank (MRR), and Hits@K re-470 spectively. MR computes the arithmetic mean over all indi-471 vidual ranks. MRR takes the reciprocal of the rank to guaran-472 tee robustness against outliers. It is more affected by changes 473 of low ranked triples than high ranked triples. Hits@K cal-474 culates how many positive triples in a test set show up in the 475 top k position against synthetic negative triples. This measure 476 describes the fraction of positive facts that appear in the first 477 k triples of the sorted rank list, and thus does not differenti-478 ate between triples that rank higher than k. Table 2 exhibits 479 the model losses and the three main metrics of evaluation for 480 the KGEM: MR, MRR, and Hits@K at three levels - 1, 3, 481 and 10 hits. From this table, we can see that the TransE scor-482 ing function performs better than other model algorithms. An 483 epoch size of 5 also outperformed the competing epoch size 484 of 2 for all metrics in this model. The TransE model with 5 485 epochs achieved the best results with a mean rank of 31.376, 486 mean reciprocal rank of 0.312, and 0.254, 0.335, and 0.524 487 Hits@1, 3, and 10, respectively. 488

5 Discussion and Perspectives

5.1 Contributions of the study

The present study demonstrated the ability to apply knowledge graph technology, which is extremely novel in public health, to sparse air quality data sources and health data. The PANORAMA knowledge graph provides data coherency and a holistic presentation of several knowledge sources. Further, the graph is able to fill gaps that may exist in discrete

489

⁶https://pykeen.github.io/

	Epochs	Losses	MR	MRR	10	Hit@ 3	1
TransE	2	0.00437,0.00206	928.582	0.142	0.217	0.182	0.158
	5	0.00138,0.00106	31.376	0.312	0.524	0.335	0.254
ComplEx	2	0.01091,0.00580	170058.469	0.0197	0.078	0.0026	0.0038
-	5	0.00580,0.00576	165203.078	0.0249	0.0182	0.0078	0.00396
RotatE	2	0.00782,0.00781	157281.492	0.0358	0.00468	0.00156	0.00156
	5	0.00614,0.00568	26690.462	0.0313	0.000377	0.000223	0.000108

Table 2: Metrics of the knowledge graph embedding model with the following parameters: Scoring functions TransE, ComplEx, RotatE, epoch size 2, 5, embedding size 50, and batch size 128. The TransE epoch 5 model produced optimal results for the embedding (highlighted) - a MR of 31.376, MRR of 0.312, and Hits@1, 3, and 10 of 0.254, 0.335, and 0.524 respectively.

knowledge sources, such as missing measurements on par-497 ticular dates or in particular locations, that may be achieved 498 with the integration of other air quality data. This graph can 499 recognize ambiguous entities in context, connect entities to 500 similar data sources, leverage related information in a dataset, 501 and provides a foundation for use in other settings and con-502 texts. Moreover, promising results of the embedding model 503 demonstrate inferential reasoning abilities about health from 504 air quality and pollution. A mean rank of 31.376 and mean 505 reciprocal rank of 0.312 for an initial embedding model us-506 ing a TransE scoring function show that future renditions of 507 the model will likely be able to produce insightful prediction 508 results. A Hits@10 result of 0.524 indicates that there is a 509 proportion of 0.524 of correct positive facts (plausible asso-510 ciation between pollutants and health outcomes) appearing 511 in the top 10 entries of a sorted ranked list again synthetic 512 negatives in the model, demonstrating further aptitude of this 513 knowledge graph's embeddings and link prediction. These 514 hits indicate that over half of results in a sorted ranked list 515 are relevant, and as the model matures with the fine-tuning of 516 hyperparameters, inferential reasoning is probable. 517

518 5.2 Limitations of the study

The study presents some limitations and challenges. Firstly, 519 the ability to create a semantic knowledge graph from a vari-520 ety of heterogeneous sources is entirely reliant on the ability 521 to align and harmonize the discrete data. In this study, it was 522 possible to merge the data sources on location (which acts 523 as a proxy), simply due to the availability of these different 524 sources in the same region. However, this will be a continued 525 526 challenge in applying this knowledge graph without adaptation in other settings, as the presence of harmonizing vari-527 ables is fully dependent on data availability in the context of 528 interest. In a global health setting, in which the PANORAMA 529 knowledge graph may be of important use, such harmonizing 530 variables may not be always present. In the context of LMICs, 531 relying on Open Data could be a relevant alternative for pre-532 diction tasks at the macro level. A major challenge with us-533 ing complex and heterogeneous data is the varying granular-534 ity (particularly spatial coverage of satellite Vs in-situ sen-535 sors based data) and frequency of data variables, particularly 536 those which will be used for harmonization. When only a 537 limited number of measurement locations are available, or 538 539 when the measurement locations are sparse and not spatially representative of the target region, taking into account only 540

