

RESEARCH

# Ontology based mining of pathogen-disease associations from literature

Şenay Kafkas<sup>1,2\*</sup> and Robert Hoehndorf<sup>1,2</sup>

Correspondence:

benay.kafkas@kaust.edu.sa

Computational Bioscience

Research Center, King Abdullah

University of Science and

Technology, 23955-6900 Thuwal,

Saudi Arabia

Computer, Electrical and

Mathematical Sciences and

Engineering Division, King

Abdullah University of Science and

Technology, 23955-6900 Thuwal,

Saudi Arabia

Full list of author information is

available at the end of the article

## Abstract

## Background

Infectious diseases claim millions of lives especially in the developing countries each year, and resistance to drugs is an emerging threat worldwide. Identification of causative pathogens accurately and rapidly plays a key role in the success of treatment. To support infectious disease research and mechanisms of infection, there is a need for an open resource on pathogen–disease associations that can be utilized in computational studies. A large number of pathogen–disease associations is available from the literature in unstructured form and we need automated methods to extract the data.

## Results

We developed a text mining system designed for extracting pathogen–disease relations from literature. Our approach utilizes background knowledge from an ontology and statistical methods for extracting associations between pathogens and diseases. In total, we extracted a total of 3,420 pathogen–disease associations from literature. We integrated our literature–derived associations into a database which links pathogens to their phenotypes for supporting infectious disease research.

## Conclusions

To the best of our knowledge, we present the first study focusing on extracting pathogen–disease associations from publications. We believe the text mined data can be utilized as a valuable resource for infectious disease research. All the data is publicly available from

<https://github.com/bio-ontology-research-group/padimi> and through a public SPARQL endpoint from <http://patho.phenomebrowser.net/>.

**Keywords:** text mining; relationship extraction; pathogen-disease association; pathogen; infectious disease

### **3      Background**

4      Each year, millions of people die due to infectious diseases. The World Health  
5      Organisation (WHO)[1] reported that 11 million deaths were due to HIV/AIDS in  
6      2015 alone. Infectious diseases cause devastating results not only on global public  
7      health but also on the countries economies. Developing countries, especially the  
8      ones in Africa, are the most affected by infectious diseases.

9      Several scientific resources have been developed to support infectious disease re-  
10     search. A large number of these resources focus on host–pathogen interactions [2, 3]  
11     as well as particular mechanisms of drug resistance [4]. Additionally, there are sev-  
12     eral resources that broadly characterize different aspects of diseases [5]. Relatively  
13     little structured information is available about the relationships between pathogens  
14     and disease, information that is also needed to support infectious disease research.  
15     Currently, such associations are mainly covered by proprietary databases such as the  
16     Kyoto Encyclopedia of Genes and Genomes (KEGG) [6] as well as the biomedical  
17     literature and public resources such as Wikipedia [7], MedScape [8], or the Human  
18     Disease Ontology [5] in natural language form. Automated methods are needed to  
19     extract the associations from natural language.

20     Here, we further developed and evaluated a text mining system for extracting  
21     pathogen–disease associations from literature[9]. While most of the existing text  
22     mining studies related to infectious disease focus on extracting host–pathogen in-  
23     teractions from text [10, 11] and archiving this data [2, 3], to the best of our knowl-  
24     edge, we present the first text mining system which focuses on extracting pathogen-  
25     disease associations. Our literature-extracted associations are available for download  
26     from <https://github.com/bio-ontology-research-group/padimi> and through  
27     a public SPARQL endpoint at <http://patho.phenomebrowser.net/>.

## 28 Materials & Methods

### 29 Ontologies and Resources Used

30 We used the latest archived version of the Open Access full text articles ([http://  
europepmc.org/ftp/archive/v.2017.12/](http://europepmc.org/ftp/archive/v.2017.12/), containing approximately 1.8 million  
31 articles) from the Europe PMC database [12]. We used the NCBI Taxonomy [13]  
32 (downloaded on 22-08-2017) and the Human Disease Ontology (DO) [5] (February  
33 (February 2018 release) to provide the vocabulary to identify pathogen and infectious disease  
34 mentions in text. We generated two dictionaries from the labels and synonyms in  
35 the two ontologies and refined them before applying text mining. In the refinement  
36 process, we filtered out terms which have less than three characters and terms  
37 that are ambiguous with common English words (e.g., “Arabia” as a pathogen  
38 name). We extracted only the species labels and synonyms belonging to fungi, virus,  
39 bacteria, worms, insects, and protozoa from NCBI Taxonomy to form our pathogen  
40 dictionary. The final pathogen and disease dictionaries cover a total of 1,250,373  
41 distinct pathogens and 438 distinct infectious diseases.

