

1 Semantic Annotation in Biomedicine:
2 the Current Landscape

3
4 Jelena Jovanović¹, Ebrahim Bagheri²

5
6 ¹Department of Software Engineering, University of Belgrade, 154 Jove Ilica Street, Serbia

7 ²Department of Electrical Engineering, Ryerson University, 245 Church Street, Toronto,
8 Canada

9 jeljov@fon.bg.ac.rs, bagheri@ryerson.ca

10
11
12
13 Corresponding author: Ebrahim Bagheri

19

20

21 **Abstract**

22 The abundance and unstructured nature of biomedical texts, be it clinical or research
23 content, impose significant challenges for the effective and efficient use of information and
24 knowledge stored in such texts. Annotation of biomedical documents with machine intelligible
25 semantics facilitates advanced, semantics-based text management, curation, indexing, and
26 search. This paper focuses on annotation of biomedical entity mentions with concepts from
27 relevant biomedical knowledge bases such as UMLS. As a result, the meaning of those
28 mentions is unambiguously and explicitly defined, and thus made readily available for
29 automated processing. This process is widely known as semantic annotation, and the tools
30 that perform it are known as semantic annotators.

31 Over the last dozen years, the biomedical research community has invested significant
32 efforts in the development of biomedical semantic annotation technology. Aiming to establish
33 grounds for further developments in this area, we review a selected set of state of the art
34 biomedical semantic annotators, focusing particularly on general purpose annotators, that is,
35 semantic annotation tools that can be customized to work with texts from any area of
36 biomedicine. We also examine potential directions for further improvements of today's
37 annotators which could make them even more capable of meeting the needs of real-world
38 applications. To motivate and encourage further developments in this area, along the
39 suggested and/or related directions, we review existing and potential practical applications
40 and benefits of semantic annotators.

41 Keywords

42 Natural Language Processing (NLP), Biomedical Ontologies, Semantic Technologies,
43 Biomedical Text Mining, Semantic Annotation

44 1. Background

45 Over the last few decades, huge volume of digital unstructured textual content have been
46 generated in biomedical research and practice, including a range of content types such as
47 scientific papers, medical reports, and physician notes. This has resulted in massive and
48 continuously growing collections of textual content that need to be organized, curated and
49 managed in order to be effectively used for both clinical and research purposes. Clearly,
50 manual curation and management of such “big” corpora are infeasible, and hence, the
51 biomedical community has long been examining and making use of various kinds of Natural
52 Language Processing (NLP) methods and techniques to, at least partially, facilitate their use.

53 In this paper, we focus on a specific NLP task, namely the extraction and disambiguation of
54 entities mentioned in biomedical textual content. Early efforts in biomedical information
55 extraction were devoted to Named Entity Recognition (NER), the task of recognizing specific
56 types of biomedical entities mentioned in text [1]. For instance, in the sentence “The patient
57 was diagnosed with upper respiratory tract infection”, a NER tool would recognize that the
58 phrase “respiratory tract infection” denotes a disease, but would not be able to determine
59 what particular disease it is. Semantic annotation, the NLP task of interest to this paper,
60 makes a significant advance, by not only recognizing the type of an entity, but also uniquely
61 linking it to its appropriate corresponding entry in a well-established knowledge base. In the
62 given example, a semantic annotator would not only recognize that the phrase “respiratory
63 tract infection” represents a disease, but would also identify what disease it is by connecting
64 the phrase with the concept C0035243 denoting ‘Respiratory Tract Infections’ from the UMLS
65 Metathesaurus (see Table 1). This way, the semantics of biomedical texts is made

66 accessible to software programs so that they can facilitate various laborious and time
67 consuming tasks such as search, classification, or organization of biomedical content.

68 While a suite of biomedical semantic annotation tools is available for practical use, the
69 biomedical community is yet to heavily engage in and leverage the benefits of such tools.

70 The goal of this paper is to introduce (i) some of the benefits and application use cases of
71 biomedical semantic annotation technology, (ii) a selection of the publicly available general
72 purpose semantic annotation tools for the biomedical domain, i.e., semantic annotators that
73 are not specialized for a particular biomedical entity type, but can detect and normalize
74 entities of multiple types in one pass, and (iii) potential areas where the work in the
75 biomedical semantic annotation domain can be strengthened or expanded. While the
76 overview of application cases and state of the art tools can be of relevance to practitioners in
77 the biomedical domain, with the summary of potential areas for further research, we are also
78 targeting researchers who are familiar with NLP, semantic technologies, and semantic
79 annotation in general, but have not been dealing with the biomedical domain, as well as
80 those who are well aware of biomedical semantic technologies, but have not been working
81 on semantic annotation. By providing researchers with an insight into the current state of the
82 art in biomedical semantic annotation in terms of the approaches and tools, as well as the
83 research challenges, we aim to offer them a basis for engagement with semantic annotation
84 technology within the biomedical domain and thus support even further developments in the
85 field.

86 The following section provides several examples of practical benefits achievable through
87 semantic annotation of biomedical texts (see also Table 2). The paper then examines the
88 available tool support, focusing primarily on general purpose biomedical annotators (Table
89 3). Still, considering the relevance and large presence of entity-specific biomedical
90 annotators, i.e., tools developed specifically for semantic annotation of a particular type of
91 biomedical entities such as genes or chemicals, we provide an overview of these tools, as
92 well. While examining the available tool support, we also consider biomedical knowledge

93 resources required for semantic annotation (Table 1), as well as resources used for
94 evaluating the tools' performance (Table 4). This is followed by a discussion of the
95 challenges that are preventing current semantic annotators from achieving their full potential.

96

97 **Table 1.** An overview of ontologies, thesauri and knowledge bases used by biomedical
98 semantic annotation tools discussed in the paper

99 [Table 1 should be placed approximately here]

100

101 2. Benefits and Use Cases

102 2.1 Better use of electronic medical record (EMR) in clinical practice

103 Electronic medical records (EMRs) are considered valuable source of clinical information,
104 ensuring effective and reliable information exchange among physicians and departments
105 participating in patient care, and supporting clinical decision making. However, EMRs largely
106 consist of unstructured, free-form textual content that require manual curation and analysis
107 performed by domain experts. A recent study examining the allocation of physician time in
108 ambulatory practice [2] confirmed the findings of previous similar studies (e.g. [3]), namely
109 that physicians spend almost twice as much time on the management of EMRs and related
110 desk work than on direct clinical face time with patients. Considering the inefficiency of
111 manual curation of EMRs, automation of the process is required if the potentials of EMRs are
112 to be exploited in clinical practice [4].

113 Semantic annotators provide the grounds for the required automation by extracting clinical
114 terms from free-form text of EMRs, and disambiguating the extracted terms with concepts of
115 a structured vocabulary, such as UMLS Metathesaurus. The identified concepts can be

116 subsequently used to search a repository of biomedical literature or evidence-based clinical
117 resources, in order to enrich EMRs with information pertinent to the patient's state. The
118 extracted terms can also be used for making summaries of clinical notes and articles [5].

119 2.2 Improved search and retrieval of resources for biomedical research

120 Publicly available biomedical data, tools, services, models and computational workflows
121 continuously increase in number, size and complexity. While this ever-growing abundance of
122 valuable resources opens up unprecedented opportunities for biomedical research, it is also
123 making it ever more challenging for researchers to efficiently discover and use the resources
124 required for accomplishing their tasks [6]. Hence, automation of the search and discovery
125 processes has turned into a necessity.

126 Clinical information stored in EMRs is also important in medical research, e.g., for
127 comparative effectiveness research, and epidemiological and clinical research studies [7].
128 Considering the unstructured nature of EMRs and their abundance, automated document
129 classification and information extraction methods are essential for assuring the efficiency and
130 effectiveness of search and retrieval of relevant information from large EMR collections.
131 Semantic annotation techniques can play a significant role in this regard. For instance, they
132 can be used to extract domain-specific concepts that could serve as discriminative features
133 for building automated classifiers of clinical documents [8]. Based on such classification,
134 clinical documents can be searched more effectively [7]. Furthermore, semantic concepts
135 extracted from biomedical literature can also be used for semantic indexing and retrieval of
136 biomedical publications [9] or biomedical tools and services [6]. In particular, biomedical
137 Information Retrieval systems use semantic annotators to expand the users' queries with
138 concepts and terms from vocabularies/ontologies (mapping the query text to the appropriate
139 ontology concepts, and then expanding the query with the terms associated with the mapped
140 concepts), as well as to classify the retrieved documents based on their content or the
141 occurrence of specific topics in the documents [1].

142 2.3 Disambiguation of abbreviations

143 Polysemous abbreviations are frequently present in biomedical literature and clinical texts
144 making it difficult for researchers and clinical practitioners to understand texts that are
145 outside the strict area of their expertise [10]. According to Chang et al. [11], in biomedical
146 journal articles, abbreviations with six or less characters have on average 4.61 possible
147 meanings. For instance, “ANA” has numerous possible meanings among which the most
148 frequent ones are “antinuclear antibodies”, “American Nurses Association”, “Alko Non-
149 Alcohol”, and “anandamide”.

150 Semantic annotators combined with general-purpose, machine-readable knowledge bases,
151 such as DBpedia (Table 1), can be used to disambiguate polysemous abbreviations and
152 unambiguously describe abbreviated terms based on the context in which they appear [10].
153 This can help researchers and practitioners better understand the meaning of such
154 abbreviations.

155 2.4 Seamless integration of data from disparate sources

156 Biomedical data are stored and maintained in disparate repositories. For instance, according
157 to the 2016 Molecular Biology Database Update, there are 1,685 biological databases [12].
158 This complicates the tasks of data management, retrieval and exploitation since one needs
159 to, first, locate the repositories that contain the required data; then, to familiarize oneself with
160 the meaning of the attributes and data types used in each repository; and, finally, learn how
161 to access and query the repositories [13].

162 For this reason, data integration can be highly useful for medical researchers. Jonquet et al.
163 [14] have nicely illustrated the need for seamless integration of data from various medical
164 sources: “a researcher studying the allelic variations in a gene would want to know all the
165 pathways that are affected by that gene, the drugs whose effects could be modulated by the
166 allelic variations in the gene, and any disease that could be caused by the gene, and the

167 clinical trials that have studied drugs or diseases related to that gene. The knowledge
168 needed to address such questions is available in public biomedical resources; the problem is
169 finding [and connecting] that information.”

170 Ontologies that are used to semantically annotate items in biomedical repositories allow for
171 weaving semantic links both within and across repositories thus establishing a semantic
172 network of biomedical items [13]. If both ontologies and resources they connect (through
173 semantic annotations) are in public domain, the resulting network takes the form of Linked
174 Open Data, as has already been shown in the Linked Life Data initiative (Table 1).

