

Coronavirus Knowledge Graph: A Case Study

Chongyan Chen
University of Texas at Austin
Austin, Texas
chongyanchen_hci@utexas.edu

Islam Akef Ebeid
University of Texas at Austin
iaebeid@utexas.edu

Yi Bu
Peking University
Beijing, China
buyi@pku.edu.cn

Ying Ding
University of Texas at Austin
ying.ding@ischool.utexas.edu

ABSTRACT

The emergence of the novel COVID-19 pandemic has had a significant impact on global healthcare and the economy over the past few months. The virus's rapid widespread has led to a proliferation in biomedical research addressing the pandemic and its related topics. One of the essential Knowledge Discovery tools that could help the biomedical research community understand and eventually find a cure for COVID-19 are Knowledge Graphs. The COVID-19 dataset is a collection of publicly available full-text research articles that have been recently published on COVID-19 and coronavirus topics. Here, we use several Machine Learning, Deep Learning, and Knowledge Graph construction and mining techniques to formalize and extract insights from the PubMed dataset presented in [8] and the COVID-19 dataset [1] to identify COVID-19 related experts and bio-entities. Besides, we suggest possible techniques to predict related diseases, drug candidates, gene, gene mutations, and related compounds as part of a systematic effort to apply Knowledge Discovery methods to help biomedical researchers tackle the pandemic.

CCS CONCEPTS

• **Applied computing** → **Health informatics; Biological networks; • Computing methodologies** → **Information extraction; Knowledge representation and reasoning; • Information systems** → Information systems applications;

KEYWORDS

corona virus, named entity Recognition, BioBERT, knowledge Graph, drug discovery

ACM Reference Format:

Chongyan Chen, Islam Akef Ebeid, Yi Bu, and Ying Ding. 2020. Coronavirus Knowledge Graph: A Case Study. In *KDD2020: ACM Knowledge Discovery in Databases, August 23–27, 2020, San Diego, California*. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '2020, August 23–27, 2020, San Diego, California

© 2020 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

1.1 Knowledge Graphs

A Knowledge Graph (KG) is a graph-based data structure used to represent unstructured information so that a machine can read it. Emerging from Knowledge Bases, KGs now represent a ubiquitous set of methods for representing and integrating knowledge in various domains. A KG contains descriptions of entities and their relationships as information in the form of first-order logical facts such as <Mount Fuji is located in Japan> that can be retrieved and queried heuristically. KGs emerged to power what was known in the eighties and the nineties as Expert Systems - an early form of Artificial Intelligence and Decision Support Systems [25]. The number of entity types and relationships in a KG is finite and is usually but not necessarily organized in a schema or an ontology. In 2012 Google introduced the Google Knowledge Graph [21], a technology that converts multiple information sources to a graph structure where the nodes represent real-life entities and types. The edges represent the relationships between those entities and types, and the technology was aimed at enhancing the users' search experience through predicting the users' search intent and introducing a Knowledge Panel on the right of page [21]. Knowledge Graph technology today has been adopted in many domains and fields to store, integrate, and represent unstructured information in a structured format that is more flexible and machine-readable than the traditional entity-relationship data model [6].

Other examples of widely used and adopted KGs in different domains such as Social Networks and Life Sciences include the Facebook Social Graph [30] and Chem2Bio2RDF [5]. The advancement of KG construction and mining was powered by the already established research fields of Machine Learning, Deep Learning, Graph Mining, and Complex Networks. The techniques and methods developed by researchers in such fields are used to mine data and information in KGs to extract insights crucial to advancing knowledge in various domains. Research in Information Networks also played a role in the construction and mining of KGs mainly through relying on statistical methods and machine learning techniques [28]. Information networks are widely heterogeneous graphs of nodes and edges representing meta-information about a published corpus of literature such as authors, papers, publications, and venues. Hence, information networks are KGs in which graph mining techniques can be applied to extract insights about author collaboration patterns and their topics of interest. Mining information networks as KGs has to lead to understanding trends such as

collaboration patterns and potential drug re-purposing opportunities in a specific domain without reading the entire literature in a field.

KGs have been curated manually, yet, over the years, KG construction techniques have changed. For example, Cyc [18] is a KG that was manually curated, while Freebase [3] and Wikidata [32] were crowd-sourced. KGs can also be extracted using Natural Language Processing (NLP) techniques such as in DBpedia [17] and YAGO [27]. Alternatively, KGs can be constructed using a combination of manual curation and automatic extraction like in NELL [4] and Knowledge Vault [11]. Regardless of the approach of how a KG has been constructed, KGs need to be queried and mined to map complex real-world phenomena and eventually be exploited to solve important research questions. For example, Facebook’s Social Graph needs to be mined to suggest new friends for users. Life Sciences KGs like Chem2Bio2RDF need to be mined to answer research questions related to biomedical science.