### 43 Pathogen and Disease class Recognition

44 A class is an entity in an ontology that characterizes a category of things with  
45 particular characteristics. Classes usually have a set of terms attached as labels or  
46 synonyms [14]. We used the Whatizit text mining pipeline [15] to annotate pathogen  
47 and disease classes in text with the two dictionaries for diseases and pathogens.  
48 Because disease name abbreviations can be ambiguous with some other names (e.g.,  
49 ALS is an abbreviation both for “Amyotrophic Lateral Sclerosis” and “Advanced  
50 Life Support”), we used a disease abbreviation filter for screening out the ambiguous  
51 abbreviations that could be introduced during the annotation process [16]. Briefly,  
52 this filter operates based on rules utilizing heuristic information. First, it identifies  
53 abbreviations and their long forms in text by using regular expressions. Second,  
54 it utilizes several rules to decide whether to keep the abbreviation annotated as a

55 disease name or filter out. The rules cover keeping the abbreviation either if any of  
56 its long forms from DO exists in the document or its long form contains a keyword  
57 such as “disease”, “disorder”, “syndrome”, “defect”, etc that describes a disease  
58 name.

## 59 Pathogen–Disease Association Extraction

60 Our association extraction method is based on identification of pathogen-disease  
61 co-occurrences at the sentence level and applying a filter based on co-occurrence  
62 statistics and Normalized Point-wise Mutual Information (NPMI) [17] association  
63 strength measurement to reduce noise possibly introduced by the high recall, low  
64 precision co-occurrence method. We selected the associations having an NPMI value  
65 above 0.2 and co-occurring at least 10 times in the literature.

66 We extended NPMI, which is a measure of collocation between two terms, to  
67 a measure of collocation between two classes. Hence, we reformulated the NPMI  
68 measure for our application. First, we identify, for every class, the set of labels  
69 and synonyms associated with the class ( $Labels(C)$  denotes the set of labels and  
70 synonyms of  $C$ ). We then define  $Terms(C)$  as the set of all terms that can be used  
71 to refer to  $C$ :  $Terms(C) := \{x | x \in Labels(S) \wedge S \sqsubseteq C\}$ .

72 We calculate the NPMI between classes  $C$  and  $D$  as

$$73 npmi(C, D) = \frac{\log \frac{n_{C,D} \cdot n_{tot}}{n_C \cdot n_D}}{-\log \frac{n_{C,D}}{n_{tot}}} \quad (1)$$

74 where  $n_{tot}$  is the total number of sentences in our corpus (i.e., 4,427,138),  $n_{C,D}$  is  
75 the number of sentences in which both a term from  $Terms(C)$  and a term from  
76  $Terms(D)$  co-occur,  $n_C$  is the number of sentences in which a term from  $Terms(C)$   
77 occurs, and  $n_D$  is the number of sentences in which a term from  $Terms(D)$  occurs.

78 **Results**

79 **Statistics on Extracted Pathogen–Disease Associations**

80 We extracted a total of 3,420 distinct pathogen–disease pairs from over 1.8 million  
81 Open Access full text articles. To identify the associations, we used a combination  
82 of lexical, statistical, and ontology-based rules. We used lexical matches to identify  
83 whether the label or synonym of a pathogen or disease is mentioned in a document;  
84 we used a statistical measure, the normalized point-wise mutual information, to  
85 determine whether pathogen and disease mentions co-occur significantly often in  
86 literature; and we used ontologies as background knowledge to expand sets of terms  
87 based on ontology-base inheritance.

88 **Performance Evaluation**

89 To evaluate the pathogen–disease associations we obtain, we used the KEGG [6]  
90 database as reference and compare our results to the information contained in  
91 KEGG. KEGG contains pathogen–disease associations obtained through manual  
92 curation. We could identify 744 pathogen–disease pairs in KEGG for which we  
93 could map the pathogen and disease identifiers from NCBI Taxonomy and DO to  
94 their identifiers in KEGG. Figure 1 shows the overlapping and distinctly identified  
95 pathogen-disease associations from Kegg and literature. We covered 29.4% (219) of  
96 the pathogen–disease associations from KEGG and extracted many more associa-  
97 tions from literature (3,201). There are 525 pairs which we could not cover by text  
98 mining. The main reason we cannot identify an association is due to limitations in  
99 our named entity and normalization procedure.

100 We further evaluated the performance of our system manually on 50 randomly  
101 selected pathogen–disease associations. In our manual evaluation, we achieve a pre-  
102 cision of 64%, a recall of 84% and an F-score of 73%. The false positives were mainly  
103 due to ambiguous abbreviations and pathogen names. For example, “bronchus” was  
104 annotated as an insect name by our method.

