# DISNET: A framework for extracting phenotypic disease information from public sources

**Gerardo Lagunes García** [1] , **Alejandro Rodríguez González** [Corresp., 1, 2] , **Lucia Prieto Santamaría** [1] , **Eduardo P García del Valle** [1] , **Massimiliano Zanin** [1] , **Ernestina Menasalvas Ruiz** [1]

[1] Centro de Tecnología Biomédica, Universidad Politécnica de Madrid, Pozuelo de Alarcón, Madrid, Spain

[2] Escuela Técnica Superior de Ingenieros Informáticos, Universidad Politécnica de Madrid, Boadilla del Monte, Madrid, Spain

Corresponding Author: Alejandro Rodríguez González
Email address: alejandro.rg@upm.es

**Background.** Within the global endeavour of improving population health, one major challenge is the identification and integration of medical knowledge spread through several information sources. The creation of a comprehensive dataset of diseases and their clinical manifestations based on information from public sources is an interesting approach that allows, not only to complement and merge medical knowledge, but also to increase it and thereby to interconnect existing data and analyse and relate diseases to each other. In this paper, we present DISNET (disnet.ctb.upm.es), a web-based system designed to periodically extract the knowledge from signs and symptoms retrieved from medical databases, and to enable the creation of customisable disease networks.

**Methods.** We here present the main features of the DISNET system. We describe how information on diseases and their phenotypic manifestations is extracted from Wikipedia and PubMed websites; specifically, texts from these sources are processed through a combination of text mining and natural language processing techniques.

**Results.** We further present the validation of our system on Wikipedia and PubMed texts, obtaining the relevant accuracy. The final output includes the creation of a comprehensive symptoms-disease dataset, shared (free access) through the system's API. We finally describe, with some simple use cases, how a user can interact with it and extract information that could be used for subsequent analyses.

**Discussion.** DISNET allows retrieving knowledge about the signs, symptoms and diagnostic tests associated with a disease. It is not limited to a specific category (all the categories that the selected sources of information offer us) and clinical diagnosis terms. It further allows to track the evolution of those terms through time, being thus an opportunity to analyse and observe the progress of human knowledge on diseases. We further discussed the validation of the system, suggesting that it is good enough to be used to extract diseases and diagnostically-relevant terms. At the same time, the evaluation also revealed that improvements could be introduced to enhance the system's reliability.

1  # DISNET: A framework for extracting phenotypic
2  # disease information from public sources
3
4
5  Gerardo Lagunes García[1], Alejandro Rodríguez González[1,2], Lucía Prieto Santamaría[1], Eduardo
6  P. García del Valle[2], Massimiliano Zanin[1], Ernestina Menasalvas Ruíz[1,2]
7
8  [1] Centro de Tecnología Biomédica, Universidad Politécnica de Madrid, 28223 Pozuelo de
9  Alarcón, Madrid, Spain
10  [2] ETS de Ingenieros Informáticos, Universidad Politécnica de Madrid, 28660 Boadilla del
11  Monte, Madrid, Spain
12
13  Corresponding Author:
14  Alejandro Rodríguez González[1]
15  Centro de Tecnología Biomédica, Universidad Politécnica de Madrid, 28223 Pozuelo de
16  Alarcón, Madrid, Spain
17  Email address: alejandro.rg@upm.es
18

19  ## Abstract
20
21  **Background.** Within the global endeavour of improving population health, one major challenge
22  is the identification and integration of medical knowledge spread through several information
23  sources. The creation of a comprehensive dataset of diseases and their clinical manifestations
24  based on information from public sources is an interesting approach that allows, not only to
25  complement and merge medical knowledge, but also to increase it and thereby to interconnect
26  existing data and analyse and relate diseases to each other. In this paper, we present DISNET
27  (disnet.ctb.upm.es), a web-based system designed to periodically extract the knowledge from
28  signs and symptoms retrieved from medical databases, and to enable the creation of customisable
29  disease networks.
30  **Methods.** We here present the main features of the DISNET system. We describe how
31  information on diseases and their phenotypic manifestations is extracted from Wikipedia and
32  PubMed websites; specifically, texts from these sources are processed through a combination of
33  text mining and natural language processing techniques.
34  **Results.** We further present the validation of our system on Wikipedia and PubMed texts,
35  obtaining the relevant accuracy. The final output includes the creation of a comprehensive
36  symptoms-disease dataset, shared (free access) through the system's API. We finally describe,
37  with some simple use cases, how a user can interact with it and extract information that could be
38  used for subsequent analyses.

39   **Discussion.** DISNET allows retrieving knowledge about the signs, symptoms and diagnostic
40   tests associated with a disease. It is not limited to a specific category (all the categories that the
41   selected sources of information offer us) and clinical diagnosis terms. It further allows to track
42   the evolution of those terms through time, being thus an opportunity to analyse and observe the
43   progress of human knowledge on diseases. We further discussed the validation of the system,
44   suggesting that it is good enough to be used to extract diseases and diagnostically-relevant terms.
45   At the same time, the evaluation also revealed that improvements could be introduced to enhance
46   the system's reliability.
47

48   **Introduction**
49   In 1796, Edward Jenner found an important link between the variola virus, which affected only
50   humans and was highly lethal, and the bovine smallpox virus, which attacked cows and was
51   transmitted to humans by physical contact with infected animals, and which, despite its severity,
52   rarely resulted in death. He found that people who became infected with the latter (also called
53   cowpox) did not subsequently catch the former; and thus, that something in the bovine smallpox
54   virus made humans immune to variola virus. This led him to thoroughly investigate the
55   relationship between these diseases and understand their behaviour for more than twenty years;
56   to be finally able to find a cure for the variola virus, saving thousands of humans lives
57   worldwide.
58

59   This discovery illustrates the importance of the knowledge that we can get from diseases and,
60   more specifically, from how they are related. Despite the fact that in the last 200 years our
61   understanding of diseases has greatly increased, and valuable advances have been made in this
62   area (Botstein & Risch, 2003), the number of those without treatment or cure is still extremely
63   high (e.g. Alzheimer's disease, small cell lung cancer, HIV, etc.). It is thus imperative to explore
64   new approaches and tools to tackle them and, therefore, improve the health of the world's
65   population.
66

67   It is almost a truism that the search for new drugs requires a better understanding about diseases.
68   This includes finding new insights on the relationship between diseases (which diseases are
69   related and how), as well as the creation of public and easy-to-access large databases of diseases
70   knowledge (Pérez-Rodríguez et al., 2019). During the last decade, such endeavour has been
71   vastly facilitate by the World Wide Web. On one hand, it is possible to find free biomedical
72   vocabularies like Unified Medical Language System (UMLS) (Bodenreider, 2004), Human
73   Phenotype Ontology (HPO) (Robinson et al., 2008; Köhler et al., 2017), Disease Ontology (DO)
74   (Schriml et al., 2012) or MeSH (Lipscomb, 2000), all of them offering disease classifications,
75   disease coding standards and associated medical resources. On the other hand, one can find
76   bioinformatic databases created by complex medical system, like DiseaseCard (Oliveira et al.,
77   2004; Dias et al., 2005; Lopes & Oliveira, 2013), MalaCards (Rappaport et al., 2013, 2014; Espe,

78 2018), GeneCard (Safran et al., 2002), Diseases Database (DD)[1], DISEASES (Pletscher-Frankild
79 et al., 2015), SIGnaling Network Open Resource (SIGNOR) (Perfetto et al., 2016), Kyoto
80 Encyclopedia of Genes and Genomes (KEGG) (Kanehisa & Goto, 2000), MENTHA (Calderone,
81 Castagnoli & Cesareni, 2013), PhosphositePlus (Hornbeck et al., 2015), PhosphoELM
82 (Hornbeck et al., 2015), UniProtKB (UniProt Consortium, 2014), Human Gene Mutation
83 Database (HGMD) (Stenson et al., 2014), Comparative Toxicogenomics Database (CTD)
84 (Mattingly et al., 2006), and the database for Pediatric Disease Annotation and Medicine
85 (PedAM) (Jia et al., 2018). These datasets have generally been created by processing the
86 information from several sources, and they usually offer simple search engines; yet, not all of
87 them have a systematic and structured form of sharing their knowledge. In this context, it is
88 important to relate the quantity of available medical sources and systems on one hand, and the
89 need of health professionals for quality information on the other, helping them performing their
90 work with higher precision and lower time (Russell-Rose, Chamberlain & Azzopardi, 2018).
91 Therefore, diagnostic systems (Chen et al., 2018) have become more relevant and researchers
92 such as Xia et al. attempt to take on the challenge through the mining of information from
93 sources such as DO, Symptom Ontology (SYMP) and MEDLINE/PubMed citation records (Xia
94 et al., 2018). We can also observe in the literature a large volume of studies that use the mining
95 of texts from different unstructured or semi-structured medical information sources (Frunza,
96 Inkpen & Tran, 2011; Mazumder et al., 2016; Singhal, Simmons & Lu, 2016; Xu et al., 2016;
97 Tsumoto et al., 2017; Sudeshna, Bhanumathi & Hamlin, 2017; Aich et al., 2017; Gupta et al.,
98 2018; Rao & Rao, 2018; Zhao et al., 2018; Bou Rjeily et al., 2019).
99
100 It goes without doubt that the large amount of available bioinformatic resources allows both to
101 enhance the research in the biomedical field and to have a better understanding of how the
102 diseases behave and how can we fight them. However, most of the already mentioned sources
103 are mainly focused on retrieving and exposing the captured knowledge and are not primarily
104 focused on the analysis of the interactions and relationships that exists between different diseases
105 or different disease characteristics.
106
107 In this context, several works have attempted to understand these relationships by creating and
108 analysing disease networks. The complexity of such endeavour was soon clear, as diseases may
109 share not only symptoms and signs, but also genes, proteins, causes and, in many cases, cures
110 (Goh et al., 2007; Zanzoni, Soler-López & Aloy, 2009; Barabási, Gulbahce & Loscalzo, 2010;
111 Lee et al., 2011; Zhou et al., 2014; Chen et al., 2015; Quwaider & Alfaqeeh, 2016; Piñero et al.,
112 2017; Lo Surdo et al., 2018; Hwang et al., 2019; Szklarczyk et al., 2019; García del Valle et al.,
113 2019). One of the most important works on the subject was published in 2007 by K.-I. Goh et al.
114 (Goh et al., 2007), in which the HDN (Human Disease Network) was developed, a network of
115 human diseases and disorders that links diseases based on their genetic origins and biological
116 interactions. Different diseases were then associated according to shared genes, proteins or