neighborhood-level spatial or temporal dependencies may re-541 sult in overfitting [Lin et al., 2020]. Air quality data tends to 542 be available at different specificities, and in the current study, 543 proved to be challenging for integrating Atmo, CAMS, and 544 SUDF data. CAMS data is geolocated according to geograph-545 ical coordinates, but only to the tenths place for each latitude 546 and longitude value. When examining a small region such as 547 Gironde for application of PANORAMA, it proved difficult 548 to align such high-level coordinates from CAMS with Atmo 549 in-situ sensors, which record data at the sensor's precise lo-550 cation. Thus, it is expected that there will be some variation 551 in the locations, which may impede the process of harmo-552 nization. This challenge will likely disappear when applied 553 to larger areas or with data sources that provide more precise 554 location records, such as GPS coordinates. However, it is un-555 realistic to assume that such perfect data will be obtained for 556 other applications; this represents a limitation to be consid-557 ered moving forward. Lastly, ambiguity among data sources 558 for identical variable records presents a challenge. This study 559 was fortunate to have access to health data containing symp-560 tom and disease identifiers based on the internationally rec-561 ognized ICD vocabulary. However, health data is indicated 562 differently across the globe, and this lexical gap may make 563 harmonization difficult. 564

5.3 Future works

The PANORAMA knowledge graph was initially created 566 with the primary purpose of being used in global health 567 contexts where air quality data is sparse and heterogeneous. 568 However, the graph's potential extends beyond this initial ap-569 plication, and there is an interest in testing its capabilities in 570 other settings such as the Global South. Doing so will provide 571 insights into global health contexts and help us understand 572 the graph's limitations with various data sources. One way to 573 enhance the graph's potential is by integrating data collected 574 by citizen scientists or the Internet of Things (IoT) using low-575 cost sensors. It will be relevant to investigate similarly to [Lin 576 et al., 2020] how to exploit both local autocorrelation patterns 577 and global autocorrelation trends of predictions and labeled 578 data in the learned spatio-temporal embedding space for ac-579 curate fine-scale prediction tasks with limited labeled data. 580 Additionally, it would be interesting to explore the dynamic 581 modality of measures of exposure as discussed in Letellier et 582 al. [Letellier et al., 2022], while utilizing geolocated meta-583 data provided by mobile operators' Call Data Records. 584