105 Some false negatives were due to rejections by the pipeline based on the threshold  
106 settings. For example, “Tetranychus” and “asthma” co-occurred only once in our  
107 corpus and therefore the association between them was rejected as we limited our  
108 analysis to pathogen–disease pairs that co-occurred ten or more times. Other false  
109 negatives were due to missing pathogen or disease labels in our dictionaries.

## 110 Discussion

111 By using ontologies as background knowledge to expand our sets of terms and la-  
112 bels, it is possible to identify pathogen–disease associations even if the labels and  
113 synonyms directly associated with the pathogen or disease are not directly found to  
114 co-occur in text. For example, we extracted a total of 44 distinct pathogen–disease  
115 associations relevant to *dengue disease* (DOID:11205). 12/44 of these associations  
116 are the direct associations of *dengue disease* (i.e., a label or synonym of the disease is  
117 explicitly mentioned in text) while the remaining 32/44 are indirect associations ob-  
118 tained from associations with labels and synonyms of the sub-classes *asymptomatic*  
119 *dengue* (DOID:0050143), *dengue hemorrhagic fever* (DOID:12206), and *dengue shock*  
120 *syndrome* (DOID:0050125). In total, we found 812 pathogen-disease associations  
121 which do not directly co-occur in literature.

122 The performance of our system depends on two parameters: the NPMI value  
123 and the number of co-occurrences used as a threshold. In the future, we may use  
124 these two values to automatically determine optimal threshold based on a more  
125 comprehensive evaluation set of pathogen–disease associations. While our initial  
126 text mining approach performs at a promising level (F-score 73%), there is still  
127 some room for improvements. As we found the pathogen names to be ambiguous  
128 with other domain specific names, we plan to further improve the abbreviation and  
129 name filters we apply. For improving the recall of our system, it may be possible to  
130 expand our dictionaries with other resources covering disease and pathogen names  
131 such as the Experimental Factor Ontology (EFO) [18] and the Unified Medical

132 Language System (UMLS) [19] for diseases, and the Encyclopedia of Life [20] for  
133 pathogens.

134 **Conclusion**

135 Here, we present a text mining method for extracting pathogen-disease associations  
136 from the biomedical literature. Our method performed at the state of the art level  
137 with some room for improvements. In future, we plan to improve our text mining  
138 method by developing and integrating a pathogen abbreviation filter and expanding  
139 the coverage of our pathogen and disease dictionaries. In the scope of infectious  
140 disease research, we have included our results in a database of pathogens and the  
141 phenotypes they elicit in humans. We believe that our results can further support  
142 infectious disease research.

143 **List of abbreviations**

144 World Health Organisation (WHO)  
145 Kyoto Encyclopedia of Genes and Genomes (KEGG)  
146 Human Disease Ontology (DO)  
147 Normalized Point-wise Mutual Information (NPMI)  
148 Experimental Factor Ontology (EFO)  
149 Unified Medical Language System (UMLS)

150 **Funding**

151 This work has been supported by funding from King Abdullah University of Science and Technology (KAUST)  
152 Office of Sponsored Research (OSR) under Award No. URF/1/3454-01-01 and FCC/1/1976-08-01.

153 **Availability**

154 All the data is available from <https://github.com/bio-ontology-research-group/padimi> and  
155 (<http://patho.phenomebrowser.net/>) through a public SPARQL endpoint.

156 **Author's contributions**

157 RH and SK conceived of the study; SK performed all experiments. SK and RH analyzed the results. SK drafted the  
158 manuscript, R.H. revised the manuscript. All authors have read and approved the final version of the manuscript.

159 **Competing interests**

160 The authors declare that they have no competing interests.

161 **Consent of publication**

162 Not applicable.

163 **Ethics approval and consent to participate**

164 Not applicable.

165 **Acknowledgement**

166 Authors would like to thank Mrs. Marwa Abdellatif for her help to make the data available from the SPARQL  
167 end-point.

168 **Author details**

169 <sup>1</sup> Computational Bioscience Research Center, King Abdullah University of Science and Technology, 23955-6900  
170 Thuwal, Saudi Arabia. <sup>2</sup> Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah  
171 University of Science and Technology, 23955-6900 Thuwal, Saudi Arabia.

