1 **NERO: A Biomedical Named-entity (Recognition) Ontology with a Large,**
2 **Annotated Corpus Reveals Meaningful Associations Through Text Embedding**
3

4 [1,2]Kanix Wang
5 [3]Robert Stevens
6 [4]Halima Alachram
7 [5]Yu Li
8 [6]Larisa Soldatova
9 [7]Ross King
10 [3,8]Sophia Ananiadou
11 [3,8]Maolin Li
12 [3,8]Fenia Christopoulou
13 [9]Jose Luis Ambite
14 [9]Sahil Garg
15 [9]Ulf Hermjakob
16 [9]Daniel Marcu
17 [9]Emily Sheng
18 [4]Tim Beißbarth
19 [10]Edgar Wingender
20 [9]Aram Galstyan
21 [5]Xin Gao
22 [11]Brendan Chambers
23 [2,12]Bohdan B. Khomtchouk
24 [11]James A. Evans
25 [1,2,12,13*]Andrey Rzhetsky
26

27 [1]The Committee on Genetics, Genomics, and Systems Biology, University of Chicago,
28 Chicago, IL 60637, US;
29 [2]The Institute of Genomics and Systems Biology, University of Chicago, Chicago, IL
30 60637, US;
31 [3]Depatment of Computer Science, University of Manchester, M13 9PL, UK;
32 [4]Institute of Medical Bioinformatics, University of Göttingen, Goldschmidtstrasse 1,
33 37077 Göttingen, Germany.
34 [5]Computational Bioscience Research Center;
35 Computer, Electrical and Mathematical Sciences and Engineering Division;
36 King Abdullah University of Science and Technology (KAUST)
37 Thuwal, 23955, Saudi Arabia;
38 [6]Goldsmiths, University of London, 8 Lewisham Way, New Cross, London SE14 6NW,
39 UK;
40 [7]Department of Chemical Engineering and Biotechnology, University of Cambridge,
41 Philippa Fawcett Dr, Cambridge CB3 0AS, United Kingdom
42 Alan Turing Institute, 96 Euston Rd, Somers Town, London NW1 2DB, United Kingdom
43 Department of Biology and Biological Engineering, Chalmers University of Technology,
44 SE-412 96 Göteborg, Sweden.
45 [8]National Centre for Text Mining, University of Manchester, M1 7DN, UK;

46  [9]The Information Sciences Institute, University of Southern California, Marina del Rey,
47  CA 90089, US;
48  [10]geneXplain GmbH, Am Exer19b, 38302 Wolfenbüttel, Germany;
49  [11]Knowledge Lab, Department of Sociology, University of Chicago, IL 60637, US;
50  [12]Department of Medicine, University of Chicago, Chicago, IL 60637, US;
51  [13]Department of Human Genetics, University of Chicago, Chicago, IL 60637, US;
52
53  [*]Corresponding author, andrey.rzhetsky@uchicago.edu.
54

55  **Machine reading is essential for unlocking valuable knowledge contained in the**
56  **millions of existing biomedical documents. Over the last two decades [1,2], the most**
57  **dramatic advances in machine-reading have followed in the wake of critical**
58  **corpus development[3]. Large, well-annotated corpora have been associated with**
59  **punctuated advances in machine reading methodology and automated**
60  **knowledge extraction systems in the same way that ImageNet [4] was fundamental**
61  **for developing machine vision techniques. This study contributes six**
62  **components to an advanced, named-entity analysis tool for biomedicine: (a) a**
63  **new, Named-Entity Recognition Ontology (NERO) developed specifically for**
64  **describing entities in biomedical texts, which accounts for diverse levels of**
65  **ambiguity, bridging the scientific sublanguages of molecular biology, genetics,**
66  **biochemistry, and medicine; (b) detailed guidelines for human experts annotating**
67  **hundreds of named-entity classes; (c) pictographs for all named entities, to**
68  **simplify the burden of annotation for curators; (d) an original, annotated corpus**
69  **comprising 35,865 sentences, which encapsulate 190,679 named entities and**
70  **43,438 events connecting two or more entities; (e) validated, off-the-shelf, named-**
71  **entity recognition automated extraction, and; (f) embedding models that**
72  **demonstrate the promise of biomedical associations embedded within this**
73  **corpus.**
74