585 **References**

- [Beckerman *et al.*, 2013] Bernardo S. Beckerman, Michael
 Jerrett, Randall V. Martin, Aaron Van Donkelaar, Zev
 Ross, and Richard T. Burnett. Application of the deletion/substitution/addition algorithm to selecting land use
 regression models for interpolating air pollution measurements in California. *Atmospheric Environment*, 77:172–
- ⁵⁹² 177, October 2013.
- [Bordes *et al.*, 2013] Antoine Bordes, Nicolas Usinier, Alberto Garcia-Dur\'{a}n, Jason Weston, and Oksana
 Yakhnenko. Translating Embeddings for Modeling Multi-
- Relational Data. In *Proceedings of the 26th International*
- 597 Conference on Neural Information Processing Systems -
- Volume 2, pages 2787–2795. Curran Associates Inc., 2013.
- [Ferrer-Cid *et al.*, 2021] Pau Ferrer-Cid, Jose M. BarceloOrdinas, and Jorge Garcia-Vidal. Graph Learning Techniques Using Structured Data for IoT Air Pollution Monitoring Platforms. *IEEE Internet of Things Journal*,
 2021
- 603 8(17):13652–13663, September 2021.
- [Guan et al., 2018] Saiping Guan, Xiaolong Jin, Yuanzhuo
 Wang, and Xueqi Cheng. Shared Embedding Based Neural Networks for Knowledge Graph Completion. In Pro-*ceedings of the 27th ACM International Conference on In- formation and Knowledge Management*, pages 247–256,
 Torino Italy, October 2018. ACM.
- [Hogan et al., 2022] Aidan Hogan, Eva Blomqvist, Michael 610 Cochez, Claudia D'amato, Gerard De Melo, Clau-611 dio Gutierrez, Sabrina Kirrane, José Emilio Labra 612 Gayo, Roberto Navigli, Sebastian Neumaier, Axel-613 Cyrille Ngonga Ngomo, Axel Polleres, Sabbir M. Rashid, 614 Anisa Rula, Lukas Schmelzeisen, Juan Sequeda, Steffen 615 Staab, and Antoine Zimmermann. Knowledge Graphs. 616 ACM Computing Surveys, 54(4):1–37, May 2022. 617
- ⁶¹⁸ [Holloway *et al.*, 2021] Tracey Holloway, Daegan Miller,
 ⁶¹⁹ Susan Anenberg, Minghui Diao, Bryan Duncan, Ar⁶²⁰ lene M. Fiore, Daven K. Henze, Jeremy Hess, Patrick L.
 ⁶²¹ Kinney, Yang Liu, Jessica L. Neu, Susan M. O'Neill,
 ⁶²² M. Talat Odman, R. Bradley Pierce, Armistead G. Russell,
 ⁶²³ Daniel Tong, J. Jason West, and Mark A. Zondlo. Satellite
- Monitoring for Air Quality and Health. *Annual Review of Biomedical Data Science*, 4(1):417–447, July 2021.
- [Huang *et al.*, 2021] Xuqian Huang, Jiuyang Tang, Zhen
 Tan, Weixin Zeng, Ji Wang, and Xiang Zhao. Knowledge graph embedding by relational and entity rotation. *Knowledge-Based Systems*, 229:107310, October 2021.
- [Kampa and Castanas, 2008] Marilena Kampa and Elias
 Castanas. Human health effects of air pollution. *Envi- ronmental Pollution*, 151(2):362–367, January 2008.
- [Kerckhoffs *et al.*, 2019] Jules Kerckhoffs, Gerard Hoek,
 Lützen Portengen, Bert Brunekreef, and Roel C. H. Ver-
- meulen. Performance of Prediction Algorithms for Model-
- ing Outdoor Air Pollution Spatial Surfaces. *Environmental Science & Technology*, 53(3):1413–1421, February 2019.
- [Khojasteh *et al.*, 2021] Davood Namdar Khojasteh, Gho lamreza Goudarzi, Ruhollah Taghizadeh-Mehrjardi, Ak-
- 640 wasi Bonsu Asumadu-Sakyi, and Masoud Fehresti-Sani.

Long-term effects of outdoor air pollution on mortality and morbidity-prediction using nonlinear autoregressive and artificial neural networks models. *Atmospheric Pollution Research*, 12(2):46–56, February 2021. 644

- [Letellier *et al.*, 2022] Noémie Letellier, Steven Zamora, Chad Spoon, Jiue-An Yang, Marion Mortamais, Gabriel Carrasco Escobar, Dorothy D. Sears, Marta M. Jankowska, and Tarik Benmarhnia. Air pollution and metabolic disorders: Dynamic versus static measures of exposure among Hispanics/Latinos and non-Hispanics. *Environmental Research*, 209:112846, June 2022.
- [Li et al., 2020] Linfeng Li, Peng Wang, Jun Yan, Yao Wang,
 Simin Li, Jinpeng Jiang, Zhe Sun, Buzhou Tang, Tsung-Hui Chang, Shenghui Wang, and Yuting Liu. Real-world
 data medical knowledge graph: construction and applications. Artificial Intelligence in Medicine, 103:101817,
 March 2020.
- [Li *et al.*, 2022] Michelle M. Li, Kexin Huang, and Marinka Zitnik. Graph representation learning in biomedicine and healthcare. *Nature Biomedical Engineering*, 6(12):1353–1369, October 2022.
- [Lin et al., 2020] Yijun Lin, Yao-Yi Chiang, Meredith662Franklin, Sandrah P. Eckel, and Jose Luis Ambite. Build-663ing Autocorrelation-Aware Representations for Fine-Scale664Spatiotemporal Prediction. In 2020 IEEE International665Conference on Data Mining (ICDM), pages 352–361, Sor-666rento, Italy, November 2020. IEEE.667
- [Ma et al., 2019] Jun Ma, Yuexiong Ding, Jack C.P. Cheng,
Feifeng Jiang, and Zhiwei Wan. A temporal-spatial in-
terpolation and extrapolation method based on geographic
Long Short-Term Memory neural network for PM2.5.
Journal of Cleaner Production, 237:117729, November
2019.668
670
670
- [Ma et al., 2020] Ling Ma, Cheng-An Zhao, Tingxian Wu,
 Yingying Guo, and Ziru Zhao. The Effects of AIr Quality
 and Weather Conditions on Weather Sensitive Diseases.
 In 2020 International Conference on Computer Informa tion and Big Data Applications (CIBDA), pages 60–65,
 Guiyang, China, April 2020. IEEE.
- [Noy *et al.*, 2009] Natalya F. Noy, Nigam H. Shah, Patricia L. Whetzel, Benjamin Dai, Michael Dorf, Nicholas Griffith, Clement Jonquet, Daniel L. Rubin, Margaret-Anne Storey, Christopher G. Chute, and Mark A. Musen. BioPortal: ontologies and integrated data resources at the click of a mouse. *Nucleic Acids Research*, 37(Web Server issue):W170–173, July 2009.
- [Raheja *et al.*, 2022] Garima Raheja, Kokou Sabi, Hèzouwè
 Sonla, Eric Kokou Gbedjangni, Celeste M. McFarlane,
 Collins Gameli Hodoli, and Daniel M. Westervelt. A Network of Field-Calibrated Low-Cost Sensor Measurements
 of PM _{2.5} in Lomé, Togo, Over One to Two Years. ACS
 Earth and Space Chemistry, 6(4):1011–1021, April 2022.
- [Shaddick *et al.*, 2018] Gavin Shaddick, Matthew L. 693
 Thomas, Heresh Amini, David Broday, Aaron Cohen, 694
 Joseph Frostad, Amelia Green, Sophie Gumy, Yang 695
 Liu, Randall V. Martin, Annette Pruss-Ustun, Daniel 696