175 Table 2 provides a more structured view of some application cases for biomedical semantic
176 annotation technology.

177

178 **Table 2.** Example application cases of biomedical semantic annotation tools

179 [Table 2 should be placed approximately here]

180

181 3. Annotation Process, Tools and Resources

182 Biomedical texts have several characteristics that make them particularly challenging not
183 only for semantic annotation, but for any NLP task [6]. Some of these characteristics include:

- 184 i) Clinical text produced by practitioners often do *not fully adhere to correct grammar,*
185 *syntactic or spelling rules,* as the following triage note illustrates: “SORE THROAT pt
186 c/o sore throat x 1 week N pt states took antibiotic x 5 days after initiation of sore
187 throat and sx resolved and now back after completed antibiotics N pt tolerating po
188 fluids yet c/o pain on swallowing”;

- 189 ii) Biomedical terms are *often polysemous* and thus *prone to ambiguity*; for example, an
190 analysis of over 409K Medline abstracts revealed that 11.7% of the phrases were
191 ambiguous relative to the UMLS Metathesaurus [15].
- 192 iii) These textual corpora frequently use *abbreviations and acronyms that tend to be*
193 *polysemous* (see Section 2.3). In addition, clinical texts often contain *non-standard*
194 *shorthand phrases*, laboratory results and notes on patients' vital signs, which are
195 often filled with periods and thus can complicate typically straightforward text
196 processing tasks such as sentence splitting [16].
- 197 iv) Biomedical texts about or related to gene and protein mentions are particularly
198 challenging for semantic annotation. This is because *every protein* (e.g., SBP2), *has*
199 *an associated gene, often with the same name* [17]. Furthermore, *multiple genes*
200 *share symbols and names* (e.g. 'CAT' is the name of different genes in several
201 species, namely in cow, chicken, fly, human, mouse, pig, deer and sheep [18]).

202 To address these and other challenges of unstructured biomedical text, state-of-the-art
203 semantic annotators often rely on a combined use of text processing, large-scale knowledge
204 bases, semantic similarity measures and machine learning techniques [19]. In particular, in
205 the biomedical domain, semantic annotation is typically based on one of the following two
206 general approaches [20]: term-to-concept matching approach and approach based on
207 machine learning (ML) methods.

208 The term-to-concept matching approach, also referred to as dictionary lookup, is based on
209 matching specific segments of text to a structured vocabulary/dictionary or knowledge base
210 (e.g. UMLS or some of the OBO ontologies, Table 1). The drawback of some of the
211 annotators that implement this approach, e.g., NCBO Annotator [14] and ConceptMapper
212 [21], is the lack of disambiguation ability, meaning that the terms recognized in texts are
213 connected with several possible meanings, i.e., dictionary entries / concepts, instead of being
214 associated with a single meaning that is most appropriate for the given context. For example,
215 in the absence of disambiguation, in the sentence "In patients with DMD, the infiltration of

216 skeletal muscle by immune cells aggravates disease”, the term DMD would be associated
217 with several possible meanings, including Duchenne muscular dystrophy, dystrophin, and
218 DMD gene, whereas only the first one is correct for this given context.

219 The approaches based on ML methods are often found in annotators developed for specific,
220 well-defined application areas such as annotating drugs in medical discharge summaries [22]
221 or recognizing gene mentions in biomedical papers [23]. These annotators unambiguously
222 detect domain-specific concepts in text, and are typically highly performant on the specific
223 tasks they were developed for. However, as they are often based on supervised ML
224 methods, their development, namely, training of a ML model, requires large expert annotated
225 corpora, which are very expensive to develop. Another drawback of such annotators is that
226 they are only able to recognize specific categories of entities they are trained for, such as
227 genes or diseases, and cannot be applied to recognize concepts from broader vocabularies
228 [24]. The high costs associated with these approaches has led to a shift towards
229 unsupervised or semi-supervised ML methods that require few or no manually labelled data
230 [25]. Furthermore, several recent approaches have considered the idea of *distant supervision*
231 to generate ‘noisy’ labeled data for entity recognition [26] and entity typing [27].

232 3.1 Semantic Biomedical Annotation Tools

233 A large number of semantic annotation tools have been developed for the biomedical domain
234 [20, 24]. Many of them have resulted from research projects. Our focus in this paper is on a
235 subset of these tools that have the following characteristics:

- 236 • Semantic annotators that have been applied in practice or at least in research
237 projects other than those they originated from. In other words, we are not considering
238 research prototypes, but semantic annotators that have evolved from a research
239 prototype and have demonstrated their robustness for practical use.

240 • Semantic annotation tools that are available either as software libraries, web services
241 or web applications.

242 • General-purpose biomedical annotators, i.e., those semantic annotators that are not
243 tied to any particular biomedical task or entity type, but can be configured to work with
244 texts from different biomedical subdomains. This capacity originates from the fact that
245 they are either fully or at least partially grounded in the term-to-concept annotation
246 approach, which is flexible with respect to the annotation terminology.

247 Table 3 gives an overview of the semantic annotation tools that fulfilled the above given
248 criteria and thus were selected for inclusion in our study¹. The table compares the selected
249 tools with respect to several characteristic, including those related to the underlying
250 annotation method (configurability and disambiguation), the vocabulary (terminology) the tool
251 relies on, the tool's speed², its implementation aspects, and availability. The table also points
252 to some of the tools' specific features, which are further examined in the tool descriptions
253 given below.

254 As shown in Table 3 and further discussed below, all the tools are configurable in several
255 and often different ways, making it very difficult, if possible at all, to give a fair general
256 comparison of the tools. In other words, we believe that the only way to properly compare
257 these (and similar) annotation tools is in the context of a specific application case, where
258 each tool would be configured based on the application requirements. We expand on this in
259 Section 5.4 where we discuss the need for a benchmarking toolkit that would facilitate this
260 kind of application-specific tool benchmarking. Still, to offer some general insight into the
261 annotation capabilities of the selected tools, in Section 3.2 we briefly report on the
262 benchmarking studies that included several of the examined semantic annotators. In the
263 following, we introduce the selected semantic annotation tools and discuss their significant
264 features. The tools are presented in the order that corresponds to their order in Table 3.

265

266 **Table 3.** General purpose biomedical semantic annotation tools

267 [Table 3 should be placed approximately here]

268

269 ***Clinical Text Analysis and Knowledge Extraction System (cTAKES)*** [4] is a well-known
270 toolkit for semantic annotation of biomedical documents in general, and clinical research
271 texts in particular. It is built on top of two well-established and widely used open-source NLP
272 frameworks: Unstructured Information Management Architecture - UIMA [28] and OpenNLP
273 [29]. cTAKES is developed in a modular manner, as a pipeline consisting of several text
274 processing components that rely on either rule-based or ML techniques. Recognition of
275 concept mentions and annotation with the corresponding concept identifiers is done by a
276 component that implements a dictionary look-up algorithm. For building the dictionary,
277 cTAKES relies on UMLS. The concept recognition component does not resolve ambiguities
278 that result from identifying multiple concepts for the same text span. Disambiguation is
279 enabled through the integration of YTEX [7] in the cTAKES framework and its pipelines.
280 YTEX is a knowledge-based word sense disambiguation component that relies on the
281 knowledge encoded in UMLS. In particular, YTEX implements an adaptation of the Lesk
282 method [30], which scores candidate concepts for an ambiguous term by summing the
283 semantic relatedness between each candidate concept and the concepts in its context
284 window.

285 ***NOBLE Coder*** [20] is another open-source, general-purpose biomedical annotator. It can be
286 configured to work with arbitrary vocabularies. Besides enabling users to annotate
287 documents with existing vocabularies (terminologies), NOBLE Coder also provides them with
288 a Graphical User Interface where they can create custom terminologies by selecting one or
289 more branches from a set of existing vocabularies, and/or filtering vocabularies by semantic
290 types. It also allows for the dynamic change of the terminology (adding new concepts,
291 removing existing ones) while processing. The flexibility of this annotator also lies in the

292 variety of supported concept matching strategies, aimed at meeting the needs of different
293 kinds of NLP tasks. For example, the 'best match' strategy aims at high precision, and thus
294 returns few candidates (at most); as such, it is suitable for concept coding and information
295 extraction NLP tasks. The supported matching strategies allow for annotation of terms
296 consisting of single words, multiple words, and abbreviations. Thanks to its greedy algorithm,
297 NOBLE Coder can efficiently process large textual corpora. To disambiguate terms with
298 more than one associated concept, this tool relies on a set of simple heuristic rules such as
299 giving preference to candidates that map to a larger number of source vocabularies, or
300 candidates where the term is matched in its 'original' form, i.e., without being stemmed or
301 lemmatized.

302 **MetaMap** [31] is probably the most well-known and most widely used biomedical annotator.
303 It was developed by the U.S. National Library of Medicine. It maps biomedical entity
304 mentions of the input text to the corresponding concepts in the UMLS Metathesaurus. Each
305 annotation includes a score that reflects how well the concept matches the biomedical
306 term/phrase from the input text. The annotation process can be adapted in several ways by
307 configuring various elements of the annotation process such as the vocabulary used, the
308 syntactic filters applied to the input text, and the matching between text and concepts, to
309 name a few. Besides the flexibility enabled by these configuration options, another strong
310 aspect of MetaMap is its thorough and linguistically principled approach to the lexical and
311 syntactic analyses of input text. However, this thoroughness is also the cause of one of
312 MetaMap's main weaknesses, namely its long processing time, and thus its inadequacy for
313 annotating large corpora. Another weakness lies in its disambiguation approach which is not
314 able to effectively deal with ambiguous terms [32]. In particular, for disambiguation of terms,
315 MetaMap combines two approaches: i) removal of word senses deemed problematic for
316 (literature-centric) NLP usage, based on a manual study of UMLS ambiguity, and ii) a word
317 sense disambiguation algorithm that chooses a concept with the most likely semantic type for
318 a given context [33].

319 **NCBO annotator** [14] is provided by the U.S. National Center for Biomedical Ontology
320 (NCBO) as a freely available Web service. It is based on a two-stage annotation process.
321 The first stage relies on a concept recognition tool that uses a dictionary to identify mentions
322 of biomedical concepts in the input text. In particular, NCBO annotator makes use of the
323 MGrep tool [34], which was chosen over MetaMap due to its better performance along
324 several examined dimensions [35]. The dictionary for this annotation stage is built by pulling
325 concept names and descriptions from biomedical ontologies and/or thesauri relevant for the
326 domain of the corpus to be annotated (typically UMLS Metathesaurus and BioPortal
327 ontologies, Table 1). In the second stage, the initial set of concepts, referred to as direct
328 annotations, is extended using the structure and semantics of relevant biomedical ontologies.
329 For instance, semantic distance measures are used to extend the direct annotations with
330 semantically related concepts; the computation of semantic distance is configurable, and can
331 be based, for instance, on the distance between the concepts in the ontology graph.
332 Semantic relations between concepts from different ontologies, established through ontology
333 mappings, serve as another source for finding semantically related concepts that can be
334 used to extend the scope of direct annotations. The NCBO annotator is unique in its
335 approach to associate concept mentions with multiple concepts, instead of finding one
336 concept that would be the best match for the given context.

337 **BioMedical Concept Annotation System (BeCAS)** [36] is a Web-based tool for semantic
338 annotation of biomedical texts, primarily biomedical research papers. Besides being available
339 through a Web-based user interface, it can be programmatically accessed through a Web-
340 based (RESTful) Application Programming Interface (API), and a widget, easily embeddable in
341 Web pages. Like majority of the aforementioned annotation tools, BeCAS is an open-source
342 modular system, comprising of several modules for text preprocessing including, e.g.,
343 sentence splitting, tokenization, lemmatization, among others, as well as modules for
344 concept detection and abbreviation resolution. Most of the concept detection modules in
345 BeCAS apply a term-to-concept matching approach to identify and annotate mentions of

346 several types of biomedical entities, including species, enzymes, chemicals, drugs, diseases,
347 etc. This approach relies on a custom dictionary, i.e., a database of concepts and associated
348 terms, compiled by pulling concepts from various meta-thesauri and ontologies such as
349 UMLS Metathesaurus, NCBI BioSystems database, ChEBI, and the Gene Ontology (Table
350 1). For the identification of gene and protein mentions and their disambiguation with
351 appropriate concepts, BeCAS makes use of Gimli, an open source tool that implements
352 Conditional Random Fields (CRF) for named entity recognition in biomedical texts [37] (see
353 Section 4).