1.2 Named Entity Recognition in Knowledge Graph Construction

The move towards natural language understanding through semantic technologies has gained much ground in the past decade, promoting Named Entity Recognition (NER) to a central NLP task. NER has been crucial for building and constructing KGs as the primary method of extracting entities and possibly relations from free text. Also, tasks such as link prediction, relation extraction, and graph completion on KGs are aided by NER. NER can be impactful when applied to mine domain-specific scientific literature such as the biomedical literature to extract bio entities aiding in constructing KGs and advancing downstream knowledge discovery tasks in biomedicine.

Although research in NER has been advancing since the nineties [20], early efforts in domain-specific biomedical NER came later in the early 2000s [26]. Those methods in biomedical NER relied on feature engineering and graphical models such as Hidden Markov Models (HMM) and Conditional Random Fields (CRF) [26]. When applying CRF models to the biomedical text, the objective is to construct a chain out of the words then predict the assigned labels based on a conditionally trained finite state machine where the probability of each label assigned to a word is correlated with a feature set. The objective was then to maximize the log-likelihood of the label given the word directly. The accuracy of the recognition of bioentities in CRF and HMM models were quite low when compared to state of the art today. The current state of the art relies on the latest in Deep Learning in contextual embedding such as BERT. BERT is a deep learning model developed in [9] by a team at Google to be fine-tuned for machine translation tasks. The model was based on the transformer architecture described in [31]. Multiple attention heads are used to train a contextual embedding where the task is to predict masked words of the input sentences. The sophisticated inner architecture of BERT based on multiple encoder-decoder layers allows for learning high-quality embedding from a large corpus of data where the learned weights can be later transferred and fine-tuned to downstream tasks. In [16], the authors trained a BERT model on the corpus of PubMed and PMC named BioBERT. The result was a biomedical contextual embedding model

that was later fine-tuned and used in a biomedical NER task producing high accuracy tagging and extraction of bio entities such as drugs, diseases, and genes. The high accuracy of the BioBERT model allowed and aided in the construction of the PubMed KG presented in [33].

1.3 The COVID-19 Knowledge Graph

Several months amid the emergence of the acute respiratory syndrome COVID-19 caused by the novel coronavirus Sars-CoV-2 in China, the disease has risen to a global pandemic level affecting almost every country on earth and infecting more than 6 million people across the globe and killing more than 350000 [2]. As a result, researchers from every domain have reoriented their efforts towards finding ways and solutions to tackle the pandemic. Specifically, the biomedical literature on COVID-19 and SARS-CoV-2 and other related acute respiratory syndromes that have reached an epidemic level such as SARS and MERS have increased exponentially since the virus’s appearance back in December 2019. As a result, government-backed calls and research institutions like the Allen Institute for AI have released a COVID-19 Open Research Dataset CORD-19 [1]. The dataset contains over 65000 full-text scholarly articles related to COVID-19, SARS-CoV-2, and other related topics. This effort aims to encourage the NLP and KG researcher community to mine the dataset to generate insights through text mining techniques and methods to help point biomedical researchers in the right direction to fight against the virus.

We see the release of the CORD-19 dataset of machine-readable scientific literature as an opportunity to extract a comprehensive and cohesive COVID-19 KG of the entities and relationships though cooccurrence within the corpus of articles. The extracted KG will help understand the relationships between the diseases, the genes, the viruses, and the cures involved in and related to COVID-19 so that future graph and network mining efforts can be applied to extract insights from the dataset. Here we present our vision in contributing to that effort.

We demonstrate methods of entity extraction and KG building to harvest a COVID KG capable of being a useful dataset for future mining in the hope that it will help biomedical researchers find a cure and tackle the pandemic through generating deep insights. We first introduce how to use BioBERT for named entity recognition in the PubMed and CORD-19 datasets. Then we built several Coronavirus Knowledge Graphs based on two different kinds of measurements. One measure the relationship between source node and each target node based on co-occurrence frequency. The other is to use Cosine Similarity to measure the similarity between the source node and each target node.

2 RELATED WORK

Previous efforts and trials to build a comprehensive COVID-19 KG have lacked in several areas. For example [10] built a COVID-19 related KG based on 145 articles and provided a web application for ease of use and access. This COVID-19 KG contains 3954 nodes and 9484 relations, covering ten entity types. It reveals host-pathogen interactions, comorbidities, symptoms, and discovered over 300 candidate drugs for COVID-19. Nevertheless, the effort was limited in terms of the number of publications included in constructing the

KG. [12] applied a machine learning model (BERE [14]) to integrate and mine KG to also aid in the effort of identifying candidate drugs for COVID-19. Besides, [24] used a pre-built KG for COVID-19 drug discovery and identified the drug "baricitinib" to protect lung cells from being infected by the virus. Previously mentioned efforts though promising, yet they lacked the large scale KG construction and mining approaches necessary to extract more profound and in-depth insights about the disease and possible cures, treatments, and genetic influences.