---

[1] http://www.diseasesdatabase.com

117  protein interactions. The hypothesis that different diseases, with potentially different causes, may
118  share characteristics allows the design of common strategies regarding how to deal with the
119  diagnosis, treatment and prognosis of a disease.
120
121  Within this line of research it is worth mentioning the Human Symptoms-Disease Network
122  (HSDN) (Zhou et al., 2014), an HDN network in which similarities between diseases were
123  estimated through common symptoms. This is an important change in perspective with respect to
124  previous works, in which the focus was centred on the genetic and biological origin of the
125  diseases. In (Zhou et al., 2014), diseases are defined by their clinical phenotypic manifestations,
126  i.e. signs and symptoms; this is not surprising, as these manifestations are basic medical
127  elements, and crucial characteristics in the diagnosis, categorization and clinical treatment of the
128  diseases. It was then proposed to use these as a starting point to understand the existing
129  relationships between different diseases.
130
131  Building on top of these previous works and stemming from the necessity of having exhaustive
132  and accurate sources of disease-based information, in this paper we present the DISNET
133  (Diseases Networks) system. DISNET aims at going one step further in improving human
134  knowledge about diseases, not only by seeking and analysing the relations between them, but
135  most importantly, by finding real connections between diseases and drugs, thus potentially
136  enabling novel drug repositioning strategies.
137
138  Therefore, the objectives of this research work are:
139
140      • Present the first version of the web-based DISNET (phenotypic information) system.
141      • Describe the characteristics of its recovery and generation process of phenotypic
142        knowledge.
143      • Provide an indicator of the accuracy of the information generated by DISNET, through a
144        manual information validation process.
145      • Provide free access to the DISNET dataset with structured information about diseases and
146        symptoms through the system's API.
147
148  The current version of the DISNET system is focused on phenotypic information and allows to
149  capture knowledge about diseases from heterogeneous textual sources. We have five main
150  advantages with respect to the previously described research. Firstly, the use of Wikipedia as the
151  main source of knowledge. While this encyclopaedia has been the object of study of numerous
152  research works, to the best of our knowledge DISNET is the first system to mine texts that
153  describe how the disease manifests itself, and to recover disease codes, leading to a more
154  extensive mapping between several biomedical information sources. Secondly, DISNET offers a
155  public API, that enables the free and structured sharing of the knowledge generated by the
156  system; it is worth noting that having an appropriate method for information sharing, while being

157   a basic element, is not common among the previously reviewed systems. Thirdly, the proposed
158   system allows to recover the temporal evolution of phenotypic information. This is especially
159   relevant for sources like Wikipedia, which is constantly edited, and whose medical articles are
160   frequently corrected and updated. This allows an analysis of the dynamics of diseases, in terms
161   of how their description has been evolving through a collective effort. Fourthly, DISNET has
162   been designed to be able to store and integrate information from heterogeneous sources; this
163   allows to cross-validate and enrich medical knowledge of diseases and symptoms. Future content
164   to be introduced includes genetic and drug information to create a complex multilayer network,
165   where each layer represents the different type of information (phenotypical, biological, drugs).
166   Finally, we also provide an evaluation of the DISNET extracted content, with examples on how
167   diseases can be analysed and their relationships described through a network structure.
168
169   Beyond this introduction, this paper is organised as follows: Section 2 explains the technologies
170   used in the creation of DISNET phenotypical features repository. Section 3 presents the main
171   results obtained in the validation of the system and discussion about them, describes several
172   simple use cases. Finally, Section 4 draws some conclusions and discusses future work.
173

174   ## Materials & Methods
175   This Section discusses the technical characteristics of the DISNET system, focusing on two
176   aspects: the sources of information hitherto considered, and the DISNET workflow. More
177   specifically, the last point describes how the system retrieves phenotypic information, in the
178   form of raw texts, from the discussed sources; how these texts are processed to obtain diagnostic
179   terms; and how these terms are validated to compile a final list of valid symptom-type terms[2].
180   The source code of the entire DISNET platform and their components is available online[3].
181

182   ## Information Source
183   As it has previously been shown, it is customary for works aimed at unveiling relationships
184   between diseases to focus on single source of information, in most cases just *abstracts* of
185   Medline articles. On the other hand, the proposed system aims at obtaining inputs from as many
186   sources as possible, to guarantee the recovery of as much knowledge as possible. By bringing
187   together information from different sources, we expect them to complement each other, creating
188   a network with a higher capacity of relating diseases. The rationale for this is that the different
189   sources of textual knowledge, such as Wikipedia or PubMed, are written in different styles and
190   by people with different backgrounds; the information they contain may therefore be
191   complementary. In order to take advantage of such richness, the DISNET system allows the user
192   to query the symptoms according to different rules: for instance, from one or multiple sources,
193   by applying filters based on prevalence information, or on percentages of similarity among
194   others. This clearly comes at a cost: the system should be flexible enough to be able to process

---

[2] Study approved by Ethics Committee of Universidad Politécnica de Madrid.

[3] https://github.com/disnet-project/

195    sources with different structures. In the remainder of this Section we discuss the patterns used to
196    select data sources, how they have been mined, and finally the challenges involved in such tasks.
197

198    **Source Selection**
199    Traditionally, in order to obtain the whole body of knowledge that mankind has accumulated
200    about a given disease, one would refer to medical books. Although books usually contain much
201    of the information available, they also present some important limitations: they are not constantly
202    updated; the automatic access to their content is difficult, especially when digital versions are not
203    available; and they are usually written for study, thus the information they contain is not
204    structured for data mining tasks. On the other hand, one has the World Wide Web, whose main
205    characteristic is to be (mostly) free accessible to anyone with an internet connection. It mainly
206    offers three sources of information. Firstly, the abstract, and in some cases, the full text, of
207    medical papers, which can be accessed through platforms like PubMed. Secondly, specialized
208    sources of information, such as MedlinePlus, MayoClinic, or CDC. Finally, good medical data
209    can be obtained in sources of knowledge that are not specialized, such as Wikipedia or Freebase.
210    Note that all of them have different characteristics, in terms of comprehensiveness, degree of
211    structure of the information, and up-to-datedness.
212

213    The criteria used for the selection of the sources of information in DISNET are: i) open access,
214    ii) recognised quality and reliability, and iii) availability of substantial quantities of data
215    (structured or not). This suggested to include the following three web sites in the system, which
216    are described below: i) Wikipedia and ii) PubMed. It is important to note that the system is not
217    closed; on the contrary, thanks to its flexibility, new sources could (and will) be incorporated in
218    the future.
219

220    **Wikipedia**
221    Wikipedia is an online, open and collaborative source of information. It was created by the
222    Wikimedia Foundation and its English edition is the largest and most active one. The
223    monumental and primary task of editing, revising and improving the quality of all articles is not
224    performed by a core of administrators: it is instead the collaborative result of thousands of users.
225    Consequently, this encyclopaedia is considered the greatest collective project in the history of
226    humanity (Mehdi et al., 2017; Aibar, 2017).
227

228    Wikipedia contains more than 155,000 articles in the field of medicine (Azzam et al., 2017) and
229    is one of the most widely used medical sources (Friedlin & McDonald, 2010) by the general
230    community (Aibar, 2017) and also by medical specialists (Azer, 2014; Shafee et al., 2017), the
231    latter ones having deeply been involved in its enrichment (Azzam et al., 2017)(Cohen, 2013).
232    One of the initiatives is the Cochrane/Wikipedia, which aims at increasing reliability in articles
233    with medical content (Matheson & Matheson, 2017). In 2014 Wikipedia was referred to as "*the*
234    *single leading source of medical information for patients and health care professionals*" by the

235  Institute of Medical Science (IMS) (Heilman & West, 2015). This stems from the fact that an
236  increasing number of people in the medical field are becoming aware of the importance of
237  collaborating and generating quality content in the world's largest online encyclopaedia.
238  We have focused on Wikipedia in its English edition, and specifically on those articles
239  categorized as diseases. In order to obtain a list of such articles we resort to conventional
240  DBpedia and DBpedia-Live (DBpedia), an open and free Web repository that stores structured
241  information from Wikipedia and other Wikimedia projects. By containing structured
242  information, this source allows complex questions to be asked through SPARQL queries
243  ("SPARQL Query Language for RDF," 2017). We developed a query[4] that is able to get all the
244  articles of Wikipedia in English referring to human diseases and run it in the **Virtuous**
245  **environment SPARQL Query Editor of DBpedia**[5]. This first approach to detecting and
246  extracting Wikipedia's web links can be addressed in different ways and in the **Error! Reference**
247  **source not found.** section we will talk about them.

248

249  Even though disease articles have a standard structure, due to the very nature of Wikipedia,
250  articles can be edited by anyone; consequently, it is possible to find articles that do not comply
251  with the standard form that the creators of the encyclopaedia propose ("Wikipedia," 2018). The
252  structure is organized in sections, of which we have selected those whose content is related to the
253  phenotypic manifestations of the disease. The essential sections mined by DISNET are: "*Signs*
254  *and symptoms*", "*signs and symptoms*", "*Symptoms and causes*", "*Signs*", "*Symptoms*",
255  "*Causes*", "*Cause*", "*Diagnosis*", "*Diagnostic*", "*Causes of injury*", "*Diagnostic approach*",
256  "*Presentation*", "*Symptoms of ...*", "*Causes of ...*" , and *infobox.*

257

258  The data retrieved from these sections are: i) the texts (paragraphs, lists and tables) contained in
259  the previously described sections; ii) the links contained in these texts; and iii) the disease codes
260  of vocabularies external to Wikipedia, which can be found in the *infoboxes* of the article. Note
261  there are two types of *infobox*. Figure 1 shows an example of the external vocabulary codes
262  retrieved in a vertical *infobox*, usually located at the beginning of the document; Figure 2 shows
263  an example of a horizontal *infobox*, generally located at the foot of the document. These disease
264  codes in different vocabulary are relevant elements when searching for diseases in the system's
265  database. The list of external vocabularies to DISNET can be found at [6].

266

267  **PubMed**
268  PubMed[7] comprises more than 28 million biomedical literature citations from MEDLINE, life
269  science journals and online books. Quotations may include links to full text content from

---

270   PubMed Central[8] and editorial websites (pubmeddev). As in other studies, we here only
271   considered the abstracts of the articles, as, firstly, it is not always possible to access the full text,
272   and secondly, the full text of articles does not follow a standard format. However, we are aware
273   of the limitations of the extraction of information only for abstracts (Westergaard et al., 2018),
274   and future versions of DISNET platform will focus in extracting the content from the full paper
275   when possible.  Note that in PubMed the information about one single disease is spread among
276   multiple documents – as opposed to Wikipedia, in which there is a bijective relationship between
277   articles and diseases.
278
279   Obtaining the list of diseases in PubMed involves two main steps. Firstly, one should extract the
280   list of MeSH terms (DMTL) relating to human diseases *C*, which are categorized from *C01* to
281   *C20* (excluding those categories such as "Animal Diseases" or "Wounds and Injuries") as shown
282   in the classification tree in Fig. 3[9]; and map each disease with Human Disease Ontology[10] to
283   obtain disease codes of the vocabulary ICD-10, OMIM, MeSH, SNOMED_CT and UMLS. Note
284   that the use of multiple vocabularies aims at obtaining the greatest amount of means (identified
285   codes) to identify diseases in different sources of information. As a second step, it is necessary to
286   extract all PubMed articles whose terms are associated with each of the elements of the
287   previously extracted disease list DMTL, through PubMed's Entrez API (AEPM) it is possible to
288   carry out this task, because this allows access to all Entrez databases including PubMed, PMC,
289   Gene, Nuccore and Protein. An important feature to mention of the AEPM, and also used in our
290   work, has been the sorting of articles by their relevance (Information et al., 2019), managing to
291   focus the efforts on those articles with better quality. Thus, this configuration has given us the
292   possibility to obtain, if they exist, the 100 most relevant articles of each MeSH term consulted.
293   Specifically, for each article we retrieve: 1) abstract, 2) authors' names, 3) unique identifier in
294   PubMed and PubMed Central, 4) doi (digital object identifier), 5) title, 6) associated MeSH
295   terms and 7) keywords. The workflow for extracting texts from PubMed documents is shown in
296   Fig. 4.
297

298   **Challenges**
299   Mining information from the sources previously described entails several computational
300   challenges, which may be boiled down to one requirement for the DISNET system: the need of a
301   high versatility in data acquisition. We here review such challenges, as these partly explain the
302   adopted software solution.
303

304   First of all, the mapping disease-webpage may take different forms. Specifically, it is one to one
305   for Wikipedia, as all the information of a disease is included in a single page; but it becomes one
306   to many for PubMed, in which multiple articles are available for each single concept. Consulting
307   the latter thus requires a more complex procedure.