354 **Whatizit** is a freely available Web service for annotation of biomedical texts with concepts
355 from several ontologies and structured vocabularies [38]. Like previously described tools, it is
356 also developed in a modular way so that different components can be combined into custom
357 annotation pipelines, depending on the main theme of the text being processed. For
358 example, *whatizitGO* is a pipeline for identifying Gene Ontology (GO) concepts in the input
359 text, while *whatizitOrganism* identifies species defined in the NCBI taxonomy. In Whatizit,
360 concept names are transformed into regular expressions to account for morphological
361 variability in the input texts [39]. Such regular expressions are then compiled into Finite State
362 Automata, which assure quick processing regardless of the length of the text and the size of
363 the used vocabulary; therefore, processing time is linear with respect to the length of the text.
364 Whatizit also offers pipelines that allow for the recognition of biomedical entities of a specific
365 type based on two or more knowledge sources. For instance, *whatizitSwissprotGo* is the
366 pipeline for the annotation of protein mentions based on the UniProtKb/Swiss-Prot
367 knowledge base (Table 1) and the Gene Ontology. Finally, there are more complex pipelines
368 that combine simpler pipelines to enable detection and annotation of two or more types of
369 biomedical entities. For instance, *whatizitEbiMed* incorporates *whatizitSwissprotGo*,
370 *whatizitDrug* and *whatizitOrganism* to allow for the detection and annotation of proteins,
371 drugs and species.

372 **ConceptMapper** [21] is a general purpose dictionary lookup tool, developed as a component
373 of the open-source UIMA NLP framework. Unlike the other annotators that have been
374 examined so far, ConceptMapper is the only one that was not specifically developed for the
375 biomedical domain, but is rather generic and configurable-enough to be applicable to any
376 domain. Its flexibility primarily stems from the variety of options for configuring its algorithm
377 for mapping dictionary entries onto input text. For instance, it can be configured to detect
378 entity mentions even when they appear in the text as disjoint multi-word phrases, e.g., in the
379 text “intraductal and invasive mammary carcinoma”, it would recognize “intraductal
380 carcinoma” and “invasive carcinoma” as diagnosis. It can also deal with a variety of ways a
381 concept can be mentioned in the input text, e.g., synonyms and different word forms. This is
382 enabled by a dictionary that for each entry stores several possible variants, and connects
383 them to the same concept. For instance, the entry with the main (canonical) form “spine”
384 would also include variants such as “spinal”, “spinal column”, “vertebral column”, “backbone”,
385 and others, and associates them all with the semantic type AnatomicalSite. Even though
386 ConceptMapper is not originally targeted at the biomedical domain, if properly configured, it
387 can even outperform state-of-the-art biomedical annotators [24]. However, the task of
388 determining the optimal configuration and developing a custom dictionary might be
389 overwhelming for regular users; we return to this topic in Section 5.3.

390 **Neji** [40] is yet another open source and freely available software framework for annotation
391 of biomedical texts. Its high modularity is achieved by having each text processing task
392 wrapped in an independent module. These modules can be combined in different ways to
393 form different kinds of text processing and annotation pipelines, depending on the
394 requirements of specific annotation tasks. The distinct feature of Neji is its capacity for multi-
395 threaded data processing, which assures high speed of the annotation process. Neji makes
396 use of existing software tools and libraries for text processing, e.g., tokenization, sentence
397 splitting, lemmatization, with some adjustments to meet the lexical specificities of biomedical
398 texts. For concept recognition, Neji supports both dictionary-lookup matching and ML-based

399 approaches by customizing existing libraries that implement these approaches. For instance,
400 like BeCAS, it uses the CRF tagger implemented in Gimli. Hence, various CRF models
401 trained for Gimli can be used in Neji, each model targeting a specific type of biomedical
402 entities such as genes or proteins. Since Gimli does not perform disambiguation, Neji has
403 introduced a simple algorithm to associate each recognized entity mention with a unique
404 biomedical concept.

405 3.2 Summary of benchmarking results

406 Tseytlin et al. [20] have conducted a comprehensive empirical study that includes five state-
407 of-the-art semantic annotators that were compared based on the execution time and
408 standard annotation performance metrics (precision, recall, F1-measure). Four of the
409 benchmarked tools, namely cTAKES, MetaMap³, ConceptMapper, and NOBLE Coder have
410 been directly covered in the previous section, whereas the fifth tool - MGrep - was
411 considered as a service used by NCBO Annotator in the first stage of its annotation process.
412 The benchmarking was done on two publicly available, human-annotated corpora (see Table
413 4): one (ShARe) consisting of annotated clinical notes, the other (CRAFT) of annotated
414 biomedical literature. Documents from the former corpus (ShARe) were annotated using the
415 SNOMED-CT vocabulary (Table 1), while for the annotation of the latter corpus (CRAFT), a
416 subset of OBO ontologies were used as recommended by the corpus developers.

417 The study showed that all the tools performed better on the clinical notes corpus (ShARe)
418 than on the corpus of biomedical literature (CRAFT). The results demonstrated that on the
419 ShARe corpus, NOBLE Coder, cTAKES, MGrep, and MetaMap were of comparable
420 performance, while only ConceptMapper somewhat lagged behind. On the CRAFT corpus,
421 NOBLE Coder, cTAKES, MetaMap, and ConceptMapper were quite aligned, whereas MGrep
422 performed significantly worse, due to very low recall. In terms of speed, on both corpora,
423 ConceptMapper proved to be the fastest one. It was followed by cTAKES, NOBLE Coder,

424 and MGrep, respectively, whose speed was more-or-less comparable. However, MetaMap
425 was by far the slowest (about 30 times slower than the best performing tool).

426 Another comprehensive empirical study that compared several semantic annotators with
427 respect to their speed and the quality of the produced annotations is reported in [40]. The
428 study included five contemporary annotators - Whatizit, MetaMap, Neji, Cocoa, and
429 BANNER, which were compared on three manually annotated corpora of biomedical
430 publications, namely NCBI Disease corpus, CRAFT, and AnEM (see Table 4). Evaluation on
431 the CRAFT corpus considered 6 different biomedical entity types (e.g. species, cell, cellular
432 component, gene and proteins), while on the other two corpora only the most generic type
433 was considered, i.e., anatomical entity for AnEM, and disorder for NCBI. Two of the
434 benchmarked annotators are either no longer available (Cocoa) or no longer maintained
435 (BANNER⁴), whereas the other three were covered in the previous section. Benchmarking
436 was done for each considered type of biomedical concept separately, and also using
437 different configurations of the examined tools (e.g., five different term-to-concept matching
438 techniques were examined).

439 The study showed that the tools' performance varied considerably between various
440 configuration options, in particular, various strategies for recognizing entity mentions in the
441 input text. This variability in the performance associated with different configurations was
442 also confirmed by Funk et al [24]; we return to this topic in Section 5.4.

443 Overall, Neji had the best results, especially on the CRAFT corpus, with significant
444 improvements over the other tools on most of the examined concept types. Whatizit proved
445 to have the most consistent performance across different configuration options, with an
446 average variation of 4% in F1-measure. In terms of speed, Neji significantly outpaced the
447 other tools.

448

449 **Table 4.** Corpora used for the evaluation of biomedical semantic annotators. The table
450 includes corpora that were used in the reported use cases (Section 2, Table 2), and/or used
451 for benchmarking the discussed tools (Section 3.2 and Section 4)

452 [Table 4 should be placed approximately here]

453 4. Entity-specific biomedical annotation tools

454 While the primary focus of this paper is on biomedical semantic annotation, and in particular
455 general purpose biomedical semantic annotators, the work in the closely related area of
456 biomedical Named Entity Recognition (NER) also deserves to be mentioned given that it has
457 been a precursor to the biomedical semantic annotation technology. Early work in biomedical
458 NER were mainly focused on developing dictionary-based, rule-based, or heuristics-based
459 techniques for identifying entity mentions within chemical, biological, and medical corpora.
460 Some of the earlier works include the work by Fukuda et al. [41] that used rules for extracting
461 protein names, MedLEE [42] that employed contextual rules to perform mapping to an
462 encoding table extracted from UMLS, and EDGAR [43] that extracted drugs and genes
463 related to cancer. However, more advanced techniques based on machine learning (ML)
464 models, more specifically Hidden Markov Models (HMM), Conditional Random Fields (CRF),
465 and Support Vector Machines (SVM), have become more prominent in the recent years.

466 ABNER [44] was one of the earlier works that benefited from CRF models and was trained
467 for five specific entity types, namely Protein, DNA, RNA, Cell Line, and Cell Type. ABNER
468 extracted features based on regular expressions and n-grams to train a CRF model, and did
469 not introduce any syntactic or semantic features in this process. Gimli [37] is a more recent
470 NER toolkit that is also based on CRF models. The main advantage of Gimli is its
471 introduction of a wide range of features, namely: orthographic, linguistic, morphological,
472 external, and local context features. The orthographic features include capitalized mentions,
473 counting, and symbol type features, while the linguistic features consist of word lemmas,

474 POS tags, and products of dependency parsing. The morphological features cover n-grams
475 and word shapes. The local and external features constitute gene and protein names as well
476 as trigger words. The wide spectrum of features enables the CRF model to be highly
477 accurate on different benchmark datasets including GENETAG and JNLPBA (see Table 4).

478 The work by Leaman et al [45], known as DNorm, is a method specifically built for disease
479 mention detection in biomedical text. DNorm is based on BANNER [46] for disease mention
480 detection and subsequently uses a pairwise learning to rank framework to perform
481 normalization. Similar in objective to DNorm but with focus on genes, SR4GN [47] is a rule-
482 based system specifically built to link species with corresponding gene mentions. This tool
483 has shown better performance compared to LINNAEUS [48], which is a tool for the same
484 purpose built using a dictionary-based approach for mention detection and a set of heuristics
485 for ambiguity resolution. The authors of SR4GN subsequently proposed GNormPlus that
486 focuses on the identification of gene names and their identifiers. The distinguishing aspect of
487 this tool is that it is able to distinguish gene, gene family, and protein domains by training a
488 supervised CRF model on annotated gene corpora.

489 There have also been attempts at combining the benefits of rule-based methods and ML
490 techniques. For instance, OrganismTagger [49] uses a set of grammar rules written in the
491 JAPE language, a set of heuristics, as well as an SVM classifier to identify and normalize
492 organism mentions in text including genus, species, and strains.

493 In a later publication [50], the developers of DNorm discussed the benefits of developing an
494 entity type agnostic NER framework that could be retrained easily given sufficiently
495 annotated training data and a related lexicon. Based on this objective, the TaggerOne tool
496 was developed as an entity type independent tool that employs a semi-Markov structured
497 linear classifier and has shown favorable performance on both NCBI Disease corpus as well
498 as the chemical BioCreative 5 CDR corpus (see Table 4). In contrast to tools such as
499 TaggerOne that rely only on a single ML model, there has also been work in the literature

500 that rely on ensembles of models. For instance, tmChem, an ensemble built on BANNER
501 and tmVar [51], focuses on the recognition of seven different types of chemical mentions in
502 biomedical literature, namely Abbreviation, Family, Formula, Identifier, Multiple, Systematic
503 and Trivial.