NLP techniques have also been utilized outside of the KG construction arena, for example, [29] introduced CovidQA, a question answering dataset, which comprises 124 questions and answers of triples built by hand from knowledge collected from the CORD-19 dataset. [13] developed a self-supervised context-aware COVID-19 document exploration based on BERT. [19] used BERT to analyze a large collection of COVID-19 literature from the CORD-19 dataset [15] to extract COVID-19 related radiological findings. Though rigorous in using large datasets such as CORD-19, the previous NLP techniques were limited in terms of applications and the impact of those applications on the COVID-19 oriented biomedical research field.

3 DATASETS

3.1 PubMed dataset

The PubMed database contains more than 30 million citations within the various fields of life sciences. The PubMed citation database archived by the MEDLINE archive has always been the desired datasets for biomedical text and graph mining research communities.

We select PubMed dataset because it is a popular dataset in biomedical area and reflect general biomedical knowledge.

3.2 The PubMed Knowledge Graph

[33] built a PubMed KG which connects disambiguated author names, their articles, and bio-entities using the PubMed database were they parsed 29 million PubMed abstracts from 1781 till 2019. In addition to funding extracted from the National Institutes of Health using ExPORTER, and affiliations were extracted from ORCID and MapAffil.

3.3 The CORD-19 Dataset

The CORD-19 dataset was released in response to COVID-19, where the US Government has issued requests for research groups and institutions to combine efforts to release the COVID-19 Open Research Dataset (CORD-19). The datasets contain more than 135000 articles with over 68000 full texts on topics related to Coronavirus and the COVID-19 pandemic. The data set was released to help the biomedical research community by applying the latest in NLP to extract deep insights and understandings of the pandemic patterns and the possible drugs, cures, and genes that might be involved and identified [1].

Here we perform our analysis on the entities and relationships extracted from the three datasets and we show the potential in knowledge discovery.

4 EXPERIMENT-0 IDENTIFY EXPERTS ON CORONAVIRUS TOPICS

We would like to identify the experts for COVID to encourage collaboration. To do that, we analyzed COVID-19 44k dataset and ranked the researchers according to the number of articles they published in the COVID-19 44K dataset. Part of the results are shown in Table 1.

Table 1: COVID-19 related researchers

Author	# of articles published in COVID-19 dataset
Perlman, Stanley	142
Drosten, Christian	137
Yuen, Kwok-Yung	136
Baric, Ralph S	132
Jiang, Shibo	120
Enjuanes, Luis	116
Snijder, Eric J	104
Weiss, Susan R	92

5 EXPERIMENT-1 NAMED ENTITY RECOGNITION WITH BIOBERT

5.1 Model

BERT (Bidirectional Encoder Representations from Transformers) [9] is a highly influential Natural Language Processing model that proposed back in 2018. BERT was inspired by many advanced Deep Learning models, such as semi-supervised sequence learning[7], ELMo [23] and the Transformer architecture [31].

The input representation of BERT is the sum of a token embedding using WordPiece, a segmentation embedding indicating whether each token belongs to sentence A or sentence B, and a position embedding. A [CLS] flag is added before the first word of the sentence, and a [SEP] flag is added as a separator token.

BERT has two tasks for pre-training: Masked Language Model task and Next Sentence Prediction task. Considering most of traditional NLP model, instead of training a left-to-right or right-to-left model based on the input language, it is better to use the bidirectional model. However, the bidirectional model is not suitable for the conditional task. Thus, inspired by the Cloze task, a masked language model is adapted as the first task for BERT pre-training. The second task for BERT pre-training is Next Sentence Prediction (NSP), which allows the model to understand sentence relationships.

BioBERT (Bidirectional Encoder Representations from Transformers for Biomedical Text Mining) [16] is a biomedical language representation model based on BERT [9]. It is proposed because directly adapting BERT to text mining in the biomedical area was not promising, given the word shift from generic domain to the biomedical domain. BioBERT is pre-trained on PubMed abstracts and PubMed Central full-text articles (PMC).

We select BioBERT-Base v1.1 (+PubMed 1M) based on the BERT-base-Cases model. For our fine-tuning section, we fine-tuned it on the NCBI disease dataset. The input to the BioBERT model is a

Table 2: Named Entity Recognition using BioBERT in PubMed dataset

Word	Label
Asymptomatic	O-MISC
carrier	O-MISC
state	O-MISC
,	O-MISC
acute	B-MISC
respiratory	I-MISC
disease	I-MISC
,	O-MISC
and	O-MISC
pneumonia	B-MISC
due	O-MISC
to	O-MISC
severe	O-MISC
acute	B-MISC
respiratory	I-MISC
syndrome	I-MISC
coronavirus	I-MISC
2	I-MISC
(O-MISC
SARSCoV	O-MISC
-	O-MISC
2	O-MISC
)	O-MISC
:	O-MISC
Facts	O-MISC
and	O-MISC
myths	O-MISC

sentence embedded following BERT's embedding process. Parts of the token-level evaluation looks like: "[CLS]Ang ##iot ##ens ##in - converting enzyme 2 (AC ##E ##2) as a SA ##RS - Co ##V - 2 receptor: molecular mechanisms and potential therapeutic target. [SEP]". The output of this model will be a sentence with labels. Label "B-MISC" means the Beginning of bioentities, "I-MISC" means Inside the bio entities, and "O-MISC" means Outside the bio entity. We tested BioBERT on PubMed KG and the CORD-19 dataset.