---

[8] https://www.ncbi.nlm.nih.gov/pmc/

[9] https://b.nlm.nih.gov/treeView

[10] http://www.obofoundry.org/ontology/doid.html

308

309 Secondly, and as one may expect, the specific structure of each source of information is different
310 – i.e. a page of Wikipedia has not the same structure of a PubMed article. This requires further
311 flexibility, in terms of the development of a modular structure with specific crawlers for each
312 source.
313

314 Finally, it is worth noting that, while here we have only considered texts, much information is
315 available in different medias, like images, videos and others binary files. While not implemented
316 at this stage, the system should be flexible enough to accommodate such sources in the future.
317

318 **Data Retrieval and Knowledge Extraction**
319 This section describes the general architecture of the DISNET system, including the data
320 extraction and the subsequent knowledge extraction. In the sake of clarity, such architecture is
321 further depicted in Fig. 5.
322

323 **The Extraction Process**
324 The first step of the DISNET pipeline is in charge of retrieving the information from the sources
325 previously identified and described. For each one of this, and before running the actual web
326 crawler, the "Get Disease List Procedure" (GDLP) component is responsible for obtaining the
327 list of diseases to be mined, thus providing links to all available disease related documents. For
328 example, the GLDP associated to Wikipedia articles makes use of the SPARQL query[1];
329 similarly, the links for the PubMed's articles are retrieved through a list of MeSH terms.
330

331 Once the URL list has been collected, the "Web Crawler" (WC) module is in charge of
332 connecting to each of the hyperlinks and extracting the specific text that describes the
333 phenotypical manifestations, as well as the links (references) contained within the texts[11]. In
334 addition, and whenever possible, it attempts to extract information related to the coding of
335 diseases, i.e. the codes used to identify the disease in different databases or existing data
336 vocabularies. Currently it is able to retrieve information from more than 6,692 articles in
337 Wikipedia and from 229,160 article abstracts in PubMed. The information mined by WC is
338 stored in an intermediate database called "Raw DB", which contains the raw unprocessed text.
339 The next step within the pipeline is called "NLP Process" (NLPP). This component is
340 responsible for: i) reading all the texts of a snapshot, and ii) obtaining for each text a list of
341 relevant clinical concepts/terms, discarding any unrelated paragraphs or words. At the moment
342 NLPP uses Metamap (Aronson, 2001)(Rodríguez González et al., 2018) as a Natural Medical
343 Language Processing tool to extract clinical terms of interest – see online NLP Tools and
344 Configuration section[12]. Semantic types (SM) are important elements created by UMLS to define

---

[11] https://jsoup.org/
[12] http://disnet.ctb.upm.es/apis/disnet#NLP_Tools_and_Configuration

345 categories of concepts. Metamap uses SM to find medical elements, and a full list of them is
346 available online[13].

347

348 The output of the NLP process is stored in the "DISNET Medical DB" (DMDB) database. It
349 stores, in a structured way, the medical concepts that have been obtained by the NLPP, as well as
350 any information required to track the origin of such concepts – in order to track any error that
351 may later be detected. Therefore, and to summarize, the information stored in a structured way in
352 DMDB is: i) the medical concepts with their location, information and semantic types, ii) the
353 texts from which they were extracted and the links by them contained, iii) the sections which the
354 texts belong to, iv) the document or documents describing the disease (Web link) and v) the
355 disease identifiers codes in different vocabulary or databases. Additional information, as the day
356 of the extraction and the source, is further saved.

357

358 Before reaching the last step of the process, it is important to highlight the nature of the
359 information hitherto stored. Specifically, the system has not extracted only signs or symptoms of
360 a disease, but instead medical terms that we believe may be phenotypic manifestations of
361 disease. It is thus necessary to filter those that are not relevant for the objective initially
362 described.

363

364 Having clarified this, the next component of the pipeline, the "TVP Process" TVPP, reads all the
365 concepts of a snapshot - source pair and filters them. This process is responsible for determining
366 whether these UMLS medical terms are really phenotypic manifestations, and for storing the
367 results back in the DMDB. TVPP is based on the Validation Terms Extraction Procedure that
368 was developed, implemented and tested by Rodriguez-Gonzalez et al (Rodríguez-González et al.,
369 2015). The results of this component (a purification of concepts) are thus those validated terms
370 that we will consider as true phenotypic manifestations of diseases.

371

372 The DISNET extraction process (IEPD), i.e. the process of retrieving and storing information
373 about diseases, basically ends here. Nevertheless, for the sake of providing an accessible and
374 user-friendly way of retrieving and manipulating this information, DISNET also offers a REST-
375 based interface. This is described in detail in the system website
376 (http://disnet.ctb.upm.es/apis/disnet); also refer to Section 4 for an application example.

377

378 **Results**
379 This section describes how the medical concepts data set is built, for then validating and
380 analysing its content.

381

382 **Construction of the DB**

---

[13] https://metamap.nlm.nih.gov/SemanticTypesAndGroups.shtml

383  The database in the DISNET system contains information recovered from three sources of
384  information: Wikipedia and PubMed. From Wikipedia we have 26 snapshots, from February 1st,
385  2018 to February 15th, 2019, for PubMed we have one snapshot, that of April 3rd, 2018. Within
386  the system it is possible to consult, for each snapshot and source, the total number of articles
387  with medical terms, the total number of medical terms found, the number of processed texts, the
388  total number of retrieved codes, and the total number of semantic types found[14].
389  When summing that sources, the system counts with 6,545 diseases, 2,212 medical terms from
390  UMLS (SNOMED-CT) and 19 semantic types, which can be consulted online[15].
391  Wikipedia snapshots are built using the configurations that are available online[16]. We have
392  obtained a list of 11,074 articles catalogued as diseases in Wikipedia according to DBpedia[17],
393  from which we obtained 6,692 articles with at least one text referring to phenotypic knowledge
394  of the disease, or at least one code to an external information source, 4,798 of which were found
395  to be relevant medical concepts[18].
396
397  The snapshot for PubMed has been built using the configuration described online[19]. This
398  snapshot has been built on top of a list of 2,354 MeSH terms[19] referring to human diseases, but
399  only for 2,213 MeSH terms did we obtain information (199,013 scientific articles in total, i.e.
400  about 0.71% of the 28 million articles existing in PubMed[20]) and of each of these PubMed
401  articles obtained, only in 174,900 were abstracts found and only in 125,515 were relevant
402  medical terms found. Figure 6 and Figure 7 presents some basic database statistics at an
403  aggregated level as well as by source (for Wikipedia and PubMed). Some notable differences can
404  be observed; for instance, the five most common terms for Wikipedia are *Pain*, *Lesion*,
405  *Neoplasms*, *Magnetic resonance imaging*, *Inflammation* and *Malnutrition*, while for PubMed
406  these are *Neoplasms*, *Lesion*, *Magnetic resonance imaging*, *Malnutrition* and *Inflammation*.
407  Similarly, the three diseases with the greatest number of concepts in Wikipedia are *Kawasaki*
408  *disease*, *Cerebral palsy* and *Hypoglycemia,* while for PubMed these are *Hypercalcemia*, *Cranial*
409  *nerve palsy* and *Beriberi*.
410

411  **Data evaluation of the DB**
412  In this section, we discuss the results of the validation process we executed on the system, to
413  ensure the relevance of the diagnostic knowledge (valid medical diagnostic terms) generated
414  through our NLP process (Metamap and TVP). The evaluation has been made on both Wikipedia
415  and PubMed mined.
416

---

[14] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/knowledge_sources

[15] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/DISNET_summing_source_counts

[16] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/snapshot_settings.txt

[17] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/wikipedia_diseases_articles_by_dbpedia.txt

[18] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/wikipedia_articles_with_relevant_terms.txt

[19] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/mesh_terms_human_diseases.txt

[20] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/list_pubmed_papers.txt

417   The validation for Wikipedia was carried out on the February 1st, 2018 snapshot, selecting 100
418   diseases at random with the only condition of having at least 20 valid medical terms in order to
419   ensure that our validation procedure analyses articles with a high concentration of medical
420   knowledge. Similarly, the validation for PubMed has been done on the April 3rd, 2018 snapshot,
421   selecting a random sample of 100 article abstracts. It is necessary to highlight that the validation
422   procedure was designed to carry out on articles and due to the nature of each of the sources it is
423   necessary to remember that Wikipedia articles are composed by one or more texts, while
424   PubMed articles are composed by only one text, the abstract. And for this reason for Wikipedia,
425   to validate an article means to validate a disease, for PubMed to validate an article means to
426   validate a part of a disease. These snapshots were performed at different times, and therefore
427   with different configurations – the latter ones can be viewed online[19]. During the validation of
428   Wikipedia, we detected that the initial configuration of Metamap did not find all the necessary
429   medical concepts: for instance, Anxiety, Stress, Amnesia, Bulimia and other psychological
430   concepts were missing. We therefore decided to update the initial list of semantic types to be
431   detected (see online NLP Tools and Configuration section[16]) by adding the following elements:
432   **Intellectual Product**, **Mental Process**, **Mental or Behavioral Dysfunction**, **Pathologic**
433   **Function**, **Congenital Abnormality**.
434
435   The evaluation was conducted through a thorough manual analysis of the basic data. For each
436   disease obtained from Wikipedia or PubMed we compared: (1) the list of medical terms
437   extracted manually from the texts describing the disease; (2) the list of medical terms extracted
438   by Metamap from the same texts; (3) the value (TRUE=valid or FALSE=invalid) resulting from
439   the TVP process for each term found by Metamap; (4) the value of diagnostic relevance for a
440   disease for each term. An example of the format of the Acute decompensated heart failure
441   validation sheet for Wikipedia is shown in Fig. 8.
442
443   It is possible to note that an additional column was also present, called RELEVANT, and which
444   synthesises all the information available about the relevance of a term to a disease. The possible
445   values of this column are defined as:
446
447       (1) RELEVANT = **YES**. If (WIKIPEDIA = YES) & (METAMAP = YES) & (TVP = (YES
448           or NO)), that is, it is considered to be a valid medical concept for the diagnosis of a
449           disease.
450       (2) RELEVANT = **NO**. If (WIKIPEDIA = YES) & (METAMAP = YES) & (TVP = NO),
451           that is, it is considered to be a medical concept that is nonspecific, and thus too general to
452           be helpful in the diagnosis of a disease.
453       (3) RELEVANT = **FPREAL**. If (WIKIPEDIA = YES) & (METAMAP = YES) & (TVP =
454           YES). The term **is not relevant** because it is considered to be a nonspecific, general
455           concept that does not make sense for diagnosis, even though Metamap has detected it and
456           the TVP process has evaluated it as a diagnostic term. For example, in an excerpt from

457         Acute decompensated heart disease on Wikipedia: "*Other cardiac symptoms of heart*

458         *failure include chest pain/pressure and palpitations…*", Metamap has detected **Chest**

459         **pain** and **Pain** from "*chest pain*", both were marked as TRUE by TVP but the concept

460         dismissed by nonspecific and general was Pain.