504 While the above approaches benefit from annotated corpora and some form of (semi)
505 supervised training, such methods are task and entity dependent, and training them on new
506 entity types is time consuming and resource intensive [52]. For this reason, unsupervised
507 NER methods have started to emerge. For instance, the method proposed by Zhang and
508 Elhadad [52] uses a noun chunker to detect possible entity candidates and subsequently
509 categorizes the entity candidates based on distributional semantics. This method showed
510 reasonable performance on the i2b2 and GENIA corpora (see Table 4).

511 It is worth mentioning that given the large amount of biomedical documents and texts that
512 need to be processed by NER tools, several researchers have looked at optimizing the
513 parallel capabilities of these tools. The work by Tang et al. [53] and Li et al. [54] are two
514 notable recent work in this respect. These two works contend that given the sequential
515 nature of CRF models, their parallelization is not trivial. On this basis, they show how the
516 MapReduce framework can be used to efficiently train CRFs for biomedical NER.

517 It is also important to note that research and development of biomedical named entity
518 recognition and normalization tools have been fostered through different initiatives of the
519 biomedical research community. A notable one is the BioCreative initiative
520 (<http://www.biocreative.org/tasks/>), a series of yearly challenges focused on text mining and
521 information retrieval tasks relevant to the life science domain, including recognition of
522 chemicals, genes, drugs, and diseases in biomedical texts. For instance, one of the tasks at
523 the BioCreative IV challenge [55] was to automatically identify terms in a given article that
524 refer to the concepts from the Gene Ontology (GO; see Table 1), that is, to semantically
525 annotate articles with GO concepts. Benchmarking of the proposed solutions was done on

526 the BC4GO corpus (see Table 4). The best performing team on this task applied a
527 supervised classification method that relies on a knowledge base built by leveraging a large
528 database of (over 100K) MEDLINE abstracts annotated with GO terms [56]. In particular, the
529 tool developed by this team, known as the GOCat tool, relies on similarities between an input
530 text and already curated instances in the tool's knowledge base, to annotate the input with
531 the most prevalent GO terms among the instances from the knowledge base. The
532 BioCreative V challenge hosted a Disease Named Entity Recognition (DNER) task [57],
533 where the participating systems were given PubMed titles and abstracts and asked to return
534 normalized disease concept identifiers. The benchmarking of the submitted solutions was
535 done on the Chemical-Disease Relation (CRD) corpus (see Table 4). The best system
536 (based on a CRF model with post-processing) achieved an F-score of 86.46%, a result that
537 approaches the human inter-annotator agreement (0.8875). A large majority of the proposed
538 solutions relied on ML (only 3 out of 16 were based exclusively on a dictionary-lookup
539 method); one third of these solutions (4) used ML only, while others (8) exploited a
540 combination of ML with dictionaries and/or pattern matching.

541 For a more comprehensive list of biomedical NER tools, in-depth discussion on the
542 techniques, features and corpora used, the entity types that are covered and a comparative
543 performance analysis, we refer the interested reader to the work by Campos et al. [58].

544 5. Challenges

545 Even though significant efforts have been devoted to the development of sophisticated
546 semantic annotation tools, there are still challenges that need to be resolved if these tools
547 are to reach their full potential. This section points to some of those challenges, as well as to
548 some of the existing research work that offers potential solutions.

549 5.1 The lack of sufficient context for understanding entity mentions

550 An often cited source of difficulty associated with the recognition of entities in biomedical
551 texts is the lack of sufficient context for interpreting the entity mentions [59]. For instance,
552 Tseytlin et al. [20] reported that the largest proportion of annotation errors made by their
553 NOBLE Coder annotator was due to the missing or incomplete context or background
554 knowledge.

555 The collective annotation approach was proposed as a way of dealing with this challenge
556 [59]. It relies on the global topical coherence of entities mentioned in a piece of text and is
557 done by disambiguating a set of related mentions simultaneously. The basic idea is that if
558 multiple entity mentions co-occur in the same sentence or paragraph, they can be
559 considered semantically related. In particular, the approach proposed by Zheng et al. [59]
560 consists of creating a document graph (G_d) with entity mentions recognized in a document as
561 nodes, while edges are established between those pairs of nodes (entity mentions) that co-
562 occur in the same sentence or paragraph of the document. Each entity mention is then
563 connected with one or more entity candidates from the knowledge base (KB) based on the
564 name variants associated with entities in the KB. Finally, for each entity mention (m) - entity
565 candidate (c) pair (m,c), a score is computed based on i) the general popularity of the
566 candidate entity c in the KB, that is, its level of connectedness to other entities in the KB
567 (non-collective score), and ii) level of connectedness of candidate c only with candidate
568 concepts of entity mentions that are connected to mention m in the G_d graph. The candidate
569 entity c from the (m,c) pair with the highest score is selected as the appropriate entity for the
570 given entity mention m . A similar approach was proposed and proved effective for general
571 purpose semantic annotators in work such as [60].

572 5.2 Scaling to very large document sets

573 One of the weaknesses of today's biomedical semantic annotators lies in their speed, that is,
574 the time required for completing the annotation task on very large corpora (with tens and
575 hundreds of millions of documents) [20]. Note that speed estimates given in Table 3
576 (qualifying almost all examined tools as suitable for real-time processing) are based on the
577 experimental results reported in the literature, where experiments were done with small
578 corpora (up to 200 documents).

579 Divita et al. [61] aimed at using semantic annotation to improve information extraction and
580 retrieval of clinical notes from the Veterans Informatics and Computing Infrastructure (VINCI)
581 hosting huge and continuously growing amounts of medical notes. However, they found
582 today's annotators unapt for that task, as, based on the Divita et al., even when running on
583 several multi-core machines, today's annotators would need multiple years to index VINCI
584 notes with semantic concepts. As a solution to this challenge, they proposed Sophia, an
585 UMLS-based annotation tool, that deals with high throughput by replicating either certain
586 components of the annotation pipeline or the entire pipeline [61]. Sophia is built from the
587 components of the v3NLP framework [62], a suite of middleware text-processing components
588 aimed for building various kinds of NLP applications. In particular, Sophia makes use of the
589 v3NLP components for dealing with the idiosyncrasies of clinical texts, as well as the
590 framework's scaling-up and scaling-out functionalities for efficiently handling huge quantities
591 of texts.

592 5.3 Adaptation to new document type(s) and/or terminologies specific to 593 particular biomedical subdomain

594 Another challenge originates in the variety of biomedical texts and differences among
595 different kinds of text, particularly differences between biomedical literature and clinical text
596 [20, 25]. According to Garla and Brandt [7], "clinical text is often composed of semi-

597 grammatical ‘telegraphic’ phrases, uses a narrower vocabulary than biomedical literature,
598 and is rife with domain-specific acronyms.” In addition, common to both clinical texts and
599 scientific papers is the presence of local dialects, such as specific jargon developed within a
600 medical center, or particular, idiosyncratic protein nomenclatures created within research
601 laboratories [63]. Due to these issues, an annotation tool developed and/or configured for a
602 particular type of medical texts or even one application case, tied to a particular medical
603 institution/center, cannot be directly ported to some other text type and/or application case
604 without, often significant, drop in performance.

605 A potential solution to this diversity in text types and terminologies is the use of flexible
606 general-purpose annotation tools that can be configured to work with different text types and
607 vocabularies [19]. In fact, Funk et al. [24] have demonstrated that if properly tuned, a generic
608 annotation tool can offer better performance than tools designed specifically for particular
609 biomedical task or domain. The key is in the modularity and flexibility of a tool so that one
610 can precisely control how terms in the text are to be matched against the available
611 terminologies.

612 While the majority of the annotators listed in Table 3 were developed to be modular and
613 flexible, their configuration is a complex task for users lacking expertise in NLP and not
614 knowing the intricacies of the tool’s internal functioning. The latter is especially relevant as
615 not all parameters equally affect the performance; also, the interaction of the parameters
616 need to be considered.

617 Besides configuring the tool’s annotation method, e.g., kinds of text processing and term
618 matching options, adaptation to a different biomedical (sub)domain also requires either
619 development or, at least, customization of the dictionary that the tool uses to recognize
620 concept mentions in the input text. While there are numerous ontologies, knowledge bases,
621 thesauri, and similar kinds of biomedical resources that can be used for dictionary
622 development, that task is often overly complex for regular users. This is because each tool

623 has its own idiosyncratic structure and format for vocabulary representation and storage,
624 designed to optimally match the tool's annotation algorithm. To alleviate the task of dictionary
625 development / customization, Tseytlin et al. [20] have developed an interactive terminology
626 building tool, as a component of the NOBLE Coder annotator. The tool allows users to import
627 existing terminologies (of various kinds), and then customize them by selecting only certain
628 segments (branches) of the imported terminologies, and/or to filter them by semantic types.
629 A tool of this type would be a useful complement to any semantic annotator that relies on a
630 dictionary-lookup approach.

631 5.4 Application-specific tool benchmarking

632 As argued in Section 3.1, benchmarking of semantic annotators requires that each annotator
633 is configured based on the specificities of the benchmarking task, so that it demonstrates its
634 optimal performance on the task. The effect of configuration on the annotators' performance
635 was well demonstrated by Funk et al [24] in their comprehensive empirical study that
636 included MetaMap, ConceptMapper, and NCBO Annotator (see Section 3.1). The
637 researchers examined over 1,000 parameter combinations in the context of the CRAFT
638 evaluation corpus (Table 4) and 8 different terminologies (ontologies). They found that
639 default parameter values often do not lead to the best performance, and that by appropriately
640 setting parameters, F-measure can be significantly increased (even by 0.4 points). This
641 suggests that if it is to be used for making a decision on the annotator to adopt in a particular
642 application case, the benchmarking studies should not be based on the tools' default
643 configuration, but should include tools customized to the specific features of the application
644 case.

645 Another requirement for application-specific benchmark study is the selection of appropriate
646 evaluation corpora and terminology source for building or customizing the tools' dictionaries.
647 While numerous annotated corpora have been developed (Table 4), including both manually
648 annotated gold standard corpora and corpora annotated in an automated or semi-automated

649 way known as silver standards, the information about these resources are dispersed on the
650 web and it takes time and effort to collect information about the available evaluation corpora
651 and their features.

652 Considering the above stated difficulties associated with the setup of application-specific
653 benchmarking studies, we point to the need for a benchmarking ‘toolkit’ that would facilitate
654 the task of tool benchmarking in the context of a specific application case. An important
655 component of such a toolkit would be a searchable registry of existing annotated corpora.
656 For each corpus, the registry should include basic qualitative and quantitative information,
657 e.g., the sources and types of documents that it includes, and the vocabularies or ontologies
658 that were used for annotating the corpus, among others. In addition, the registry of annotated
659 corpora would need to contain guidelines for how each corpus should be used, references to
660 the studies where the corpus was previously used, and any additional information that might
661 be of relevance for effective use of the given corpus.

662 Another important component of the benchmarking toolkit would be guidelines and/or tools
663 for optimal configuration of annotation tools. The starting point for such guidelines could be
664 the set of suggestions that Funk et al [24] derived from their study, related to the selection of
665 optimal parameter values based on the terminology (ontology) to be used for annotation.
666 Tools enabling semi-automated or automated parameter tuning would greatly facilitate this
667 task. Algorithmic procedures and tools developed for general purpose semantic annotators,
668 like the one proposed in [64], could be adapted to tune parameters of biomedical annotators.