5.2 Results

Examples of entity-level recognized names from the PubMed dataset are shown in Table 2. The recognized bio-entities are "acute respiratory disease", "pneumonia", and "acute respiratory syndrome coronavirus 2". "SARSCOV-2" should, but is not recognized as bio-entities.

Since we do not have detailed labels for BioBERT fine-tuning training and thus cannot predict detailed labels directly from BioBERT. To get more detailed labels, we trained a Random Forest model on around 100,000 PubMed bio-entities labeled with five categories: species, gene, diseases, drug, gene mutation, and tested on the bio-entities recognized by BioBERT. The F1-score is shown in Table 3. The F1-Score of disease, gene, and drug recognition are all over 75%. However, the model poorly predicts when it comes to label

Table 3: Specific named entity classification using Random Forest

Label	Precision	Recall	F1-Score	Data Size
Species	0.72	0.31	0.43	15519
Disease	0.94	0.61	0.74	12077
Gene	0.95	0.64	0.76	18678
Drug	0.66	0.99	0.79	53523
Gene mutation	0.5	0.17	0.25	36

Table 4: named entity Recognition using BioBERT in COVID-19 44K dataset

Word	Label
Thrombocytopenia	O-MISC
is	O-MISC
associated	O-MISC
with	O-MISC
severe	O-MISC
coronavirus	B-MISC
disease	I-MISC
2019	I-MISC
COVID	B-MISC
-	I-MISC
19	I-MISC
infections	O-MISC
A	O-MISC
meta	O-MISC
-	O-MISC
analysis	O-MISC

"Species" and labels "Gene mutation," probably because the dataset was very unbalanced (Gene mutation only has 17 samples, and Species has only 1395 samples).

Examples of named entity recognized from CORD-19 44K dataset are shown in Table 4. In the example, the recognized bio-entities are "coronavirus disease 2019", and "COVID-19". "Thrombocytopenia," a kind of disease, should, but is not recognized. From these cases, we find that BioBERT does not perform greatly, despite its high accuracy. It may be because BioBERT can easily recognize the easy and common bio-entities with a high occurrence rate but fail to recognize rare biomedical terms.

6 EXPERIMENT-2.1 CO-OCCURRENCE FREQUENCY BASED KNOWLEDGE GRAPH

6.1 Method

We used Gephi to build the co-occurrence frequency based Knowledge Graph. Co-occurrence frequency is an above-chance frequency of occurrence of two entities from an article. The data is from the PubMed Knowledge Graph. For each target node (related bio-entities), we calculated the times it shows up with the source node and treat the times(co-occurrence frequency) as target node's weight. The higher the co-occurrence frequency, the closer the target node is to the source node.