461      (4) RELEVANT = **FPCONTEXT**. If (WIKIPEDIA = YES) & (METAMAP = YES) &

462         (TVP = YES). The term **is not relevant** because it is outside the diagnostic context, even

463         though Metamap has detected it and the TVP process has evaluated it as a diagnostic

464         term. In other words, this term has been obtained from texts whose content is outside the

465         diagnostic context. For example, in an excerpt from *Acute decompensated heart failure*

466         disease on Wikipedia: "*Other well recognized precipitating factors include anemia and*

467         *hyperthyroidism…*", Metamap has detect **Anemia** and **Hyperthyroidism** which are

468         medical terms but in context we dismiss them because they are risk factors for that

469         disease.

470      (5) RELEVANT = **FN**. If (WIKIPEDIA = YES) & (METAMAP = NO) & (TVP = NO).

471         These terms were manually detected in the texts, but Metamap failed in recognising

472         them.

473

474 The cases (3) and (4) above define situations in which the detected term is esteemed to be of no

475 relevance, and as such represent cases of false positives. It is nevertheless necessary to

476 discriminate the reason behind such error, which can be because: i) the term is a very general,

477 nonspecific concept whose definition does not represent and contributes nothing to the diagnosis

478 (FP_REAL), or ii) because the term is a medical term that is out of place with respect to the

479 context that is narrated in the text – in other words, it could be a valid diagnostic term but not for

480 the disease that is under validation or in the context in which have been described and therefore

481 should be discarded (FP_CONTEXT).

482

483 Using this information for all diseases and terms, true positive (**TP**), false positive (**FP**), true

484 negative (**TN**) and false negative (**FN**) rates were computed in order to calculate precision, recall

485 and F1 score values as metrics to measure the performance of DISNET system. The mean values

486 for these parameters are depicted in **Error! Reference source not found.**. The **TP** is all terms

487 with (WIKIPEDIA = YES) & (METAMAP = YES) & (TVP = YES) & (RELEVANT = YES).

488 As previously explained, the **FP** is composed of two parts, being the total FP the sum of

489 **FP_REAL** + **FP_CONTEXT**:

490

491     •   **FP_REAL** = (WIKIPEDIA = YES) & (METAMAP = YES) & (TVP = YES) &

492        (RELEVANT = FPREAL).

493     •   **FP_CONTEXT** = (WIKIPEDIA = YES) & (METAMAP = YES) & (TVP = YES) &

494        (RELEVANT = FPCONTEXT).

495

496 **FN** is also composed of two parts, i.e. **FN_METAMAP** + **FN_TVP**.

497

- **FN_METAMAP** = (WIKIPEDIA = YES) & (METAMAP = NO) & (TVP = NO) &
  (RELEVANT = FN). These are terms that Metamap has not found.
- **FN_TVP** = (WIKIPEDIA = YES) & (METAMAP = YES) & (TVP = NO) &
  (RELEVANT = YES). These are terms that TVP has validated as false while being
  relevant.

503

504 Finally, the **TN** measures the TVP process (WIKIPEDIA = YES) & (METAMAP = YES) &
505 (TVP = NO) & (RELEVANT = NO).  In the Table 1 are reported the values obtained for
506 Wikipedia and PubMed.

507

508 Detailed results for each disease are available online, for Wikipedia[21] and for PubMed[22],
509 including the list of terms manually extracted from the relevant texts of the articles, the matching
510 with the list of terms provided by Metamap, the result of the TVP process for each term and the
511 value of relevance as annotated by our researchers.

512

513 Results indicate that our NLP (Metamap + TVP) process is sufficiently reliable, with an accuracy
514 of 0.731 (confidence interval of [0.710, 0.753], calculated through a Wilson's score interval with
515 continuity correction and a confidence level of 99%) for Wikipedia and of 0.640 (confidence
516 interval of: [0.606, 0.680]) for PubMed (**Error! Reference source not found.**). The results of
517 the calculations of these parameters for each disease can be viewed online for Wikipedia[23] and
518 for each abstract in PubMed[24].

519

520 About the results for **FP** presented in Table 1, we can say that they are mainly due to the
521 configuration used for Metamap for the extraction of terms, extended in successive extractions to
522 avoid leaving out terms that are relevant for the detection of diseases.

523

524 Thus, one of the last extensions in the search terms added the semantic types Mental or
525 Behavioral Dysfunction and Intellectual Product; thanks to this extension, important symptoms
526 have been detected for certain diseases, which were not detected before, such as: *Anxiety*,
527 *Bulimia*, *Anorexy*, *Stress*, etc. We believe that it is better to discard those terms that are not
528 relevant than to omit those that are relevant to a disease.

529

530 It is further interesting to observe the large difference in the false positive rates between
531 Wikipedia (11.41%) and PubMed (17.54%). We speculate that this is due to the concretion of
532 articles. Accordingly, in Wikipedia, articles referring to one disease refer almost exclusively to
533 that particular disease, and thus include no irrelevant terms – with a few exceptions related to

---

[21] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/tree/master/wikipedia_validation_sheets
[22] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/tree/master/pubmed_validation_sheets
[23] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/wikipedia_individual_validation_results.csv
[24] https://midas.ctb.upm.es/gitlab/disnet/paperdisnet/blob/master/pubmed_individual_validation_results.csv

534    differential diagnoses. Nevertheless, this is not the case of PubMed articles as a significant part
535    of them are not so specific. Many are the articles describing real medical cases, where the
536    symptoms are those displayed by a given patient, plus others referring to congenital diseases of
537    the patient, or even diseases that he/she previously possessed. Consequently, the same PubMed
538    article includes symptoms of many different diseases that, although being true medical terms and
539    thus being recognized by Metamap, are not relevant to the disease under analysis.
540
541    For **TN,** we must also take into account that most of the terms extracted by Metamap as relevant
542    have been purged by TVP, which has been in charge of determining which terms are relevant
543    and which are not, so that the vast majority of terms extracted by Metamap that are not relevant
544    to the disease have been classified in this way by TVP (35.78% for Wikipedia and 32.84% for
545    PubMed).
546
547    In addition, we have observed that most of the true negative terms in both Wikipedia and
548    PubMed are constant, and include: *indicated*, *syndrome*, *disease*, *illness*, *infected*, *sing*,
549    *symptoms*, *used to*, etc.
550
551    Finally, **FN** are those terms that are relevant to the disease in question, but that have not been
552    detected by Metamap; note that these have been manually extracted for the validation process.
553    The vast majority of **FN** are formed by complex expressions of the language, so their detection is
554    challenging for any NLP tool. We can further observe that the difference in the ratio of false
555    negative between Wikipedia (21.68%) and PubMed (18.40%) is 3.28%. We believe that this
556    difference is mainly due to the forms of expression used in both sources, with Wikipedia being
557    more discursive, as opposed to the scientific style of PubMed.
558
559    In synthesis, we can conclude that a clear relationship can be observed between the performance
560    of the system and the nature of the underlying data source. Specifically, while PubMed is an
561    exclusively medical source, created, written and edited by specialists in the field, Wikipedia is a
562    source of public information, written by anyone who has access to the web, so that the articles in
563    it contained can be written by medical students or just users with some knowledge in the field,
564    whose expressions cannot be assimilated to those of specialists who write PubMed. Considering
565    that the tool used by DISNET for the extraction of medical terms (Metamap) is a medical tool, it
566    is not surprising that it displays a greater capacity for the recognition of medical terms, as
567    opposed to more colloquial terms formed by more complex phrases; thus, there are terms such as
568    "*Swollen lymph glads under the jaw*", or "*sensation of swelling in the area of the larynx*", that
569    Metamap cannot recognize.
570
571    It is true that the validation percentages do not seem very high, but we must take into account the
572    following facts, firstly, that there is no other system that extracts and generates phenotypic
573    information using an approach as proposed in this document and secondly, the objective of the

574 document is not clinical, but purely research, and thus allows all the knowledge generated to be
575 put within the reach of other researchers and for the scientific community in general. Therefore,
576 the use of DISNET medical information is in the hands of all types of people and they are
577 therefore responsible for the use they give to such data. It is also important to mention that
578 despite the complex and inherent nature of the texts from different sources, the percentages
579 reflect good research work.
580

## A use case

582 To illustrate the possible use of the DISNET system, we here present a simple use case, which
583 consists of the creation of several basic DISNET queries, and the visualization of the
584 corresponding results.
585

586 The DISNET API has the capacity to create a variety of queries and in this section only a couple
587 of queries have been created in order to provide a small example of the capacity to support
588 research into the proposed system.
589

### Creation of DISNET queries

591 For the sake of simplicity, we will here focus on two of the most important characteristics of
592 DISNET: **i)** the ability to create relationships between diseases according to their phenotypic
593 similarity (**C1**) and **ii)** the ability to increase/improve the phenotypic information of diseases by
594 means of periodic extractions of knowledge (**C2**).
595

596 The scenario C1 implies obtaining data for two diseases, which we suspect may share symptoms;
597 we will here use "Influenza" and "Gastroenteritis". The resulting DISNET queries are:
598

(1) disnet.ctb.upm.es/api/disnet/query/**disnetConceptList**?source=**wikipedia**&version=**20
18-08-15**&diseaseName=**Influenza**&matchExactName=**true**
(2) disnet.ctb.upm.es/api/disnet/query/**disnetConceptList**?source=**pubmed**&version=**2018
-04-03**&diseaseName=**Influenza**&matchExactName=**true**
(3) disnet.ctb.upm.es/api/disnet/query/**disnetConceptList**?source=**wikipedia**&version=**20
18-08-15**&diseaseName=**Gastroenteritis**&matchExactName=**true**
(4) disnet.ctb.upm.es/api/disnet/query/**disnetConceptList**?source=**pubmed**&version=**2018
-04-03**&diseaseName=**Gastroenteritis**&matchExactName=**true**

608 We have here used the DISNET query "**disnetConcepList**", which allows retrieving the list of
609 "**DISNET Concepts**" associated with a given disease. The parameters of this query include:
610 "**diseaseName**", with the name of the disease; "**matchExactName**", to indicate that the search
611 by disease name is exact; and "**source**" and "**snapshot**", to respectively indicate the source and
612 snapshot we want to consult. In this case, we selected to consult the two sources Wikipedia and
613 PubMed, and respectively the snapshots of August 15[th], 2018 and April 3[rd], 2018. Note that the

614  result will consists of four total lists, two for each disease. To illustrate, Fig. 11 shows an extract
615  of the response from the query (3).