669 With such a benchmarking toolkit, it would be also possible to evaluate the performance of
670 general purpose biomedical annotators on the tasks of recognizing and normalizing specific
671 types of biomedical entities, e.g., chemicals, genes, drugs, or diseases. This would allow for
672 evidence-based recommendation of appropriate semantic annotators for entity-specific tasks.
673 While some initial work in this direction has been done by Campos et al. [40] (see Section
674 3.2), only a small number of the current tools have been examined (some of the tools

675 evaluated in their study are no longer available), and they were not tuned to the entity
676 specific annotation tasks. Henceforth, new studies with current general purpose annotators,
677 customized for the entity-specific task at hand, are needed in order to obtain conclusive
678 evidence on the performance of the current tools for specific biomedical entity types.

679 5.5 Semantic annotation in languages other than English

680 Large majority of tools, ontologies, and corpora developed for biomedical semantic
681 annotation, and biomedical NLP in general, are for the English language. Semantic
682 annotators discussed in the previous sections fall in this category of "English-only" tools.
683 However, the development of NLP resources and tools for semantic annotation in languages
684 other than English has started receiving increasing attention both in research and practice.

685 The CLEF (Conference and Labs of the Evaluation Forum) conference series have been
686 hosting eHealth Labs where one of the tasks has been entity recognition and normalization,
687 i.e., semantic annotation, in languages other than English, primarily French. Systems
688 developed to face this challenge varied greatly [65, 66]. The team with the best performance
689 at the latest eHealth Lab, held in conjunction with CLEF 2016, proposed a system that could
690 be qualified as a general purpose semantic annotator [67]. In particular, to perform the entity
691 recognition task, this system used Peregrine [68], a dictionary-based concept recognition
692 tool, in conjunction with a dictionary consisting of French vocabularies from UMLS
693 supplemented with automatically translated English UMLS terms. Several post-processing
694 steps were implemented to reduce the number of false positives, such as filtering based on
695 precision scores derived from the training data. Entity normalization relied on the
696 <entity_mention, semantic_group, CUI⁶> combinations extracted from the training set.

697 Another important initiative was the CLEF-ER challenge that took place in 2013 as part of the
698 Mantra project aimed at providing multilingual documents and terminologies for the
699 biomedical domain [69]. For this challenge, Medline and biomedical patent documents were

700 released in five languages: English, German, French, Spanish, and Dutch. Mappings to
701 English documents were provided for all documents that were in a language other than
702 English, though the mappings were not available between all pairs of languages, e.g.,
703 between Spanish and German. The organizers also released the CLEF-ER terminology, a
704 multilingual vocabulary with term synonyms in the above mentioned five languages⁵. The
705 challenge received several submissions dealing with various challenges of multilingual
706 biomedical NLP, including semantic annotation, e.g. [70, 71], and the creation of multilingual
707 corpora, e.g., [72, 73]. An interesting approach to multilingual semantic annotation was
708 proposed by Attardi et al. [71]. The method starts from the English language Silver Standard
709 Corpus (SSC) provided by the CLEF-ER organizers [72], which is first translated into a target
710 language corpus, and then entity annotations are 'transferred' to it. The translation is done
711 using an open-source toolkit for statistical phrase-based machine translation. The word
712 alignment information produced by the translation tool is used to determine the
713 correspondence between entities in the source and the target language sentences. The
714 resulting annotated corpus is referred to as the Bronze Standard Corpus (BSC). In addition,
715 a dictionary of entities is also created, which associate each <entity_mention,
716 semantic_group> pair with all the corresponding CUIs that appeared in the SSC. The BSC is
717 used to train a Named Entity detection model, which is aimed at associating entity mentions
718 with their semantic groups. This model is then used for tagging entity mentions in the target
719 language sentences with the proper semantic group. Finally, after entity mentions have been
720 assigned to their semantic group, each mention is linked to corresponding CUIs by looking
721 up CUIs associated with the <entity_mention, semantic_group> pairs in the previously built
722 dictionary.

723 6. Conclusions

724 In this paper we have analyzed the current state of the art in the domain of general purpose
725 biomedical semantic annotators, and pointed to some of the areas where further research

726 and development is needed to improve the performance of the current solutions and make
727 them robust to the requirements of real-world biomedical applications. In conclusion, we can
728 say that the majority of the analyzed tools proved to be highly modular and configurable, thus
729 fulfilling the promise of general purpose biomedical annotators as annotators adaptable to
730 different areas of biomedicine. In addition, the majority of the examined tools are made
731 publicly available as open-source software libraries, thus bootstrapping further developments
732 in biomedical semantic annotation. As areas that require further research and development,
733 we have identified: i) finding new, more effective ways of dealing with the often terse context
734 of biomedical entity mentions, especially in clinical texts; ii) improving the scalability of
735 annotators so that they can efficiently process biomedical corpora with tens and hundreds of
736 millions of documents; iii) development of auxiliary tools that would facilitate the task of
737 customizing an annotator to the requirements of a particular annotation task, iv) development
738 of a toolkit for benchmarking semantic annotators in the context of a specific application
739 case, and thus enabling users to make well-informed decisions regarding the annotator to
740 use in their particular application setting, and vi) continuing and intensifying research efforts
741 aimed at multilingual biomedical semantic annotation.

742 We have also pointed to some of the potential benefits and application cases of biomedical
743 semantic annotation technology in order to demonstrate and exemplify the opportunities that
744 this technology can bring about, and thus encourage the research community to put efforts in
745 overcoming the identified challenges and bring the tools to their full potential. We believe that
746 there is a tremendous potential in using biomedical semantic annotation technology for
747 processing, analyzing and structuring unstructured textual biomedical content both in the
748 form of clinical and research material, and hope that this review paper provides the means to
749 encourage the community to further investigate and adopt this technology.

750 Endnotes

751 ¹ At the time of writing this manuscript, the given list of tools could be considered exhaustive with
752 respect to the given selection criteria, i.e., we included all the tools that met the given set of
753 criteria and were reported in the literature. However, considering the pace of new developments,
754 it is reasonable to expect new tools with the given characteristics soon to emerge.

755 ² Speed is characterized only from the perspective of the tool's usability for real-time text
756 annotation; we do not report exact operation time since it can vary considerably depending on the
757 tool's configuration, the characteristics of the corpora, the machine the tool is running on.

758 ³ The study used MMTX (<https://mmtx.nlm.nih.gov/>), Java implementation of MetaMap, which
759 produces only slightly different results than MetaMap [20] due to differences in tokenization and
760 lexicalization procedures.

761 ⁴ Source code is available from <http://banner.sourceforge.net/> but the last update was in year
762 2011.

763 ⁵ This terminology can be accessed from the project output page of the Mantra project website:
764 <https://sites.google.com/site/mantraeu/project-output>

765 ⁶ CUI stands for Concept Unique Identifier, that is, a unique identifier of a concept in a knowledge
766 base, such as UMLS Metathesaurus.

767 List of abbreviations

768 CRF - Conditional Random Fields

769 CUI - Concept Unique Identifier

770 KB - Knowledge Base

771 NLP - Natural Language Processing

772 NER - Named Entity Recognition

773 ML - Machine Learning

774 Declarations

775 Ethics approval and consent to participate

776 Not applicable

777 Consent for publication

778 Not applicable

779 Availability of data and material

780 Data sharing is not applicable to this article as no datasets were generated or analysed
781 during the current study.

782 Competing interests

783 The authors declare that they have no competing interests.

784 Funding

785 The second author graciously acknowledges funding from The Natural Sciences and
786 Engineering Research Council of Canada (NSERC).

787 Authors' contributions

788 The manuscript was written jointly by the two authors. Both authors read and approved the
789 final manuscript.

790 Acknowledgements

791 Not applicable

792 References

- 793 1. Fleuren WWM, Alkema W. Application of text mining in the biomedical domain. *Methods*. 2015;
794 74: 97–106.
- 795 2. Sinsky C, Colligan L, Li L, Prgomet M, Reynolds S, Goeders L, et al. Allocation of Physician Time
796 in Ambulatory Practice: A Time and Motion Study in 4 Specialties. *Ann Intern*
797 *Med*. 2016;165(11):753-760.
- 798 3. Hill RG, Sears LM, Melanson SW. 4000 Clicks: a productivity analysis of electronic medical
799 records in a community hospital ED. *The American Journal of Emergency Medicine*. 2013;
800 31(11):1591-1594.
- 801 4. Savova GK, Masanz JJ, Ogren PV, Zheng J, Sohn S, Kipper-Schuler KC, Chute CG. Mayo
802 clinical Text Analysis and Knowledge Extraction System (cTAKES): architecture, component
803 evaluation and applications. *J Am Med Inform Assoc*. 2010; 17(5):507–513.
- 804 5. Demner-Fushman D, Seckman C, Fisher C, Hauser SE, Clayton J, Thoma GR. A Prototype
805 System to Support Evidence-based Practice. In: *Proceedings of the 2008 Annual Symposium of*
806 *the American Medical Information Association (AMIA 2008)*. Washington, DC: 2008. p.151-155.
- 807 6. Sfakianaki P, Koumakis L, Sfakianakis S, Iatraki G, Zacharioudakis G, Graf N, *et al*. Semantic
808 biomedical resource discovery: a Natural Language Processing framework. *BMC Medical*
809 *Informatics and Decision Making*. 2015; 15:77.
- 810 7. Garla VN, Brandt C. Knowledge-based biomedical word sense disambiguation: an evaluation and
811 application to clinical document classification. *Journal of the American Medical Informatics*
812 *Association*. 2013; 20(5):882–886.
- 813 8. Garla V, Re VL, Dorey-Stein Z, et al. The Yale cTAKES extensions for document
814 classification: architecture and application. *Journal of the American Medical Informatics*
815 *Association: JAMIA*. 2011;18(5):614-620. doi:10.1136/amiajnl-2011-000093.
- 816 9. Mork JG, Yepes AJJ, Aronson AR. The NLM medical text indexer system for indexing biomedical
817 literature. In: *Proceedings of the first Workshop on Bio-Medical Semantic Indexing and Question*
818 *Answering*. Valencia, Spain: 2013.
- 819 10. Yamamoto Y, Yamaguchi A, Yonezawa A. Building Linked Open Data towards integration of
820 biomedical scientific literature with DBpedia. *Journal of Biomedical Semantics*. 2013; 4:8.
- 821 11. Chang J, Schutze H, Altman R. Creating an Online Dictionary of Abbreviations from MEDLINE.
822 *The Journal of the American Medical Informatics Association*. 2002; 9(6):612–620.