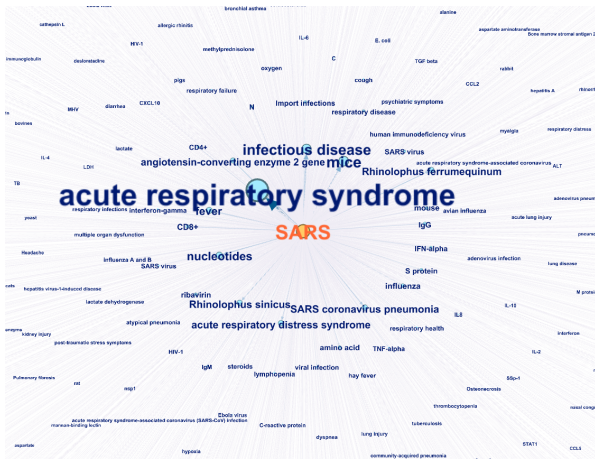


Figure 3: SARS centered KG based on co-occurrence frequency

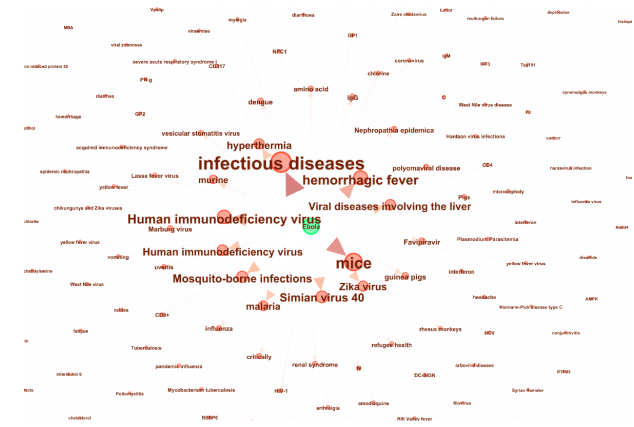


Figure 5: Ebola centered KG based on co-occurrence frequency

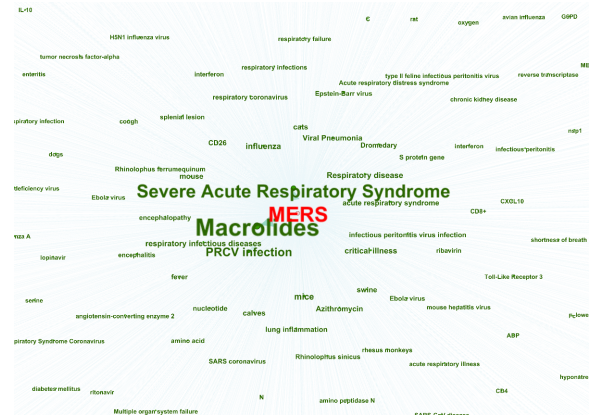


Figure 4: MERS centered KG based on co-occurrence frequency

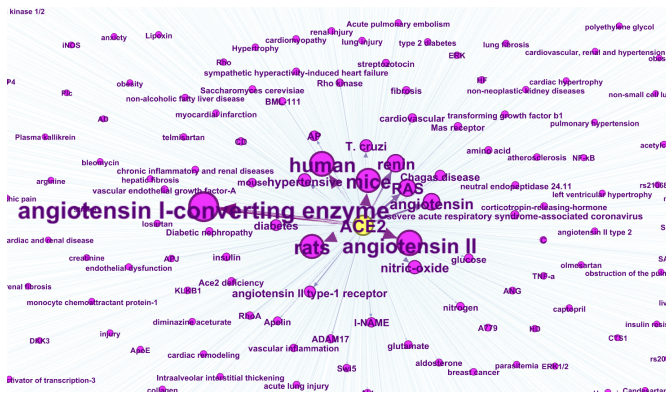


Figure 6: ACE2 centered KG based on co-occurrence frequency

Figure 5 shows that Ebola’s highly related diseases are hemorrhagic fever, hyperthermia, malaria, and mosquito-borne infections. MERS’s highly related drugs are favipiravir and amodiaquine. MERS’s highly related genes/chemicals are CD8+, CD4, DC-SIGN, CD317, GP2, IFN-g, IRF3, RBBP6, and etc.

Figure 6 shows Angiotensin-converting enzyme 2 (ACE2) centered Knowledge Graph. ACE2 is an enzyme, which lowers blood pressure by catalysing the hydrolysis of angiotensin II into angiotensin (1-7). ACE2 is the receptor that COVID-19 uses to infect lung cells. It also serves as receptor for other coronaviruses such as HCoV-NL63, SARS-CoV. As shown, ACE2’s related genes/chemicals are renin, RAS, angiotensin, insulin, Mas receptor, vascular endothelial growth factor-A, and etc. ACE2’s related diseases are diabetes, hypertensive, chagas disease, severe acute respiratory syndrome-associated coronavirus. ACE2’s related drugs are streptozotocin, nitric-oxide, and aldosterone .ACE2’s related gene mutations are “rs2106809” and “rs2074192”.

From the results we believe the PubMed Knowledge Graph is very promising. However, this kind of KG has a entity name disambiguation issue. For example, “Favipiravir” could also be shown as “favipiravir”. Another case is “ACE-2”, which is the abbreviation of Angiotensin-converting enzyme 2. Besides, the co-occurrence frequency cannot reflect the relationship between the source node and the target node well. For example, if “A has nothing to do with B” mentioned lots of times in different documents, its co-occurrence frequency will be very high.

7 EXPERIMENT-2.2 COSINE SIMILARITY BASED KNOWLEDGE GRAPH

7.1 Method

To deal with problems with co-occurrence frequency based KG, we first normalized the entity using some human designed rules to deal with entity name disambiguation issue. We mainly focus on case sensitive, singular and plural, and disambiguation. For example, “SIAsNN” will be normalized as “siann” and “respiratory illnesses” will be normalized as “respiratory_illness”. Then we used Word2Vec to convert the normalized entity to vector with length of 100. We