616

617  As for the scenario C2, it requires retrieving data for a disease whose list of symptoms may have
618  changed with time, i.e. either increased or decreased. As an example, we considered the disease
619  "Acrodynia", and executed the following DISNET queries:

620

621      (1) disnet.ctb.upm.es/api/disnet/query/**disnetConceptList**?source=**wikipedia**&version=**20**
622          **18-02-01**&diseaseName=**Acrodynia**&matchExactName=**true**
623      (2) disnet.ctb.upm.es/api/disnet/query/**disnetConceptList**?source=**wikipedia**&version=**20**
624          **18-02-15**&diseaseName=**Acrodynia**&matchExactName=**true**

625

626  Note that, as in C1, we have here used the query "**disnetConceptList**"; nevertheless, we have
627  here executed it twice, on the same disease (**Acrodynia**) and two different snapshots: February
628  1st, 2018 and February 15th, 2018.

629

630  **Visualization of the result of the DISNET queries**
631  Once the results of the query have been retrieved, the next natural step is their visualization;
632  while the actual output format may vary according to the needs of each specific project, for the
633  sake of clarity we here created a graph representation by using the external tool Cytoscape[25]. In
634  both scenarios (i.e. C1 and C2) we generated relationships between diseases and their symptoms,
635  with the aim of visualizing the value and scope of the medical data stored and processed by
636  DISNET. In Figure 10(b) we see the relationship between the Influenza and Gastroenteritis
637  diseases on one hand (highlighted in white rectangles), and the set of symptoms on the other.
638  Symptoms were obtained from two different sources, specifically Wikipedia and PubMed:
639  relationships are then respectively represented by red and blue edges. Common symptoms are
640  merged by the layout algorithm in the center of the graph; the medical terms that are not
641  common among the two diseases, on the contrary, form a peripheral shell. Note that "**Influenza**"
642  has 59 DISNET Concepts and "**Gastroenteritis**" has 47, 19 of which are in common.

643

644  In Figure 10(a) we observe the network representation of the disease "**Acrodynia**" and of its 18
645  medical terms, 15 of which were found in the snapshot of February 1st, 2018 and three new ones
646  in that of February 15th, 2018. This is thus an example of an increase in phenotypic knowledge.

647

648  This simple use case illustrates how the DISNET system allows generating a network of diseases
649  and their symptoms on a large scale, and that it provides the right environment to know how
650  diseases are related according to their phenotypic manifestations. By applying similarity
651  algorithms, such as Cosine (Zhou et al., 2014)(Li et al., 2014)(van Driel et al., 2006) or the
652  Jaccard index (Hoehndorf, Schofield & Gkoutos, 2015), it is possible to estimate the similarity

---

[25] http://www.cytoscape.org

653  between two diseases, and thus to focus further medical analyses on those pairs showing a large
654  overlap. These features will be also implemented as native features in next DISNET release.
655

## Discussion

657  This work presented the DISNET system, starting from its underlying conception, up to its
658  technical structure and data workflow. DISNET allows retrieving knowledge about the signs,
659  symptoms and diagnostic tests associated with a disease. It is not limited to a specific category
660  (all the categories that the selected sources of information offer us) and clinical diagnosis terms.
661  It further allows to track the evolution of those terms through time, being thus an opportunity to
662  analyse and observe the progress of human knowledge on diseases. We also presented the
663  DISNET REST API, which aims at sharing the retrieved information with the wide scientific
664  community. We further discussed the validation of the system, suggesting that it is good enough
665  to be used to extract diseases and diagnostically-relevant terms. At the same time, the evaluation
666  also revealed that improvements could be introduced to enhance the system's reliability.
667

## Conclusions

669  Among the potential lines of future works, priority will be given to increasing the number of
670  information sources, by including other websites like Medline Plus or CDC. Secondly, we are
671  considering the possibility of extending the TVP procedure, by adding new data sources, with the
672  aim of increasing the number of validation terms and hence of reducing the number of false
673  negatives. Note that this could also partly be achieved by resorting to a different NLP tool to
674  process the input texts, as for example to Apache cTakes (Savova et al., 2010). Other potential
675  options for future work are the improvement of the ambiguity of medical terms and the
676  implementation of tools that allow the representation of the knowledge extracted and generated.
677  Also, future implementations of DISNET also aim to provide ways to automatically compute the
678  similarity between diseases (by using already mentioned and well-known similarity metrics),
679  extending the DISNET platform to include biological and drug information and developing new
680  visualization strategies, among others.
681

## Acknowledgements

692

## References

Aibar E. 2017.La ciencia de la Wikipedia. *Available at https://metode.es/revistas-metode/article-revistes/la-ciencia-de-la-wikipedia.html* (accessed February 18, 2018).

Aich S, Sain M, Park J, Choi K, Kim H. 2017. A text mining approach to identify the relationship between gait-Parkinson's disease (PD) from PD based research articles. In: *2017 International Conference on Inventive Computing and Informatics (ICICI)*. Coimbatore, India: IEEE Computer Society, 481–485. DOI: 10.1109/ICICI.2017.8365398.

Aronson AR. 2001. Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program. *Proceedings. AMIA Symposium*:17–21.

Azer SA. 2014. Evaluation of gastroenterology and hepatology articles on Wikipedia: Are they suitable as learning resources for medical students? *European Journal of Gastroenterology & Hepatology* 26:155. DOI: 10.1097/MEG.0000000000000003.

Azzam A, Bresler D, Leon A, Maggio L, Whitaker E, Heilman J, Orlowitz J, Swisher V, Rasberry L, Otoide K, Trotter F, Ross W, McCue JD. 2017. Why Medical Schools Should Embrace Wikipedia: Final-Year Medical Student Contributions to Wikipedia Articles for Academic Credit at One School. *Academic Medicine* 92:194–200. DOI: 10.1097/ACM.0000000000001381.

Barabási A-L, Gulbahce N, Loscalzo J. 2010. Network medicine: a network-based approach to human disease. *Nature Reviews Genetics* 12:nrg2918. DOI: 10.1038/nrg2918.

Bodenreider O. 2004. The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Research* 32:D267–D270. DOI: 10.1093/nar/gkh061.

Botstein D, Risch N. 2003. Discovering genotypes underlying human phenotypes: past successes for mendelian disease, future approaches for complex disease. *Nature Genetics* 33:228–237. DOI: 10.1038/ng1090.

Bou Rjeily C, Badr G, Hajjarm El Hassani A, Andres E. 2019. Medical Data Mining for Heart Diseases and the Future of Sequential Mining in Medical Field. In: Tsihrintzis GA, Sotiropoulos DN, Jain LC eds. *Machine Learning Paradigms: Advances in Data Analytics*. Intelligent Systems Reference Library. Cham: Springer International Publishing, 71–99. DOI: 10.1007/978-3-319-94030-4_4.

Calderone A, Castagnoli L, Cesareni G. 2013. mentha: a resource for browsing integrated protein-interaction networks. *Nature Methods* 10:690–691. DOI: 10.1038/nmeth.2561.

Chen J, Li K, Rong H, Bilal K, Yang N, Li K. 2018. A disease diagnosis and treatment recommendation system based on big data mining and cloud computing. *Information Sciences* 435:124–149. DOI: 10.1016/j.ins.2018.01.001.

Chen Y, Zhang X, Zhang G, Xu R. 2015. Comparative analysis of a novel disease phenotype network based on clinical manifestations. *Journal of Biomedical Informatics* 53:113–120. DOI: 10.1016/j.jbi.2014.09.007.

Cohen N. 2013. Editing Wikipedia Pages for Med School Credit. *The New York Times*.

Dias G, Oliveira JL, Vicente F-J, Martín-Sánchez F. 2005. Integration of Genetic and Medical Information Through a Web Crawler System. In: *Biological and Medical Data Analysis*. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 78–88. DOI: https://doi.org/10.1007/11573067_9.

van Driel MA, Bruggeman J, Vriend G, Brunner HG, Leunissen JAM. 2006. A text-mining analysis of the human phenome. *European journal of human genetics: EJHG* 14:535–542. DOI: 10.1038/sj.ejhg.5201585.