172 **References**

- 173 1. World Health Organisation. <http://who.int/en/>
- 174 2. Ammari, M.G., Gresham, C.R., McCarthy, F.M., Nanduri, B.: HPIDB 2.0: a curated database for  
175 host-pathogen interactions. *Database* **2016** (2016)
- 176 3. Wardehant, M., Risley, C., McIntyre, M.K., Setzkorn, C., Baylis, M.: Database of host-pathogen and related  
177 species interactions, and their global distribution. *Scientific Data* **2**(150049, eCollection2015) (2015)
- 178 4. Jia, B., Raphenya, A.R., Alcock, B., Waglechner, N., Guo, P., Tsang, K.K., Lago, B.A., Dave, B.M., Pereira,  
179 S., Sharma, A.N., Doshi, S., Courtot, M., Lo, R., Williams, L.E., Frye, J.G., Elsayegh, T., Sardar, D., Westman,  
180 E.L., Pawlowski, A.C., Johnson, T.A., Brinkman, F.S.L., Wright, G.D., McArthur, A.G.: Card 2017: expansion  
181 and model-centric curation of the comprehensive antibiotic resistance database. *Nucleic Acids Research*  
182 **45**(D1), 566–573 (2017)
- 183 5. Kibbe, W.A., Arze, C., Felix, V., Mitraka, E., Bolton, E., Fu, G., Mungall, C.J., Binder, J.X., Malone, J.,  
184 Vasant, D., Parkinson, H.E., Schriml, L.M.: Disease ontology 2015 update: an expanded and updated database  
185 of human diseases for linking biomedical knowledge through disease data. *Nucleic Acids Research*  
186 **43**(Database-Issue), 1071–1078 (2015)
- 187 6. Kanehisa, M., Goto, S.: KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Research* **28**(1), 27–30  
188 (2000)
- 189 7. List of Infectious Diseases. [https://en.wikipedia.org/wiki/List\\_of\\_infectious\\_diseases](https://en.wikipedia.org/wiki/List_of_infectious_diseases)
- 190 8. Medscape. <https://emedicine.medscape.com/>
- 191 9. Kafkas, S., Hoehndorf, R.: Ontology based mining of pathogen – disease associations from literature. In:  
192 Hoehndorf, R., Dumontier, M. (eds.) *Proceedings of Bio-Ontologies SIG@ISMB 2018*, 6–10 July 2018; Chicago,  
193 USA. (2018)
- 194 10. Thieu, T., Joshi, S., Warren, S., Korkin, D.: Literature mining of host-pathogen interactions: comparing  
195 feature-based supervised learning and language-based approaches. *Bioinformatics* **28**(6), 867–875 (2012)
- 196 11. İlknur Karadeniz, Hur, J., He, Y., Özgür, A.: Literature mining and ontology based analysis of host-brucella  
197 gene-gene interaction network. *Frontiers in Microbiology* **6**(1386) (2015)
- 198 12. Europe PMC: a full-text literature database for the life sciences and platform for innovation. *Nucleic Acids  
199 Research* **43**(Database-Issue), 1042–1048 (2015)
- 200 13. Sayers, E.W., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Chetvernin, V., Church, D.M., DiCuccio,  
201 M., Edgar, R., Federhen, S., Feolo, M., Geer, L.Y., Helmberg, W., Kapustin, Y., Landsman, D., Lipman, D.J.,  
202 Madden, T.L., Maglott, D.R., Miller, V., Karsch-Mizrachi, I., Ostell, J., Pruitt, K.D., Schuler, G.D., Sequeira,  
203 E., Sherry, S.T., Shumway, M., Sirotnik, K., Souvorov, A., Starchenko, G., Tatusova, T.A., Wagner, L.,  
204 Yaschenko, E., Ye, J.: Database resources of the national center for biotechnology information. *Nucleic Acids  
205 Research* **37**(Database-Issue), 5–15 (2009)

206 14. Hoehndorf, R., Schofield, P.N., Gkoutos, G.V.: The role of ontologies in biological and biomedical research: a  
207 functional perspective. *Briefings in Bioinformatics* **16**(6), 1069–1080 (2015)

208 15. Rebholz-Schuhmann, D., Arregui, M., Gaudan, S., Kirsch, H., Jimeno, A.: Text processing through web  
209 services: calling whatizit. *Bioinformatics* **24**(2), 296–298 (2008)

210 16. Kafkas, S., Dunham, I., McEntyre, J.R.: Literature evidence in open targets - a target validation platform. *J.*  
211 *Biomedical Semantics* **8**(1), 20–1209 (2017)

212 17. Bouma, G.: Normalized (pointwise) mutual information in collocation extraction. In: *Proceedings of the*  
213 *Biennial GSCL Conference: 2009; Potsdam, Germany*, pp. 31–40 (2009)

214 18. Malone, J., Holloway, E., Adamusiak, T., Kapushesky, M., Zheng, J., Kolesnikov, N., Zhukova, A., Brazma, A.,  
215 Parkinson, H.E.: Modeling sample variables with an experimental factor ontology. *Bioinformatics* **26**(8),  
216 1112–1118 (2010)

217 19. Bodenreider, O.: The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic*  
218 *Acids Research* **32**(Database-Issue), 267–270 (2004)

219 20. Encyclopedia of Life. <http://eol.org/>

220 **Figures**

221 Figure 1. Overlapping pathogen–disease associations between Kegg and literature

Literature

Kegg

525

219

3201