- 823 12. Rigden DJ, Fernández-Suárez XM, Galperin MY. The 2016 database issue of Nucleic Acids
824 Research and an updated molecular biology database collection. *Nucl. Acids Res. (Database*
825 *Issue)*. 2016; 44(D1):D1-D6 doi:10.1093/nar/gkv1356
- 826 13. Legaz-García MC, Miñarro-Giménez JA, Menárguez-Tortosa M, Fernández-Breis JT. Generation
827 of open biomedical datasets through ontology-driven transformation and integration processes.
828 *Journal of Biomedical Semantics*. 2016; 7:32.
- 829 14. Jonquet C, Shah N, Musen M. The Open Biomedical Annotator. *AMIA Summit on Translational*
830 *Bioinformatics*. San Francisco, CA, United States: 2009. p.56-60.
- 831 15. Weeber M, Mork J, Aronson A. Developing a test collection for biomedical word sense
832 disambiguation. In: *Proceedings of AMIA symposium*. Washington, DC: 2001. p.746–750.
- 833 16. Meystre SM, Savova GK, Kipper-Schuler KC, Hurdle JF. Extracting Information from Textual
834 Documents in the Electronic Health Record: A Review of Recent Research. *IMIA Yearbook*,
835 2008:128–144.
- 836 17. Hatzivassiloglou V, Duboué PA, Rzhetsky A. Disambiguating proteins, genes, and RNA in text: a
837 machine learning approach. *Bioinformatics*. 2001; 17:S97-S106.
- 838 18. Chen L, Liu H, Friedman C. Gene name ambiguity of eukaryotic nomenclatures. *Bioinformatics*.
839 2004; 21(2):248-256.
- 840 19. Jovanovic J, Bagheri E, Cuzzola J, Gasevic D, Jeremic Z, Bashash R. Automated Semantic
841 Annotation of Textual Content. *IEEE IT Professional*. 2014; 16(6):38-46.
- 842 20. Tseytlin E, Mitchell K, Legowski E, Corrigan J, Chavan G, Jacobson RS. NOBLE – Flexible
843 concept recognition for large-scale biomedical natural language processing. *BMC Bioinformatics*.
844 2016; 17:32.
- 845 21. Tanenblatt M, Coden A, Sominsky IL. The ConceptMapper Approach to Named Entity
846 Recognition. In: *Proc of 7th Language Resources and Evaluation Conference (LREC)*. 2010. p.
847 546–51.
- 848 22. Tikk D, Solt I. Improving textual medication extraction using combined conditional random fields
849 and rule-based systems. *J Am Med Inform Assoc*. 2010; 17(5):540–544.
- 850 23. Hsu CN, Chang YM, Kuo C-J, Lin YS, Huang HS, Chung IF. Integrating high dimensional bi-
851 directional parsing models for gene mention tagging. *Bioinformatics*. 2008; 24(13):i286–i294.
- 852 24. Funk C, Baumgartner W, Garcia B, Roeder C, Bada M, Cohen KB, Hunter LE, Verspoor K. Large-
853 scale biomedical concept recognition: an evaluation of current automatic annotators and their
854 parameters. *BMC Bioinformatics*. 2014; 15:59.
- 855 25. Chasin R, Rumshisky A, Uzuner O, Szolovits P. Word sense disambiguation in the clinical
856 domain: a comparison of knowledge-rich and knowledge-poor unsupervised methods. *Journal of*
857 *the American Medical Informatics Association*. 2014; 21(5):842–849.

- 858 26. Ling X, Weld DS. Fine-grained entity recognition. In *Proceedings of the Twenty-Sixth AAAI*
859 *Conference on Artificial Intelligence (AAAI'12)*. AAAI Press:2012. p.94-100.
- 860 27. Yaghoobzadeh Y, Schütze H. Corpus-level Fine-grained Entity Typing Using Contextual
861 Information. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language*
862 *Processing, EMNLP 2015, Lisbon, Portugal: 2015*. p.715-725.
- 863 28. Unstructured Information Management Architecture - UIMA. <https://uima.apache.org/> Accessed 7
864 December 2016.
- 865 29. OpenNLP. <https://opennlp.apache.org/> Accessed 30 November 2016.
- 866 30. Lesk M. Automatic sense disambiguation using machine readable dictionaries: how to tell a pine
867 cone from an ice cream cone. *Proceedings of the 5th Annual International Conference on*
868 *Systems Documentation; New York, NY, USA: 1986*. p.24-26.
- 869 31. Aronson AR. Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap
870 program. *Proceedings of the AMIA Symposium*. 2001:17-21.
- 871 32. Aronson AR, Lang FM. An overview of MetaMap: historical perspective and recent advances. *J*
872 *Am Med Inform Assoc*. 2010; 17(3):229–236.
- 873 33. Humphrey SM, Rogers WJ, Kilicoglu H, Demner-Fushman D, Rindfleisch TC. Word Sense
874 Disambiguation by Selecting the Best Semantic Type Based on Journal Descriptor Indexing:
875 Preliminary Experiment. *Journal of the American Society for Information Science and Technology*.
876 2006;57(1):96-113. doi:10.1002/asi.20257.
- 877 34. Dai M, Shah NH, Xuan W, Musen MA, Watson SJ, Athey B, Meng F. An Efficient Solution for
878 Mapping Free Text to Ontology Terms. *AMIA Summit on Translational Bioinformatics, San*
879 *Francisco, CA, USA: 2008*.
- 880 35. Shah NH, Bhatia N, Jonquet C, Rubin D, Chiang AP, Musen MA. Comparison of concept
881 recognizers for building the Open Biomedical Annotator. *BMC Bioinformatics*. 2009;10(Suppl
882 9):S14. doi:10.1186/1471-2105-10-S9-S14.
- 883 36. Nunes T, Campos D, Matos S, Oliveira JL. BeCAS: biomedical concept recognition services and
884 visualization. *Bioinformatics*, 2013; 29(15):1915–1916.
- 885 37. Campos D, Matos S, Oliveira, JL. Gimli: open source and high-performance biomedical name
886 recognition. *BMC Bioinformatics* 2013;14:54. <https://doi.org/10.1186/1471-2105-14-54>.
- 887 38. Rebholz-Schuhmann D, Arregui M, Gaudan S, Kirsch H, Jimeno A. Text processing through Web
888 services: calling Whatizit. *Bioinformatics*. 2008; 24(2):296–298.
- 889 39. Kirsch H, Gaudan S, Rebholz-Schuhmann D. Distributed modules for text annotation and IE
890 applied to the biomedical domain. *Int. J. Med. Inform*. 2006; 75:496–500.
- 891 40. Campos D, Matos S, Oliveira JL. A modular framework for biomedical concept recognition. *BMC*
892 *Bioinformatics*, 2013;14:281. <https://doi.org/10.1186/1471-2105-14-281>

- 893 41. Fukuda K, Tamura A, Tsunoda T, Takagi T. Toward information extraction: identifying protein
894 names from biological papers. *Pacific Symposium on Biocomputing*:1998. p.707–718.
- 895 42. Friedman C, Shagina L, Lussier Y, Hripcsak G. Automated Encoding of Clinical Documents
896 Based on Natural Language Processing. *Journal of the American Medical Informatics Association*:
897 *JAMIA*, 2004; 11(5):392–402. <http://doi.org/10.1197/jamia.M1552>
- 898 43. Rindflesch TC, Tanabe L, Weinstein JN, Hunter L. EDGAR: Extraction of Drugs, Genes and
899 Relations from the Biomedical Literature. *Pacific Symposium on Biocomputing*: 2000. p.517–528.
- 900 44. Settles B. ABNER: an open source tool for automatically tagging genes, proteins and other entity
901 names in text. *Bioinformatics*. 2005; 21(14):3191-3192.
902 <http://dx.doi.org/10.1093/bioinformatics/bti475>
- 903 45. Leaman R, Islamaj Doğan R, Lu Z. DNorm: disease name normalization with pairwise learning to
904 rank. *Bioinformatics*. 2013;29(22):2909-2917. doi:10.1093/bioinformatics/btt474.
- 905 46. Leaman R, Gonzalez G. BANNER: an executable survey of advances in biomedical named entity
906 recognition. *Pac Symp Biocomput*. 2008. p.652-663.
- 907 47. Wei C-H, Kao H-Y, Lu Z. SR4GN: A Species Recognition Software Tool for Gene Normalization.
908 *PLoS ONE*. 2012; 7(6): e38460. <https://doi.org/10.1371/journal.pone.0038460>
- 909 48. Gerner M, Nenadic G, Bergman CM. LINNAEUS: A species name identification system for
910 biomedical literature. *BMC Bioinformatics*, 2010; 11:85. <http://doi.org/10.1186/1471-2105-11-85>
- 911 49. Naderi N, Kappler T, Baker CJO, Witte R. OrganismTagger: detection, normalization and
912 grounding of organism entities in biomedical documents. *Bioinformatics* 2011; 27(19): 2721-2729.
913 doi: 10.1093/bioinformatics/btr452
- 914 50. Leaman R, Zhiyong L. TaggerOne: joint named entity recognition and normalization with semi-
915 Markov Models. *Bioinformatics* 2016; 32(18): 2839-2846. doi: 10.1093/bioinformatics/btw343
- 916 51. Wei C-H, Harris BR, Kao H-Y, Lu Z: tmVar: A text mining approach for extracting sequence
917 variants in biomedical literature. *Bioinformatics*. 2013, 29:1433-1439.
918 10.1093/bioinformatics/btt156.
- 919 52. Zhang S, Elhadad, N. Unsupervised biomedical named entity recognition. *J. of Biomedical*
920 *Informatics* 2013; 46(6): 1088-1098. DOI=<http://dx.doi.org/10.1016/j.jbi.2013.08.004>
- 921 53. Tang Z, Jiang L, Yang L, Li K, Li K. CRFs based parallel biomedical named entity recognition
922 algorithm employing MapReduce framework. *Cluster Computing* 2015; 18(2):493-505.
923 <http://dx.doi.org/10.1007/s10586-015-0426-z>
- 924 54. Li K, Ai W, Tang Z, Zhang F, Jiang L, Li K, Hwang K. Hadoop Recognition of Biomedical Named
925 Entity Using Conditional Random Fields. *IEEE Trans. Parallel Distrib. Syst.* 2015; 26(11):3040-
926 3051. <http://dx.doi.org/10.1109/TPDS.2014.2368568>

- 927 55. Mao Y, Van Auken K, Li D, et al. Overview of the gene ontology task at BioCreative IV. Database:
928 The Journal of Biological Databases and Curation. 2014; 2014:bau086.
929 doi:10.1093/database/bau086
- 930 56. Gobeill J, Pasche E, Vishnyakova D, Ruch P. Managing the data deluge: data-driven GO
931 category assignment improves while complexity of functional annotation increases. *Database:*
932 *The Journal of Biological Databases and Curation*. 2013;2013:bat041.
933 doi:10.1093/database/bat041.
- 934 57. Wei C-H, Peng Y, Leaman R, et al. Assessing the state of the art in biomedical relation extraction:
935 overview of the BioCreative V chemical-disease relation (CDR) task. Database: The Journal of
936 Biological Databases and Curation. 2016; 2016:baw032. doi:10.1093/database/baw032
- 937 58. Campos D, Matos S, Oliveira JL. Biomedical Named Entity Recognition: A Survey of Machine-
938 Learning Tools. Theory and Applications for Advanced Text Mining, InTech, 2012. doi:
939 10.5772/51066
- 940 59. Zheng JG, Howsmon D, Zhang B, Hahn J, McGuinness D, Hendler J, Ji H. Entity Linking for
941 Biomedical Literature. In: Proceedings of the ACM 8th International Workshop on Data and Text
942 Mining in Bioinformatics. New York, NY, USA: 2014. p.3-4.
- 943 60. Hoffart J, Yosef MA, Bordino I, Fürstenau H, Pinkal M, Spaniol M, et al. Robust disambiguation of
944 named entities in text. In Proc. of the Conf. on Empirical Methods in Natural Language Processing
945 (EMNLP '11). Association for Computational Linguistics, Stroudsburg, PA, USA: 2011. p.782-792.
- 946 61. Divita G, Zeng QT, Gundlapalli AV, Duvall S, Nebeker J, Samore MH. Sophia: A Expedient
947 UMLS Concept Extraction Annotator. AMIA Annual Symposium Proceedings. 2014;
948 2014:467-476.
- 949 62. Divita G, Carter MMS, Tran LT, Redd D, Zeng QT, Duvall S, Samore MH, Gundlapalli AV. v3NLP
950 Framework: Tools to Build Applications for Extracting Concepts from Clinical Text. Generating
951 Evidence & Methods to improve patient outcomes (eGEMs). 2016; 4(3).
- 952 63. Rodriguez-Esteban R. Biomedical Text Mining and Its Applications. Lewitter F, ed. PLoS
953 Computational Biology. 2009; 5(12):e1000597. doi:10.1371/journal.pcbi.1000597.
- 954 64. Cuzzola J, Jovanovic J, Bagheri E, Gasevic D. Evolutionary Fine-Tuning of Automated Semantic
955 Annotation Systems. Expert Systems with Applications. 2015; 42(20):6864–6877.
- 956 65. Goeriot L et al. Overview of the CLEF eHealth Evaluation Lab 2015. In: Mothe J. et al. (eds)
957 Experimental IR Meets Multilinguality, Multimodality, and Interaction. Lecture Notes in Computer
958 Science, vol. 9283. Springer, Cham: 2015, p.429-443.
- 959 66. Kelly L, Goeriot L, Suominen H, Névél A, Palotti J, Zuccon G. Overview of the CLEF eHealth
960 Evaluation Lab 2016. In: Fuhr N. et al. (eds) Experimental IR Meets Multilinguality, Multimodality,
961 and Interaction. CLEF 2016. Lecture Notes in Computer Science, Vol. 9822. Springer, Cham:
962 2016, p.255-266.