739  Espe S. 2018. Malacards: The Human Disease Database. *Journal of the Medical Library*
740        *Association : JMLA* 106:140–141. DOI: 10.5195/jmla.2018.253.
741  Friedlin J, McDonald CJ. 2010. An evaluation of medical knowledge contained in Wikipedia and
742        its use in the LOINC database. *Journal of the American Medical Informatics Association:*
743        *JAMIA* 17:283–287. DOI: 10.1136/jamia.2009.001180.
744  Frunza O, Inkpen D, Tran T. 2011. A Machine Learning Approach for Identifying Disease-
745        Treatment Relations in Short Texts. *IEEE Transactions on Knowledge and Data*
746        *Engineering* 23:801–814. DOI: 10.1109/TKDE.2010.152.
747  García del Valle EP, Lagunes García G, Prieto Santamaría L, Zanin M, Menasalvas Ruiz E,
748        Rodríguez-González A. 2019. Disease networks and their contribution to disease
749        understanding: A review of their evolution, techniques and data sources. *Journal of*
750        *Biomedical Informatics* 94:103206. DOI: 10.1016/j.jbi.2019.103206.
751  Goh K-I, Cusick ME, Valle D, Childs B, Vidal M, Barabási A-L. 2007. The human disease
752        network. *Proceedings of the National Academy of Sciences* 104:8685–8690. DOI:
753        10.1073/pnas.0701361104.
754  Gupta S, Dingerdissen H, Ross KE, Hu Y, Wu CH, Mazumder R, Vijay-Shanker K. 2018.
755        DEXTER: Disease-Expression Relation Extraction from Text. *Database* 2018. DOI:
756        10.1093/database/bay045.
757  Heilman JM, West AG. 2015. Wikipedia and Medicine: Quantifying Readership, Editors, and
758        the Significance of Natural Language. *Journal of Medical Internet Research* 17. DOI:
759        10.2196/jmir.4069.
760  Hoehndorf R, Schofield PN, Gkoutos GV. 2015. Analysis of the human diseasome using
761        phenotype similarity between common, genetic, and infectious diseases. *Scientific*
762        *Reports* 5:srep10888. DOI: 10.1038/srep10888.
763  Hornbeck PV, Zhang B, Murray B, Kornhauser JM, Latham V, Skrzypek E. 2015.
764        PhosphoSitePlus, 2014: mutations, PTMs and recalibrations. *Nucleic Acids Research*
765        43:D512–D520. DOI: 10.1093/nar/gku1267.
766  Hwang S, Kim CY, Yang S, Kim E, Hart T, Marcotte EM, Lee I. 2019. HumanNet v2: human
767        gene networks for disease research. *Nucleic Acids Research* 47:D573–D580. DOI:
768        10.1093/nar/gky1126.
769  Information NC for B, Pike USNL of M 8600 R, MD B, Usa 20894. 2019. *PubMed Help*.
770        National Center for Biotechnology Information (US).
771  Jia J, An Z, Ming Y, Guo Y, Li W, Li X, Liang Y, Guo D, Tai J, Chen G, Jin Y, Liu Z, Ni X, Shi
772        T. 2018. PedAM: a database for Pediatric Disease Annotation and Medicine. *Nucleic*
773        *Acids Research* 46:D977–D983. DOI: 10.1093/nar/gkx1049.
774  Kanehisa M, Goto S. 2000. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids*
775        *Research* 28:27–30. DOI: 10.1093/nar/28.1.27.
776  Köhler S, Vasilevsky NA, Engelstad M, Foster E, McMurry J, Aymé S, Baynam G, Bello SM,
777        Boerkoel CF, Boycott KM, Brudno M, Buske OJ, Chinnery PF, Cipriani V, Connell LE,
778        Dawkins HJS, DeMare LE, Devereau AD, de Vries BBA, Firth HV, Freson K, Greene D,
779        Hamosh A, Helbig I, Hum C, Jähn JA, James R, Krause R, F. Laulederkind SJ,
780        Lochmüller H, Lyon GJ, Ogishima S, Olry A, Ouwehand WH, Pontikos N, Rath A,
781        Schaefer F, Scott RH, Segal M, Sergouniotis PI, Sever R, Smith CL, Straub V,
782        Thompson R, Turner C, Turro E, Veltman MWM, Vulliamy T, Yu J, von Ziegenweidt J,
783        Zankl A, Züchner S, Zemojtel T, Jacobsen JOB, Groza T, Smedley D, Mungall CJ,

784         Haendel M, Robinson PN. 2017. The Human Phenotype Ontology in 2017. *Nucleic Acids*
785         *Research* 45:D865–D876. DOI: 10.1093/nar/gkw1039.
786  Lee I, Blom UM, Wang PI, Shim JE, Marcotte EM. 2011. Prioritizing candidate disease genes by
787         network-based boosting of genome-wide association data. *Genome Research* 21:1109–
788         1121. DOI: 10.1101/gr.118992.110.
789  Li X, Zhou X, Peng Y, Liu B, Zhang R, Hu J, Yu J, Jia C, Sun C. 2014. Network Based
790         Integrated Analysis of Phenotype-Genotype Data for Prioritization of Candidate
791         Symptom Genes. *BioMed Research International* 2014:435853. DOI:
792         10.1155/2014/435853.
793  Lipscomb CE. 2000. Medical Subject Headings (MeSH). *Bulletin of the Medical Library*
794         *Association* 88:265–266.
795  Lo Surdo P, Calderone A, Iannuccelli M, Licata L, Peluso D, Castagnoli L, Cesareni G, Perfetto
796         L. 2018. DISNOR: a disease network open resource. *Nucleic Acids Research* 46:D527–
797         D534. DOI: 10.1093/nar/gkx876.
798  Lopes P, Oliveira JL. 2013. An innovative portal for rare genetic diseases research: The semantic
799         Diseasecard. *Journal of Biomedical Informatics* 46:1108–1115. DOI:
800         10.1016/j.jbi.2013.08.006.
801  Matheson D, Matheson C. 2017. Open Medicine Journal Wikipedia as Informal Self-Education
802         for Clinical Decision-Making in Medical Practice. *Open Medicine Journal* 4:1–25. DOI:
803         10.2174/1874220301704010015.
804  Mattingly CJ, Rosenstein MC, Davis AP, Colby GT, Forrest JN, Boyer JL. 2006. The
805         Comparative Toxicogenomics Database: A Cross-Species Resource for Building
806         Chemical-Gene Interaction Networks. *Toxicological Sciences* 92:587–595. DOI:
807         10.1093/toxsci/kfl008.
808  Mazumder R, Mahmood ASMA, Vijay-Shanker K, Wu T-J. 2016. DiMeX: A Text Mining
809         System for Mutation- Disease Association Extraction. *PLOS One* 11:e0152725. DOI:
810         10.1371/journal. pone.0152725.
811  Mehdi M, Okoli C, Mesgari M, Nielsen FÅ, Lanamäki A. 2017. Excavating the mother lode of
812         human-generated text: A systematic review of research that uses the wikipedia corpus.
813         *Information Processing & Management* 53:505–529. DOI: 10.1016/j.ipm.2016.07.003.
814  Oliveira JL, Dias G, Oliveira I, Rocha P, Hermosilla I, Vicente J, Spiteri I, Martin-Sánchez F,
815         Pereira AS. 2004. DiseaseCard: A Web-Based Tool for the Collaborative Integration of
816         Genetic and Medical Information. In: *Biological and Medical Data Analysis*. Lecture
817         Notes in Computer Science. Springer Berlin Heidelberg, 409–417. DOI:
818         https://doi.org/10.1007/978-3-540-30547-7_41.
819  Pérez-Rodríguez G, Pérez-Pérez M, Fdez-Riverola F, Lourenço A. 2019. Online visibility of
820         software-related web sites: The case of biomedical text mining tools. *Information*
821         *Processing & Management* 56:565–583. DOI: 10.1016/j.ipm.2018.11.011.
822  Perfetto L, Briganti L, Calderone A, Perpetuini AC, Iannuccelli M, Langone F, Licata L,
823         Marinkovic M, Mattioni A, Pavlidou T, Peluso D, Petrilli LL, Pirrò S, Posca D,
824         Santonico E, Silvestri A, Spada F, Castagnoli L, Cesareni G. 2016. SIGNOR: a database
825         of causal relationships between biological entities. *Nucleic Acids Research* 44:D548–
826         D554. DOI: 10.1093/nar/gkv1048.
827  Piñero J, Bravo À, Queralt-Rosinach N, Gutiérrez-Sacristán A, Deu-Pons J, Centeno E, García-
828         García J, Sanz F, Furlong LI. 2017. DisGeNET: a comprehensive platform integrating

829     information on human disease-associated genes and variants. *Nucleic Acids Research*
830        45:D833–D839. DOI: 10.1093/nar/gkw943.
831  Pletscher-Frankild S, Pallejà A, Tsafou K, Binder JX, Jensen LJ. 2015. DISEASES: Text mining
832        and data integration of disease–gene associations. *Methods* 74:83–89. DOI:
833        10.1016/j.ymeth.2014.11.020.
834  pubmeddev.Home - PubMed - NCBI. *Available at https://www.ncbi.nlm.nih.gov/pubmed/*
835        (accessed February 16, 2018).
836  Quwaider M, Alfaqeeh M. 2016. Social Networks Benchmark Dataset for Diseases
837        Classification. In: *2016 IEEE 4th International Conference on Future Internet of Things
838        and Cloud Workshops (FiCloudW)*. 234–239. DOI: 10.1109/W-FiCloud.2016.56.
839  Rao AJ, Rao RS. 2018. Review On Machine Learning Approach for Detecting Disease-
840        Treatment Relations in Short Texts. *International Journal of Scientific Research in
841        Computer Science, Engineering and Information Technology* 4:122–129. DOI:
842        10.32628/CSEIT1833616.
843  Rappaport N, Nativ N, Stelzer G, Twik M, Guan-Golan Y, Iny Stein T, Bahir I, Belinky F,
844        Morrey CP, Safran M, Lancet D. 2013. MalaCards: an integrated compendium for
845        diseases and their annotation. *Database* 2013. DOI: 10.1093/database/bat018.
846  Rappaport N, Twik M, Nativ N, Stelzer G, Bahir I, Stein TI, Safran M, Lancet D. 2014.
847        MalaCards: A Comprehensive Automatically-Mined Database of Human Diseases.
848        *Current Protocols in Bioinformatics* 47:1.24.1-1.24.19. DOI:
849        10.1002/0471250953.bi0124s47.
850  Robinson PN, Köhler S, Bauer S, Seelow D, Horn D, Mundlos S. 2008. The Human Phenotype
851        Ontology: A Tool for Annotating and Analyzing Human Hereditary Disease. *American
852        Journal of Human Genetics* 83:610–615. DOI: 10.1016/j.ajhg.2008.09.017.
853  Rodríguez González A, Costumero Moreno R, Martínez Romero M, Wilkinson MD, Menasalvas
854        Ruiz E. 2018. Extracting diagnostic knowledge from MedLine Plus: a comparison
855        between MetaMap and cTAKES Approaches. *Current Bioinformatics* 13:573–582. DOI:
856        10.2174/1574893612666170727094502.
857  Rodríguez-González A, Martínez-Romero M, Costumero R, Wilkinson MD, Menasalvas-Ruiz E.
858        2015. Diagnostic Knowledge Extraction from MedlinePlus: An Application for Infectious
859        Diseases. In: *9th International Conference on Practical Applications of Computational
860        Biology and Bioinformatics*. Advances in Intelligent Systems and Computing. Springer,
861        Cham, 79–87. DOI: 10.1007/978-3-319-19776-0_9.
862  Russell-Rose T, Chamberlain J, Azzopardi L. 2018. Information retrieval in the workplace: A
863        comparison of professional search practices. *Information Processing & Management*
864        54:1042–1057. DOI: 10.1016/j.ipm.2018.07.003.
865  Safran M, Solomon I, Shmueli O, Lapidot M, Shen-Orr S, Adato A, Ben-Dor U, Esterman N,
866        Rosen N, Peter I, Olender T, Chalifa-Caspi V, Lancet D. 2002. GeneCards™ 2002:
867        towards a complete, object-oriented, human gene compendium. *Bioinformatics* 18:1542–
868        1543. DOI: 10.1093/bioinformatics/18.11.1542.
869  Savova GK, Masanz JJ, Ogren PV, Zheng J, Sohn S, Kipper-Schuler KC, Chute CG. 2010.
870        Mayo clinical Text Analysis and Knowledge Extraction System (cTAKES): architecture,
871        component evaluation and applications. *Journal of the American Medical Informatics
872        Association : JAMIA* 17:507–513. DOI: 10.1136/jamia.2009.001560.