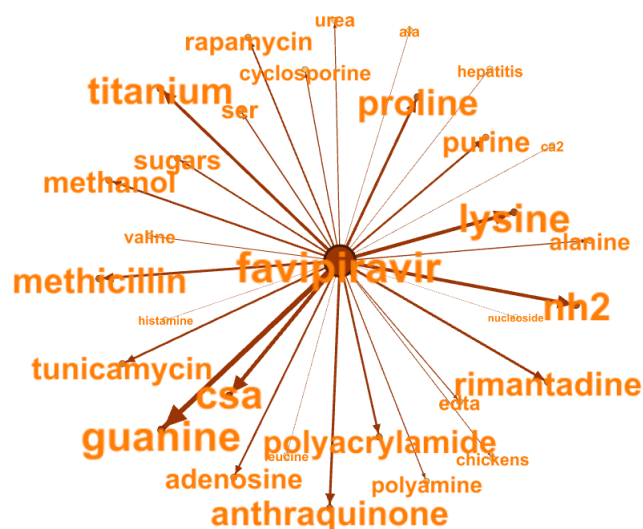


Figure 7: favipiravir-centered KG (related chemical) based on cosine similarity

then use Cosine Similarity to measure the similarity between the source node and each target node. The Cosine Similarity is defined as follows:

$$\cos(S, T) = \frac{ST}{\|S\| \|T\|} = \frac{\sum_{i=1}^n S_i T_i}{\sqrt{\sum_{i=1}^n (S_i)^2} \sqrt{\sum_{i=1}^n (T_i)^2}} \quad (1)$$

The KG based on cosine similarity is also built using Gephi software.

7.2 Results

Figure 7 shows part of the favipiravir-centered knowledge graph (chemical related). The source node is favipiravir and the target node are related chemicals. The edge is cosine similarity relations. The closer the target node is to the source node, the similar the target node is to the source node. As shown, the top 10 chemical related to favipiravir are guanine, csa, lysine, nh2, titanium, proline, methicillin, anthraquinone, rimantadine, polyacrylamide. Figure 8 shows part of the favipiravir-centered knowledge graph (gene related). As shown, the top 10 gene related to favipiravir are gm_csf, abortion, rig_i, isg15, csa, akt, mtor, p53, th1, p38, tgf_beta.

We also generate other 5 drug-centered KGs based on cosine similarity. The top 10 chemicals related to lopinavir are retinoic acid, nucleoside, tyr, glutamine, ribavirin, glycyrrhizin, co2, lopinavir, phosphonate, lymphoma, ifitm3. The top 10 genes related to lopinavir are neuraminidase, p53, eif2alpha, apod, ribavirin, infection, cox-2, ifitm3, iron.

The top 10 chemicals related to ribavirin are glycyrrhizin, lactate, corticosteroid, coronavirus, steroid, nucleoside, ribavirin, sodium, glucose, infection, oxygen, calcium, obesity. The top 10 genes related to ribavirin are toxicity, p53, swine, fibrosis, iron, neuraminidase, diabetes, ribavirin, anemia, inflammation, infection.

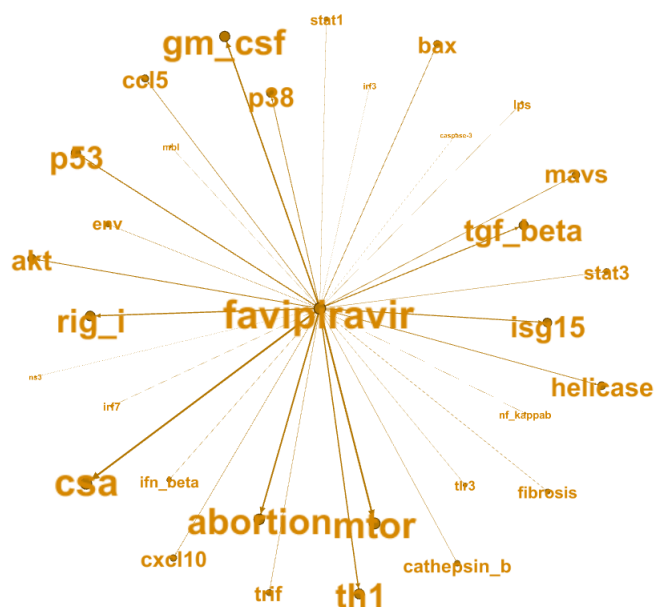


Figure 8: favipiravir-centered KG (related gene) based on cosine similarity

The top 10 chemicals related to ritonavir are atp, ritonavir, cyclophosphamide, mtt, toxicity, sialic_acid, sds, encephalitis, superoxide, sucrose, ethanol. The top 10 genes related to ritonavir are jnk, p53, rig_i, encephalitis, rnase_l, stat3, toxicity, akt, neuraminidase, stat1.

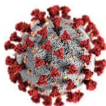
The top 10 chemicals related to tamiflu are superoxide, prednisolone, flavonol, proline, nitric_oxide, thymidine, glycyrrhizin, propidium_iodide, nitrogen, aspirin, tamiflu. The top 10 genes related to tamiflu are il-10, cd44, eif2alpha, tgf_beta1, thr2, ifn, cxcl10, tumor_necrosis_factor_tnf)-alpha, ire1, ccl2, tbk1.