- 963 67. Van Mulligen E, Afzal Z, Akhondi SA, Vo D, Kors JA (2016). Erasmus MC at CLEF eHealth 2016:
964 Concept Recognition and Coding in French Texts. CLEF 2016 Online Working Notes, CEUR
965 Workshop Proceedings 1609: 2016. URL: <http://ceur-ws.org/Vol-1609/16090171.pdf>
- 966 68. Schuemie MJ, Jelier R, Kors JA. Peregrine: Lightweight Gene Name Normalization by Dictionary
967 Lookup. Proceedings of the BioCreAtIvE II Workshop; Madrid, Spain: 2007. p.131–133.
- 968 69. Rebholz-Schuhmann D et al. Entity Recognition in Parallel Multilingual Biomedical Corpora: The
969 CLEF-ER Laboratory Overview. In: Forner P., Müller H., Paredes R., Rosso P., Stein B. (eds)
970 Information Access Evaluation. Multilinguality, Multimodality, and Visualization. CLEF 2013.
971 Lecture Notes in Computer Science, Vol. 8138. Springer, Berlin, Heidelberg: 2013. p.353-367.
- 972 70. Bodnari A, Névéol A, Uzuner O, Zweigenbaum P, Szolovits P. Multilingual Named-Entity
973 Recognition from Parallel Corpora. Working Notes for CLEF 2013 Conference, Valencia, Spain,
974 September 23-26, 2013. CEUR Workshop Proceedings 1179: 2013. URL: <http://ceur-ws.org/Vol-1179/CLEF2013wn-CLEFER-BodnariEt2013.pdf>
- 976 71. Attardi G, Buzzelli A, Sartiano D. Machine Translation for Entity Recognition across Languages in
977 Biomedical Documents. In Working Notes for CLEF 2013 Conference, Valencia, Spain,
978 September 23-26, 2013. CEUR Workshop Proceedings 1179: 2013. URL: <http://ceur-ws.org/Vol-1179/CLEF2013wn-CLEFER-AttardiEt2013.pdf>
- 980 72. Lewin I, Clematide S. Deriving an English Biomedical Silver Standard Corpus for CLEF-ER.
981 Working Notes for CLEF 2013 Conference, Valencia, Spain, September 23-26, 2013. CEUR
982 Workshop Proceedings 1179: 2013. URL: <https://doi.org/10.5167/uzh-87213>
- 983 73. Kors JA, Clematide S, Akhondi SA, van Mulligen EM, Rebholz-Schuhmann D. (2015) A
984 multilingual gold-standard corpus for biomedical concept recognition: the Mantra GSC. *J Am Med*
985 *Inform Assoc* 2015;0:1–11. doi:10.1093/jamia/ocv037
- 986 74. GoPubMed. <http://www.gopubmed.org/>. Accessed 2 December 2016.
- 987 75. RIDeM - Repository for Informed Decision Making. <http://clinicalreferences.nlm.nih.gov/ridem/>.
988 Accessed 2 December 2016.
- 989 76. Ohta T, Pyysalo S, Tsuji J, Ananiadou S. Open-domain Anatomical Entity Mention Detection. In
990 Proceedings of ACL 2012 Workshop on Detecting Structure in Scholarly Discourse (DSSD). Jeju,
991 Republic of Korea: 2012. p.27-36.
- 992 77. Van Auken K et al. BC4GO: A Full-Text Corpus for the BioCreative IV GO Task. Database: The
993 Journal of Biological Databases and Curation 2014 (2014): bau074. PMC. Web. 7 July 2017.
- 994 78. Kafkas S, Lewin I, Milward D, van Mulligen E, Kors J, Hahn U, Rebholz-Schuhmann D. Calbc:
995 Releasing the final corpora. In: Proc. of the 8th International Conf. on Language Resources and
996 Evaluation (LREC'12). Istanbul, Turkey: 2012.
- 997 79. Li J et al. Annotating chemicals, diseases and their interactions in biomedical literature. In:
998 Proceedings of the Fifth BioCreative Challenge Evaluation Workshop, Sevilla, Spain: 2015; pp.
999 173–182

- 1000 80. Bada M, Eckert M, Evans D, Garcia K, Shipley K, Sitnikov D, et al. Concept Annotation in the
1001 CRAFT Corpus. *BMC Bioinformatics*. 2012; 13:161.
- 1002 81. Tanabe L, Xie N, Thom LH, Matten W, Wilbur WJ. GENETAG: a tagged corpus for gene/protein
1003 named entity recognition. *BMC Bioinformatics* 2005; 6(Suppl 1):S3. DOI:10.1186/1471-2105-6-
1004 S1-S3
- 1005 82. Kim JD, Ohta T, Tateisi Y, Tsujii, J. GENIA corpus—a semantically annotated corpus for bio-
1006 textmining. *Bioinformatics*. 2003; 19 (Suppl_1):i180-2.
- 1007 83. Uzuner Ö, South B, Shen S, DuVall S. 2010 i2b2/VA Challenge on Concepts, Assertions, and
1008 Relations in Clinical Text. *Journal of the American Medical Informatics Association*. 2011; 18:552-
1009 556. doi:10.1136/amiajnl-2011-000203
- 1010 84. Jin-Dong K, Tomoko O, et al TY. JNLPBA '04: Proceedings of the International Joint Workshop on
1011 Natural Language Processing in Biomedicine its Applications. Stroudsburg, PA, USA: Association
1012 for Computational Linguistics; 2004. Introduction to the bio-entity recognition task at JNLPBA; pp.
1013 70–75.
- 1014 85. Doğan RI, Leaman R, Lu Z. NCBI disease corpus: a resource for disease name recognition and
1015 concept normalization. *J. Biomed. Inform.* 2014; 47:1–10.
- 1016 86. ShARe/CLEF eHealth evaluation lab. SHARE-Sharing Annotated Resources. 2013.
1017 <https://sites.google.com/site/shareclefehealth/home>. Accessed 22 Nov 2016.
- 1018
- 1019

BioPortal (http://bioportal.bioontology.org/)	A major repository of biomedical ontologies, currently hosting over 500 ontologies, controlled vocabularies and terminologies. Its Resource Index provides an ontology-based unified index of and access to multiple heterogeneous biomedical resources (annotated with BioPortal ontologies).
DBpedia (http://wiki.dbpedia.org/)	"Wikipedia for machines", that is, a huge KB developed through a community effort of extracting information from Wikipedia and representing it in a structured format suitable for automated machine processing. It is the central hub of the Linked Open Data Cloud.
LLD - Linked Life Data (http://linkedlifedata.com/)	LLD platform provides access to a huge KB that includes and semantically interlinks knowledge about genes, proteins, molecular interactions, pathways, drugs, diseases, clinical trials and other related types of biomedical entities. It is part of the Linked Open Data Cloud (http://lod-cloud.net/)
NCBI Biosystems Database (https://www.ncbi.nlm.nih.gov/biosystems)	Repository providing integrated access to structured data and knowledge about biological systems and their components: genes, proteins, and small molecules. The NCBI Taxonomy contains the names and phylogenetic lineages of all the organisms that have molecular data in the NCBI databases.
OBO - Open Biomedical Ontologies (http://www.obofoundry.org/)	Community of ontology developers devoted to the development of a family of interoperable and scientifically accurate biomedical ontologies. Well known OBO ontologies include: <ul style="list-style-type: none"> • <i>Chemical Entities of Biological Interest (ChEBI)</i> - focused on molecular entities, molecular parts, atoms, subatomic particles, and biochemical roles and applications • <i>Gene Ontology (GO)</i> - aims to standardize the representation of gene and gene product attributes; consists of 3 distinct sub-ontologies: Molecular Function, Biological Process, and Cellular Component • <i>Protein Ontology (PRO)</i> - provides a structural representation of protein-related entities
SNOMED CT (http://www.ihtsdo.org/snomed-ct)	SNOMED CT is considered the world's most comprehensive and precise, multilingual health terminology. It is used for the electronic exchange of clinical health information. It consists of concepts, concept descriptions (i.e., several terms that are used to refer to the concept), and concept relationships.
UMLS (Unified Medical Language System) Metathesaurus (https://www.nlm.nih.gov/research/umls/knowledge_sources/metathesaurus/)	The most well-known and widely used knowledge source in the biomedical domain. It assigns a unique identifier (CUI) to each medical concept and connects concepts to each other thus forming a graph-like structure; each concept (i.e. CUI) is associated with its 'semantic type', a broad category such as Gene, Disease or Syndrome; each concept is also associated with several terms used to refer to that concept in biomedical texts; these terms are pulled from nearly 200 biomedical vocabularies. Some well-known vocabularies that have been used by biomedical semantic annotators include: <ul style="list-style-type: none"> • <i>Human Phenotype Ontology (HPO)</i> contains terms that describe phenotypic abnormalities encountered in human disease, and is used for large-scale computational analysis of the human phenome. • <i>Logical Observation Identifiers Names and Codes (LOINC)</i> provides standardized vocabulary for laboratory and

	<p>other clinical observations, and is used for exchange and/or integration of clinical results from several disparate sources.</p> <ul style="list-style-type: none"> • <i>Medical Subject Headings (MeSH)</i> is a controlled vocabulary thesaurus created and maintained by U.S. National Library of Medicine (NLM), and has been primarily used for indexing articles in PubMed. • <i>RxNorm</i> provides normalized names for clinical drugs and links between many of the drug vocabularies commonly used in pharmacy management and drug interaction software.
<p>UniProtKb/Swiss-Prot (http://www.uniprot.org/uniprot/)</p>	<p>Part of UniProtKB, a comprehensive protein sequence KB, which contains manually annotated entries. The entries are curated by biologists, regularly updated and cross-linked to numerous external databases, with the ultimate objective of providing all known relevant information about a particular protein.</p>

1020 **Table 1.** An overview of ontologies, thesauri and knowledge bases used by biomedical semantic annotation tools discussed in the paper

1021

1022

1023 **Table 2.** Example application cases of biomedical semantic annotation tools