873  Schriml LM, Arze C, Nadendla S, Chang Y-WW, Mazaitis M, Felix V, Feng G, Kibbe WA.
874      2012. Disease Ontology: a backbone for disease semantic integration. *Nucleic Acids*
875      *Research* 40:D940–D946. DOI: 10.1093/nar/gkr972.
876  Shafee T, Masukume G, Kipersztok L, Das D, Häggström M, Heilman J. 2017. Evolution of
877      Wikipedia's medical content: past, present and future. *J Epidemiol Community Health*
878      71:1122–1129. DOI: 10.1136/jech-2016-208601.
879  Singhal A, Simmons M, Lu Z. 2016. Text mining for precision medicine: automating disease-
880      mutation relationship extraction from biomedical literature. *Journal of the American*
881      *Medical Informatics Association* 23:766–772. DOI: 10.1093/jamia/ocw041.
882  SPARQL Query Language for RDF. 2017. *Available at https://www.w3.org/TR/rdf-sparql-*
883      *query/* (accessed November 18, 2017).
884  Stenson PD, Mort M, Ball EV, Shaw K, Phillips AD, Cooper DN. 2014. The Human Gene
885      Mutation Database: building a comprehensive mutation repository for clinical and
886      molecular genetics, diagnostic testing and personalized genomic medicine. *Human*
887      *Genetics* 133:1–9. DOI: 10.1007/s00439-013-1358-4.
888  Sudeshna P, Bhanumathi S, Hamlin MRA. 2017. Identifying symptoms and treatment for heart
889      disease from biomedical literature using text data mining. In: *2017 International*
890      *Conference on Computation of Power, Energy Information and Commuincation*
891      *(ICCPEIC)*. 170–174. DOI: 10.1109/ICCPEIC.2017.8290359.
892  Szklarczyk D, Gable AL, Lyon D, Junge A, Wyder S, Huerta-Cepas J, Simonovic M, Doncheva
893      NT, Morris JH, Bork P, Jensen LJ, Mering C von. 2019. STRING v11: protein–protein
894      association networks with increased coverage, supporting functional discovery in
895      genome-wide experimental datasets. *Nucleic Acids Research* 47:D607–D613. DOI:
896      10.1093/nar/gky1131.
897  Tsumoto S, Kimura T, Iwata H, Hirano S. 2017. Mining Text for Disease Diagnosis. *Procedia*
898      *Computer Science* 122:1133–1140. DOI: 10.1016/j.procs.2017.11.483.
899  UniProt Consortium. 2014. Activities at the Universal Protein Resource (UniProt). *Nucleic Acids*
900      *Research* 42:D191-198. DOI: 10.1093/nar/gkt1140.
901  Westergaard D, Stærfeldt H-H, Tønsberg C, Jensen LJ, Brunak S. 2018. A comprehensive and
902      quantitative comparison of text-mining in 15 million full-text articles versus their
903      corresponding abstracts. *PLOS Computational Biology* 14:e1005962. DOI:
904      10.1371/journal.pcbi.1005962.
905  Wikipedia:Manual of Style/Medicine-related articles. 2018. *Wikipedia*.
906  Xia E, Sun W, Mei J, Xu E, Wang K, Qin Y. 2018. Mining Disease-Symptom Relation from
907      Massive Biomedical Literature and Its Application in Severe Disease Diagnosis. *AMIA*
908      *Annual Symposium Proceedings* 2018:1118–1126.
909  Xu D, Zhang M, Xie Y, Wang F, Chen M, Zhu KQ, Wei J. 2016. DTMiner: identification of
910      potential disease targets through biomedical literature mining. *Bioinformatics* 32:3619–
911      3626. DOI: 10.1093/bioinformatics/btw503.
912  Zanzoni A, Soler-López M, Aloy P. 2009. A network medicine approach to human disease.
913      *FEBS Letters* 583:1759–1765. DOI: 10.1016/j.febslet.2009.03.001.
914  Zhao N, Zheng G, Li J, Zhao H, Lu C, Jiang M, Zhang C, Guo H, Lu A. 2018. Text Mining of
915      Rheumatoid Arthritis and Diabetes Mellitus to Understand the Mechanisms of Chinese
916      Medicine in Different Diseases with Same Treatment. *Chinese Journal of Integrative*
917      *Medicine* 24:777–784. DOI: 10.1007/s11655-018-2825-x.

918    Zhou X, Menche J, Barabási A-L, Sharma A. 2014. Human symptoms-disease network. *Nature*
919         *Communications* 5:4212. DOI: 10.1038/ncomms5212.
920
921

**Table 1**(on next page)

Total values from the February 1$^{st}$, 2018 snapshot of Wikipedia and the April 3$^{rd}$, 2018 snapshot of PubMed

1       **Table 1.** Total values from the February 1st, 2018 snapshot of Wikipedia and the April 3rd, 2018 snapshot of PubMed

| Parameter | Wikipedia | PubMed |
|---|---|---|
| TP | (31.11%) 2,075 | (31.20%) 724 |
| FP | (11.41%) 761 | (17.54%) 407 |
| FPREAL | 279 | 107 |
| FPCONTEXT | 482 | 300 |
| TN | (35.78%) 2,386 | (32.84%) 762 |
| FN | (21.68%) 1,446 | (18.40%) 427 |
| FN_METAMAP | 709 | 201 |
| FN_TVP | 737 | 226 |
| TOTAL | (100%) 6,668 | (100%) 2,320 |
| PRECISION | 0.731 | 0.640 |

2

3

# Figure 1

External vocabularies in a vertical infobox in Wikipedia article on Ebstein's anomaly and Cholestasis

# Figure 2

External vocabularies in a horizontal infobox in Wikipedia article on Influenza and
Cancer

# Figure 3

Disease MeSH Term tree clasification

Diseases [C] ⊖
    Bacterial Infections and Mycoses [C01] ⊕
    Virus Diseases [C02] ⊕
    Parasitic Diseases [C03] ⊕
    Neoplasms [C04] ⊕
    Musculoskeletal Diseases [C05] ⊕
    Digestive System Diseases [C06] ⊕
    Stomatognathic Diseases [C07] ⊕
    Respiratory Tract Diseases [C08] ⊕
    Otorhinolaryngologic Diseases [C09] ⊕
    Nervous System Diseases [C10] ⊕
    Eye Diseases [C11] ⊕
    Male Urogenital Diseases [C12] ⊕
    Female Urogenital Diseases and Pregnancy Complications [C13] ⊕
    Cardiovascular Diseases [C14] ⊕
    Hemic and Lymphatic Diseases [C15] ⊕
    Congenital, Hereditary, and Neonatal Diseases and Abnormalities [C16] ⊕
    Skin and Connective Tissue Diseases [C17] ⊕
    Nutritional and Metabolic Diseases [C18] ⊕
    Endocrine System Diseases [C19] ⊕
    Immune System Diseases [C20] ⊕
    Disorders of Environmental Origin [C21] ⊕
    Animal Diseases [C22] ⊕
    Pathological Conditions, Signs and Symptoms [C23] ⊕
    Occupational Diseases [C24] ⊕
    Chemically-Induced Disorders [C25] ⊕
    Wounds and Injuries [C26] ⊕

# Figure 4

PubMed Text Extraction Procedure workflow

# Figure 5

DISNET Architecture/Workflow

# Figure 6

Basic database statistics (most common medical terms)

WIKIPEDIA

Pain: 42.912

Lesion: 29.481

Neoplasms: 26.165

Inflammation: 20.533

Malnutrition: 19.405

Magnetic resonance imaging: 22.441

Convulsions: 18.975

Ultrasonography: 17.704

Hemorrhage: 18.395

Headache: 17.508

PUBMED

Neoplasms: 13.804

Lesion: 11.261

Magnetic resonance imaging: 6.081

Malnutrition: 5.784

Inflammation: 4.026

Infection: 3.567

Pain: 3.473

Ultrasonography: 3.115

Biopsy: 3.130

Hemorrhage: 2.760

# Figure 7

Basic database statistics (diseases with more validated medical terms. Comparison of PubMed and Wikipedia)

## In Wikipedia



Bar chart titled "In Wikipedia" with y-axis labeled TERMS (0–120). Values from left to right:

- Kawasaki disease: 98
- Cerebral palsy: 77
- Hypoglycemia: 76
- Anorexia nervosa: 75
- Crohn's disease: 75
- Heart failure: 75
- Dementia: 74
- Headache: 74
- Dementia with Lewy...: 73
- Sarcoidosis: 73
- Hepatitis: 72
- Lead poisoning: 71
- Cirrhosis: 71
- Nephrotic syndrome: 68
- Hyperthyroidism: 65
- Psychosis: 63
- Uremia: 62
- Behçet's disease: 60
- Attention deficit...: 60
- Hypertension: 60

## In PubMed



Bar chart titled "In PubMed" with y-axis labeled TERMS (0–140). Values from left to right:

- Hypercalcemia: 116
- Cranial nerve palsy: 110
- Beriberi: 109
- Lipodystrophy: 105
- Subdural empyema: 105
- Brain disease: 104
- Hypervitaminosis A: 102
- Mitochondrial...: 101
- Locked-in syndrome: 101
- Noonan syndrome: 101
- Sarcoidosis: 100
- Granulomatous...: 99
- Costello syndrome: 99
- CREST syndrome: 99
- Brain stem infarction: 98
- Diabetes insipidus: 98
- Cardiac tamponade: 97
- Fanconi syndrome: 97
- Relapsing...: 96
- Pericarditis: 96

# Figure 8

Disease Acute decompensated heart failure sheet validation from the Wikipedia snapshot of February 1st, 2018