The top 10 chemicals related to umifenovir are: tacrolimus, alkyl, carbon_monoxide, ca(2, cd, nucleolin, cytosine, glycyrrhizic_acid, 2'-o, umifenovir, prostaglandin_e2. The top 10 genes related to umifenovir are parp, pd_l1, monocyte_chemoattractant_protein-1, nef, cxcr4, cd45, nucleolin, dc_ign, annexin_v, cd19, mmp-2.

8 DISCUSSION AND CONCLUSION

In this research, we first used BioBERT to recognize entities in the PubMed and the CORD-19 dataset. Our results show that most of the recognized entities are strictly biomedical. Most of the recognized entities in the CORD-19 dataset are disease lacking diversity in entity types due to a lack in finding a suitable bio-medical training dataset with detailed labeled bio-entity. For future work, we will explore more other bio-medical dataset and try other biomedical NLP models for named entity recognition, e.g., blueBERT [22].

Furthermore, we introduced the construction of Coronavirus Knowledge Graph based on two different methods: co-occurrence frequency and cosine similarity. We explored and revealed that the drug candidates recommended by drug-centered KG are promising. We will consult experts in COVID-related research to verify our



COVID-19

Type: Disease
Description: In COVID-19, 'CO' stands for 'corona,' 'VI' for 'virus,' and 'D' for disease. Formerly, this disease was referred to as "2019 novel coronavirus" or "2019-nCoV". There are many types of human coronaviruses including some that commonly cause mild upper-respiratory tract illnesses.

Bio Entities/Related	Topic distribution	Experts	Organization	Featured publications
Drugs: Remdesivir, Ritonavir, Ribavirin	clinical characterization	remuzzi g*, Istituto di Ricerche Farmacologiche Mario Negri [LINK]	World Health Organization (WHO)	Presumed asymptomatic carrier transmission of COVID-19
Gene: MHC-I, MHC-II	pathogenesis research	Andrea Remuzzi, University of Bergamo [LINK]	Recon	Pathological findings of COVID-19 associated with acute respiratory distress syndrome
Protein: ACE2 (Angiotensin-Converting Enzyme 2); Spike Protein Receptor	therapeutics research	Lauren Gardner, Johns Hopkins University, [LINK]	CDC	Wu P et al bioRxiv 2020
Species: Bats	epidemiological study	Annelies Wilder-Smith, London School of Hygiene and Tropical Medicine, [LINK]	Worldometer	
	virus transmission	Karolina George, Professor of Pharmacology, Medical School, Aristotle University, [LINK]		
	vaccines research			
	virus diagnostics, and viral genomics			

Figure 9: auto drug profiling example

drug-centered KG and conduct a more in-depth analysis for future work.

Also, we aim to build a wider COVID related KG, connecting all COVID related bio-entities rather than small drug/disease-centered KG. The extracted KG will help understand the relationships between the diseases, the genes, the viruses, and the cures involved in and related to COVID-19.

Finally, we hope to build an automatic profiling system to generate expert, drug, or disease profiling. The expected disease profiling will look like Figure 9, which includes description, related bio entities (drugs, gene, protein, species), topic distribution, related experts, organization, and featured publications.

9 CONTRIBUTIONS

Y.D. and Y.B. proposed the idea and supervised the project. C.C. wrote the paper. I.A.E wrote the Introduction and revised this paper. C.C. conducted the named entity recognition and Knowledge Graph building. I.A.E conducted the Word2Vec for Experiment 2.2.

ACKNOWLEDGMENTS

We would like express our gratitude to Prof. Jaewoo Kange's DMIS Lab team for pretraining BioBERT, Vinay Locharulu for suggestion and support, Prof. Jian Xu for providing PubMed Knowledge Graph, and Yifei Wu for conducting entity normalization.