- 697 Simpson, Aaron van Donkelaar, and Michael Brauer. Data
- 698 Integration for the Assessment of Population Exposure
- to Ambient Air Pollution for Global Burden of Disease
 Assessment. *Environmental Science & Technology*,
 52(16):9069–9078, August 2018.
- ⁷⁰¹ 52(16):9069–9078, August 2018.
 ⁷⁰² [Some *et al.*, 2019] Borlli Michel Jonas Some, Georgeta
- 703 Bordea, Frantz Thiessard, Stefan Schulz, and Gayo Di-
- allo. Design Considerations for a Knowledge Graph: The
- WATRIMed Use Case. Studies in Health Technology and Informatics, 259:59–64, 2019.
- [Somé *et al.*, 2019] Borlli Michel Jonas Somé, Georgeta
 Bordea, Frantz Thiessard, and Gayo Diallo. Enabling
- 709 West African Herbal-Based Traditional Medicine Digi-
- tizing: The WATRIMed Knowledge Graph. Studies in
- *Health Technology and Informatics*, 264:1548–1549, August 2019.
- 713 [Trouillon *et al.*, 2016] Théo Trouillon, Johannes Welbl, Se-714 bastian Riedel, Eric Gaussier, and Guillaume Bouchard.
- bastian Riedel, Eric Gaussier, and Guillaume Bouchard.
 In Proceedings of The 33rd International Conference on
- 716 Machine Learning, PMLR, pages 2071–2080, 2016.
- [Wen *et al.*, 2019] Congcong Wen, Shufu Liu, Xiaojing Yao,
 Ling Peng, Xiang Li, Yuan Hu, and Tianhe Chi. A novel
 spatiotemporal convolutional long short-term neural network for air pollution prediction. *Science of The Total En-*
- work for air pollution prediction. Science of The Total
 vironment, 654:1091–1099, March 2019.
- [Zhao *et al.*, 2019] Jiachen Zhao, Fang Deng, Yeyun Cai,
 and Jie Chen. Long short-term memory Fully connected
 (LSTM-FC) neural network for PM2.5 concentration prediction. *Chemosphere*, 220:486–492, April 2019.
- ⁷²⁶ [Zhou *et al.*, 2019] Yanlai Zhou, Fi-John Chang, Li-Chiu
- ⁷²⁶ [Zhou *et al.*, 2019] Yanlai Zhou, Fi-John Chang, Li-Chiu
 ⁷²⁷ Chang, I-Feng Kao, and Yi-Shin Wang. Explore a deep
- learning multi-output neural network for regional multi-
- step-ahead air quality forecasts. *Journal of Cleaner Pro-*
- 730 *duction*, 209:134–145, February 2019.