# Acute decompensated heart failure

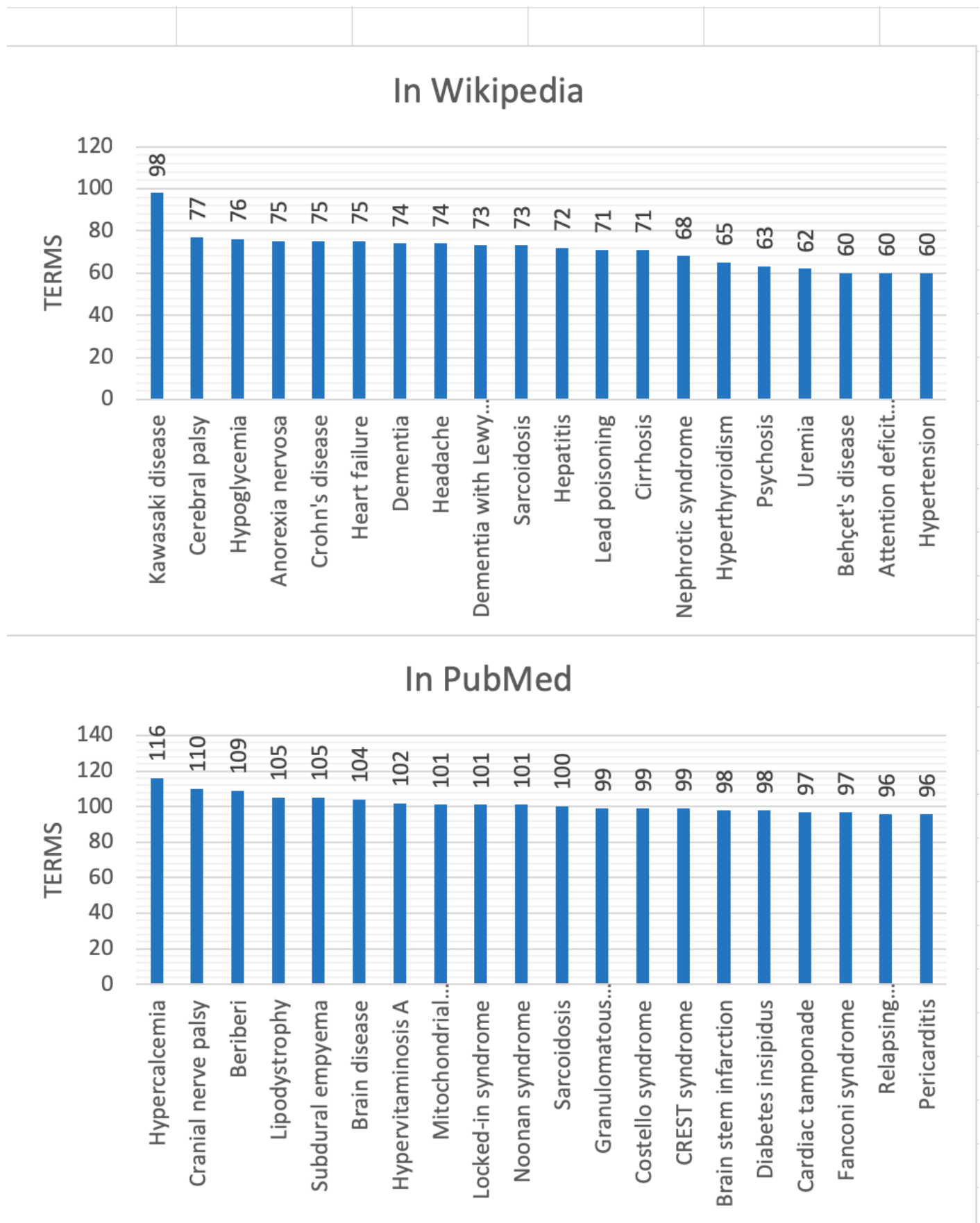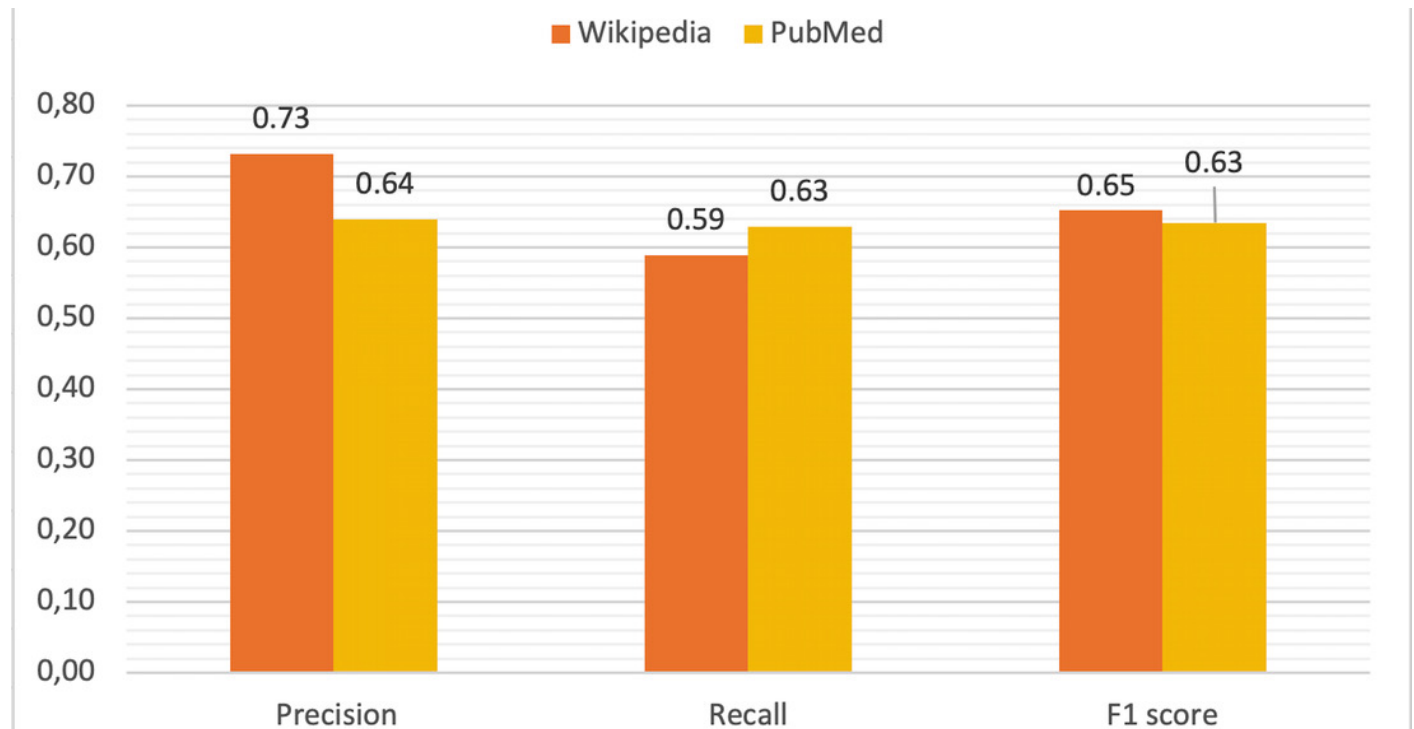| WIKIPEDIA TERMS | METAMAP TERMS | DISNET VALIDATION | | | |
|---|---|---|---|---|---|
| NAME | NAME | WIKIPEDIA | METAMAP | TVP | RELEVANT |
| 1 acute, myocardial, infarction | Acute myocardial infarction | YES | YES | YES | FPCONTEXT |
| 2 illness | Illness (finding) | YES | YES | YES | FPREAL |
| 3 hyperthyroidism | Hyperthyroidism | YES | YES | YES | FPCONTEXT |
| 4 anemia | Anemia | YES | YES | YES | FPCONTEXT |
| 5 weightloss | Weight decreased | YES | YES | YES | YES |
| 6 palpitations | Palpitations | YES | YES | YES | YES |
| 7 nausea | Nausea | YES | YES | YES | YES |
| 8 chest, pain | Chest pain NOS | YES | YES | YES | YES |
| 9 exertional, dyspnoea | Dyspnea on exertion | YES | YES | YES | YES |
| 10 pneumonia | Pneumonia | YES | YES | YES | FPCONTEXT |
| 11 high, blood, pressure | Hypertensive disease | YES | YES | YES | FPCONTEXT |
| 12 weakness | Weakness | YES | YES | YES | YES |
| 13 pain | Pain | YES | YES | YES | FPREAL |
| 14 heart, failure | Heart failure | YES | YES | YES | FPCONTEXT |
| 15 paroxysmal, nocturnal, dyspnoea | Paroxysmal nocturnal dyspnea | YES | YES | YES | YES |
| 16 orthopnoea | Orthopnea | YES | YES | YES | YES |
| 17 difficulty, breathing | Dyspnea | YES | YES | YES | YES |
| 18 heart, attack | Myocardial infarction, NOS | YES | YES | YES | FPCONTEXT |
| 19 abnormal, heart, rhythms | Cardiac arrhythmia | YES | YES | YES | FPCONTEXT |
| 20 bloating | Abdominal bloating | YES | YES | YES | YES |
| 21 chest, pressure | Pressure in chest | YES | YES | YES | YES |
| 22 low, urine, output | Oliguria | YES | YES | YES | YES |
| 23 fatigue | Fatigue | YES | YES | YES | YES |
| 24 jugular, venous, distension | Jugular venous engorgement | YES | YES | YES | YES |
| 25 atrial, fibrillation | Electrocardiographic atrial fibrillatio | YES | YES | NO | NO |
| 26 left, ventricular, failure | Left-sided heart failure | YES | YES | NO | NO |
| 27 sign, signs | Physical finding | YES | YES | NO | NO |
| 28 excess, fluid | Fluid overload | YES | YES | NO | NO |
| 29 chronic, heart, failure | Chronic heart failure | YES | YES | NO | NO |
| 30 pressure | Pressure (finding) | YES | YES | NO | NO |
| 31 acute, heart, failure | Acute heart failure | YES | YES | NO | NO |
| 32 myocardial, infarction | Electrocardiogram: myocardial infarction (finding) | YES | YES | NO | NO |
| 33 decompensation | Decompensation | YES | YES | NO | NO |
| 34 gasping | Gasping for breath | YES | YES | NO | NO |
| 35 symptom, symptoms | Symptom | YES | YES | NO | NO |
| 36 confusion | Confusion | YES | YES | YES | YES |
| 37 fluid, retention | Body fluid retention | YES | YES | YES | FPCONTEXT |
| 38 memory, impairment | Memory impairment | YES | YES | YES | YES |
| 39 sensitive | Hypersensitivity | YES | YES | NO | NO |
| 40 anxiety | Anxiety | YES | YES | YES | YES |
| Acute pulmonary edem | | YES | NO | NO | FN |
| loss of appetite | | YES | NO | NO | FN |
| waking up at night to urinate | | YES | NO | NO | FN |
| cerebral symptoms | | YES | NO | NO | FN |

# Figure 9

Validation metrics comparative

# Figure 10

a) Network of graphs representing the evolution of phenotypic knowledge in Wikipedia and b) Network of graphs representing similar medical terms between two diseases
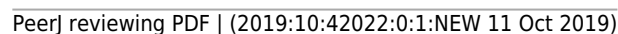
# Figure 11

Answer to the DISNET query "disnetConcepList" C1.(1)

```
        "diseaseId": "DIS006504",
        "name": "Influenza",
        "url": "http://en.wikipedia.org/wiki/Influenza",
        "disnetConceptsCount": 38,
        "disnetConceptList": [
            {
                "cui": "C0009443",
                "name": "Common cold",
                "semanticTypes": [
                    "dsyn"
                ]
            },
            {
                "cui": "C0010200",
                "name": "Coughing",
                "semanticTypes": [
                    "sosy"
                ]
            },
            {
                "cui": "C0027424",
                "name": "Nasal congestion (finding)",
                "semanticTypes": [
                    "sosy"
                ]
            },
            {
                "cui": "C0231221",
                "name": "Asymptomatic",
                "semanticTypes": [
                    "fndg"
                ]
            },
            {
                "cui": "C0015967",
                "name": "Fever",
                "semanticTypes": [
                    "fndg"
                ]
            },
            {
                "cui": "C0030193",
                "name": "Pain",
                "semanticTypes": [
                    "sosy"
                ]
            },
            {
                "cui": "C0085593",
                "name": "Chills",
                "semanticTypes": [
                    "sosy"
                ]
            },
            {
                "cui": "C0231218",
                "name": "Malaise",
                "semanticTypes": [
                    "sosy"
                ]
            },
```