Application Case (AC)	The role of semantic annotation tool in the AC	Biomedical resources relevant for the AC (or representative examples, if multiple)
Semantic search of biomedical tools and services [6]	Sematic search of biomedical tools and services enabled by semantic annotation of users' (free-form) queries with concepts from UMLS Metathesaurus	Catalogs of and social spaces created around biomedical tools and services, e.g.: - myExperiment (http://www.myexperiment.org/) - BioCatalogue (https://www.biocatalogue.org/)
Semantic search of domain specific scientific literature [74]	Semantic annotation of PubMed entries with ontological concepts related to genes and proteins	Ontologies used for the annotation of biomedical references (PubMed entries): - Gene Ontology - GO (http://geneontology.org/) - Universal Protein Resource - UniProt (http://www.uniprot.org/uniprot/)
Improved clinical decision making [75]	Extraction of key clinical concepts (UMLS-based) required for supporting clinical decision making; the concepts are extracted from biomedical literature and clinical text sources	Sources of biomedical texts used to support decision making: - PubMed Central (PMC) Open Access Subset (https://www.ncbi.nlm.nih.gov/pmc/tools/openftlist/) - MEDLINE abstracts (https://www.nlm.nih.gov/bsd/pmresources.html)
Unambiguous description of abbreviations [10]	Extended (long) forms of abbreviations are matched against both UMLS and DBpedia concepts, thus not only disambiguating the long forms, but also connecting UMLS and DBpedia KBs	Allie - a search service for abbreviations and their long forms (http://allie.dbcls.jp/)

1024

1025

1026 **Table 3.** General purpose biomedical semantic annotation tools (Part I)

	cTAKES [4]	NOBLE Coder [20]	MetaMap [31,32]	NCBO annotator [14]
Modularity / configuration options	Modular text processing pipeline	Vocabulary (terminology); Term matching options and strategies	Text processing pipeline; Vocabulary (terminology); Term matching options and strategies	Vocabulary (terminology); Term matching options
Disambiguation of terms	Enabled through integration of the YTEX component [8]	Instead of through WSD, it uses heuristics to choose one concept among candidate concepts for the same piece of input text	Supported; based on: - removal of word senses based on a manual study of UMLS ambiguity - a WSD algorithm that chooses a concept with the most likely semantic type for a given context	Not supported
Vocabulary (terminology)	Subset of UMLS, namely SNOMED CT and RxNORM	Several pre-built vocabularies, based on subsets of UMLS (e.g. SNOMED CT, MeSH, RxNORM)	UMLS Metathesaurus	UMLS Metathesaurus and BioPortal ontologies (over 330 ontologies)
Speed*	Suitable for real-time processing	Suitable for real-time processing	Better for off-line batch processing	Suitable for real-time processing
Implementation form	Software (Java) library; Stand-alone application	Software (Java) library; Stand-alone application	Software library; originally version in Prolog; Java implementation, known as MMTX, is also available	RESTful Web service
Availability	open source; available under Apache License, v.2.0	open-source; available under GNU Lesser General Public License v3	open source; terms and conditions at: https://metamap.nlm.nih.gov/MMTnCs.shtml	closed source, but freely available
Specific features	Better performance on clinical texts than on biomedical scientific literature (its NLP components are trained on clinical texts)	Offers user interface for creating custom terminologies (to be used for annotation) by selecting and merging elements from several different thesauri / ontologies	Primarily developed for annotation of biomedical literature (MEDLINE / PubMed citations); performs better on this kind of text than clinical notes	It uses MGrep term-to-concept matching tool to get primary set of annotations; these are then extended using different forms of ontology-based semantic matching
URL	http://ctakes.apache.org/	http://noble-tools.dbmi.pitt.edu/	https://metamap.nlm.nih.gov/	https://bioportal.bioontology.org/annotator

1027 **Table 3.** General purpose biomedical semantic annotation tools (Part II)

	BeCAS [36]	Whatizit [38]	ConceptMapper [21]	Neji [40]
Modularity / configuration options	Semantic types (i.e. types of entities to annotate)	pre-built pipelines for several biomedical types (see Specific features)	Text processing pipeline; Term matching options and strategies	Modular text processing pipeline
Disambiguation of terms	No information available	Not supported	Not supported	Instead of through WSD, it uses a set of heuristics rules to identify and remove annotations of lower importance
Vocabulary (terminology)	Custom built vocabulary by using concepts from multiple sources, such as UMLS, NCBI BioSystems, ChEBI, and the Gene Ontology.	The use of the vocabulary depends on the type of entity a pipeline is specialized for (e.g. NCBI KB for species, or Gene Ontology for genes)	General purpose dictionary lookup tool, not tied to any vocabulary	Not tied to any particular vocabulary
Speed*	Suitable for real-time processing	Suitable for real-time processing	Suitable for real-time processing	Suitable for real-time processing
Implementation form	Software (Python) library; RESTful Web service; Javascript widget	SOAP Web service	Software (Java) library; part of the UIMA NLP framework [28]	RESTful Web service
Availability	open source; available under Attribution-NonCommercial 3.0 Unported license	closed source, but freely available	open source; available under Apache License, v.2.0	open source; available under Attribution-NonCommercial 3.0 Unported license
Specific features	Primarily aimed for annotation of biomedical research papers; focused on annotation of several (11) types of biomedical entities, including species, microRNAs, enzymes, chemicals, drugs, diseases, etc.	Offers several pre-built pipelines for specific entity types; e.g. whatizitGO identifies proteins based on the Gene Ontology (GO), while whatizitChemical annotates chemical entities based on ChEBI	Not specifically developed for the biomedical domain, but is a general purpose dictionary lookup tool	Includes modules for both ML and dictionary-based annotation; can automatically combine annotations generated by different modules
URL	http://bioinformatics.ua.pt/becas/	http://www.ebi.ac.uk/webservices/whatizit	https://uima.apache.org/sandbox.html#concept.mapper.annotator	https://github.com/BMDSoftware/neji

1028 * Note that speed estimates are based on the experimental results reported in the literature; those experiments were done with corpora of up to 200 documents (paper
 1029 abstracts or clinical notes); the given estimates might not hold for significantly larger corpora

1030 **Table 4.** Corpora used for evaluation of biomedical semantic annotators. The table includes corpora that were used in the reported use cases
 1031 (Section 2, Table 2), and/or benchmarking of the discussed tools (Section 3.2 and Section 4).

<p>AnEM - Anatomical Entity Mention [76]</p>	<p>The corpus consists of 500 documents selected randomly from citation abstracts and full-text biomedical research papers (from PubMed); it is manually annotated (over 3,000 annotations) with anatomical entities. The corpus is available under the open CC-BY-SA license.</p> <p>URL: http://www.nactem.ac.uk/anatomy/</p>
<p>BC4GO [77]</p>	<p>The corpus, developed for the BioCreative IV shared task, consists of 200 articles (over 5,000 text passages) from Model Organism Databases; these articles were manually annotated with more than 1356 distinct GO terms. In addition to the core elements of GO annotations - a gene or gene product, a GO term, and a GO evidence code - the corpus also includes the GO evidence text.</p> <p>URL: http://www.biocreative.org/tasks/biocreative-iv/track-4-GO/</p>
<p>CALBC - Collaborative Annotation of a Large Biomedical Corpus [78]</p>	<p>A very large, publicly shared corpus of Medline abstracts automatically annotated with biomedical entities; the small corpus comprises ~175K abstracts, whereas the big one consists of more than 714K abstracts; since annotations were not made by humans but several annotation systems (and then aggregated), it is referred to as "silver standard".</p> <p>URL: http://www.ebi.ac.uk/Rebholz-srv/CALBC/corpora/resources.html</p>
<p>Chemical Disease Relation (CDR) [79]</p>	<p>The corpus, developed for the BioCreative V shared task, consists of 1,500 PubMed articles with 4,409 annotated chemicals, 5,818 diseases, and 3,116 chemical-disease interactions. MeSH is used as the controlled vocabulary.</p> <p>As BC4GO, this corpus is available exclusively for scientific, educational, and/or non-commercial purposes.</p> <p>URL: http://www.biocreative.org/tasks/biocreative-v/track-3-cdr/</p>

<p>CRAFT - the Colorado Richly Annotated Full Text corpus [80]</p>	<p>Publicly available, human annotated (gold standard) corpus of full-text biomedical journal articles; it consists of 67 document and 87,674 human annotations</p> <p>URL: http://bionlp-corpora.sourceforge.net/CRAFT/</p>
<p>GENETAG [81]</p>	<p>Publicly available corpus of 20K Medline sentences manually annotated with gene/protein names. Part of the corpus (15K sentences) was used for the BioCreative I challenge (Gene Mention Identification task), and the rest (5K sentences) was used as test data for BioCreative II competition (Gene Mention Tagging Task). URL: https://github.com/openbiocorpora/genetag</p> <p>An updated version of this corpus, named GENETAG-05, is part of a broader MedTag annotated corpus that was used in the BioCreative I challenge; it is available at: ftp://ftp.ncbi.nlm.nih.gov/pub/lsmith/MedTag/</p>
<p>GENIA [82]</p>	<p>Open access manually annotated corpora consisting of 2000 Medline abstracts (400,000+ words) with almost 100,000 annotations for biological terms. Terms are annotated with concepts from the GENIA ontology, a formal model of cell signaling reactions in humans (the ontology is provided together with the corpus).</p> <p>Available from the following repository: http://corpora.informatik.hu-berlin.de/</p>
<p>2010 i2b2/VA corpus [83]</p>	<p>The corpus consists of manually annotated de-identified clinical records (discharge summaries and progress reports) from three medical centers. It was originally created for the 2010 i2b2/VA NLP challenge to support 3 kinds of tasks: extraction of medical concepts from patient reports; assigning assertion types to medical problem concepts; and determining the type of relation between medical problems, tests, and treatments. The corpus consists of 394 annotated training reports, 477 annotated test reports, and 877 unannotated reports.</p> <p>The corpus is made available to the research community from https://i2b2.org/NLP/DataSets under data use agreements.</p>

JNLPBA [84]	<p>A publicly available manually annotated corpus originally created for the Bio-Entity Recognition Task at BioNLP/NLPBA 2004. The training set consists of 2000 Medline abstracts extracted from the GENIA Version 3.02 corpus; the data set is annotated with five entity types: Protein, DNA, RNA, Cell_line, and Cell_type. The test set consists of 404 annotated Medline abstracts, also from the GENIA project; a half of this data set is from the same domain as that of the training data, whereas the other half is from the super domain of blood cells and transcription factors.</p> <p>URL: http://www.geniaproject.org/shared-tasks/bionlp-jnlpba-shared-task-2004</p>
NCBI Disease corpus [85]	<p>Publicly available, manually annotated corpus of 793 PubMed abstracts; 6,892 disease mentions are annotated with concepts from Medical Subject Headings (MeSH) and Online Mendelian Inheritance in Man (OMIM) vocabularies.</p> <p>URL: https://www.ncbi.nlm.nih.gov/CBBresearch/Dogan/DISEASE/</p>
Mantra Gold Standard Corpus [73]	<p>Publicly available multilingual gold-standard corpus for biomedical concept recognition. It includes text from different types of parallel corpora (Medline abstract titles, drug labels, biomedical patent claims) in English, French, German, Spanish, and Dutch. It contains 5,530 annotations based on a subset of UMLS that covers a wide range of semantic groups.</p> <p>URL: http://biosemantics.org/index.php/resources/mantra-gsc</p>
ShARe - Shared Annotated Resources [86]	<p>Gold standard corpus of de-identified clinical free-text notes; it includes 199 documents and 4,211 human annotations; originally prepared for the ShARe/CLEF eHealth Evaluation Lab focused on NLP and information retrieval tasks for clinical care.</p> <p>URL: https://sites.google.com/site/shareclefehealth/data</p>