REFERENCES

- [1] [n.d.]. COVID-19 Open Research Dataset Challenge (CORD-19). <https://kaggle.com/allen-institute-for-ai/CORD-19-research-challenge>
- [2] [n.d.]. Johns Hopkins Coronavirus Resource Center. <https://coronavirus.jhu.edu/>
- [3] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. 1247–1250.
- [4] Andrew Carlson, Justin Betteglieri, Bryan Kisiel, Burr Settles, Estevam R Hruschka, and Tom M Mitchell. 2010. Toward an architecture for never-ending language learning. In *Twenty-Fourth AAAI Conference on Artificial Intelligence*.
- [5] Bin Chen, Xiao Dong, Dazhi Jiao, Huijun Wang, Qian Zhu, Ying Ding, and David J Wild. 2010. Chem2Bio2RDF: a semantic framework for linking and data mining chemogenomic and systems chemical biology data. *BMC bioinformatics* 11, 1 (2010), 255.
- [6] Peter Pin-Shan Chen. 1976. The entity-relationship model—toward a unified view of data. *ACM transactions on database systems (TODS)* 1, 1 (1976), 9–36.
- [7] Andrew M Dai and Quoc V Le. 2015. Semi-supervised sequence learning. In *Advances in neural information processing systems*. 3079–3087.
- [8] Franck Dernoncourt and Ji Young Lee. 2017. Pubmed 200k rct: a dataset for sequential sentence classification in medical abstracts. *arXiv preprint arXiv:1710.06071* (2017).
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [10] Daniel Domingo-Fernandez, Shounak Bakshi, Bruce Schultz, Yojana Gadiya, Reagon Karki, Tamara Raschka, Christian Ebeling, Martin Hofmann-Apitius, et al. 2020. COVID-19 Knowledge Graph: a computable, multi-modal, cause-and-effect knowledge model of COVID-19 pathophysiology. *BioRxiv* (2020).
- [11] Xin Dong, Evgeniy Gabrilovich, Jeremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmann, Shaohua Sun, and Wei Zhang. 2014. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 601–610.
- [12] Yiyue Ge, Tingzhong Tian, Sulin Huang, Fangping Wan, Jingxin Li, Shuya Li, Hui Yang, Lixiang Hong, Nian Wu, Enming Yuan, et al. 2020. A data-driven drug repositioning framework discovered a potential therapeutic agent targeting COVID-19. *bioRxiv* (2020).
- [13] Dusan Grucicic, Gorjan Radevski, Tinne Tuytelaars, and Matthew B Blaschko. 2020. Self-supervised context-aware Covid-19 document exploration through atlas grounding. (2020).
- [14] Lixiang Hong, Jinjian Lin, Jiang Tao, and Jianyang Zeng. 2019. BERE: An accurate distantly supervised biomedical entity relation extraction network. *arXiv preprint arXiv:1906.06916* (2019).
- [15] Scite Inc. 2020. *CORD-19_scite_citation_tallies+contexts*. <https://doi.org/10.5281/zenodo.3724818>
- [16] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36, 4 (2020), 1234–1240.
- [17] Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. DBpedia—a large-scale, multilingual knowledge base extracted from Wikipedia. *Semantic Web* 6, 2 (2015), 167–195.
- [18] Douglas B Lenat. 1995. CYC: A large-scale investment in knowledge infrastructure. *Commun. ACM* 38, 11 (1995), 33–38.
- [19] Yuxiao Liang and Pengtao Xie. 2020. Identifying Radiological Findings Related to COVID-19 from Medical Literature. *arXiv preprint arXiv:2004.01862* (2020).
- [20] David Nadeau and Satoshi Sekine. 2007. A survey of named entity recognition and classification. *Linguisticae Investigationes* 30, 1 (2007), 3–26.
- [21] Heiko Paulheim. 2017. Knowledge graph refinement: A survey of approaches and evaluation methods. *Semantic web* 8, 3 (2017), 489–508.
- [22] Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. Transfer learning in biomedical natural language processing: An evaluation of BERT and ELMo on ten benchmarking datasets. *arXiv preprint arXiv:1906.05474* (2019).
- [23] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365* (2018).
- [24] Peter Richardson, Ivan Griffin, Catherine Tucker, Dan Smith, Olly Oechsle, Anne Phelan, and Justin Stebbing. 2020. Baricitinib as potential treatment for 2019-nCoV acute respiratory disease. *Lancet (London, England)* 395, 10223 (2020), e30.
- [25] Stuart J Russell and Peter Norvig. 2016. Artificial intelligence: a modern approach. Malaysia.
- [26] Burr Settles. 2004. Biomedical named entity recognition using conditional random fields and rich feature sets. In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications (NLPBA/BioNLP)*. 107–110.
- [27] Fabian M Suchanek, Gjergji Kasneci, and Gerhard Weikum. 2007. Yago: a core of semantic knowledge. In *Proceedings of the 16th international conference on World Wide Web*. 697–706.
- [28] Yizhou Sun and Jiawei Han. 2012. Mining heterogeneous information networks: principles and methodologies. *Synthesis Lectures on Data Mining and Knowledge Discovery* 3, 2 (2012), 1–159.
- [29] Raphael Tang, Rodrigo Nogueira, Edwin Zhang, Nikhil Gupta, Phuong Cam, Kyunghyun Cho, and Jimmy Lin. 2020. Rapidly Bootstrapping a Question Answering Dataset for COVID-19. *arXiv preprint arXiv:2004.11339* (2020).
- [30] Johan Ugander, Brian Karrer, Lars Backstrom, and Cameron Marlow. 2011. The anatomy of the facebook social graph. *arXiv preprint arXiv:1111.4503* (2011).
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [32] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* 57, 10 (2014), 78–85.
- [33] Jian Xu, Sunkyu Kim, Min Song, Minbyul Jeong, Donghyeon Kim, Jaewoo Kang, Justin F Rousseau, Xin Li, Weijia Xu, Vette I Torvik, et al. 2020. Building a PubMed knowledge graph. *arXiv preprint arXiv:2005.04308* (2020).