75  Even the relatively specialized subfields of present-day biology and medicine are facing
76  a deluge of accumulating research articles, patents, and white papers. It is increasingly
77  difficult to stay up-to-date with contemporary biomedicine without the use of
78  sophisticated machine reading (MR) tools. MR tool development, in turn, has been
79  limited by the availability of biomedical corpora carefully annotated by experts. This is
80  especially true with respect to information extraction, such as named entity recognition
81  and relation or event extraction. Although several corpora have been developed for
82  specialized biomedical subdomains, the need for a corpus that can bridge biological,
83  general scientific, environmental, and clinical scientific sub-languages is greater than
84  ever before.
85        Unfortunately, the annotation of natural science texts is more challenging than in
86  other domains. Biomedical language is replete with ambiguity distinct from that
87  observed in news articles or informal text online. When a word or phrase's semantic
88  meaning is clearly separated (*the east bank of the Danube* versus *Deutsche Bank*), we
89  can implement automated sense disambiguation using machine learning tools. In
90  biomedical texts, however, alternative meanings are not always clearly separated. The
91  problem is not that a phrase can refer to several distinct, real-world entities in different

92    contexts, but that the scientists writing articles typically do not separate competing,
93    close meanings. For example, in some biomedical contexts, a named entity may refer to
94    a *gene* or a *protein* with nearly equal probability; for example, "a mutant hemoglobin $\alpha_2$"
95    can refer to either a gene or a protein. If the author meant *gene-or-protein A,* and we
96    force an annotator to choose either interpretation *gene A* or *protein A*, the resulting
97    annotation is of limited utility because the choice between *gene* and *protein* is random if
98    the meanings are equally likely based on context. Ideally, a specialized ontology of text
99    entities would allow an annotator to choose the proper level of annotation granularity
100    (*gene-or-protein*, in this example), minimizing the need for forced, random decisions. To
101    the best of our knowledge, there is no biomedical ontology that meets the requirements
102    for capturing semantic ambiguity. We aimed to fill this gap by developing a specialized,
103    variable-level meaning resolution ontology, a carefully curated corpus, along with
104    corpus annotation tools, and a collection of text embedding analyses to evaluate our
105    annotated corpus.
106          Our new ontology, called NERO, short for Named-entity Recognition Ontology,
107    attempts to minimize unwarranted, arbitrary annotative semantic label assignments in
108    text entities, see Figure 1. NERO captures named entities, starting with most broad and
109    vague concepts close to the taxonomy's root, finishing with the most narrow and
110    concrete concepts at the taxonomy's leaves. Hence, *DomainEntity*–and all ambiguous
111    semantic classes–correspond to NERO's taxonomy root. The basic division thereafter is
112    between *TextEntity* and *AbstractEntity*, where *TextEntity* further splits into *NamedEntity*,
113    *NamedEntityGroup*, *Relationship,* and *Pronoun*. After *NamedEntity*, the hierarchy
114    reflects that which is written in biological entity descriptions, rather than in those entities'
115    lexical representation. NERO defines ambiguous concepts, such as *GeneOrProtein*,
116    which subsumes both *Gene* and *Protein* using the following axiom: *EquivalentTo:*
117    *'Gene' or 'Protein.'* There are no biological entities that are either a gene or a protein,
118    but there are lexical entities that can belong to either named entity class. NERO uses
119    this pattern to express appropriate ambiguity regarding text entities, preserving
120    uncertainty from the text. In this way, NERO classes represent textual instances and not
121    the actual biological entities to which these instances refer. It is, therefore, straight-
122    forward to link between the lexical and biological entity through a relationship such as *'is*
123    *about'*. So, the NERO class *Protein 'is about'* some specific concept *'protein'* in an
124    ontology pointing to real biological entities, such as the Protein Ontology [5].
125          Striving to make the ontology practically useful, we designed guidelines for
126    annotators making decisions in annotating text entities, available in the *Supplementary*
127    *Data*. Furthermore, by recruiting a team of postdoctoral-level experts, we annotated a
128    large biomedical corpus to enable a broad range of natural language processing and
129    biomedical machine learning tasks. Our annotations span 35,865 unique sentences,
130    8,650 of which were annotated by multiple annotators with remarkably high inter-
131    annotator agreement (see Table 1). In our annotated corpus, we aimed to encompass
132    all entity types that might occur in biomedical literature. In addition to named entities,
133    our ontology captures *events* which represent relationships between biomedical
134    concepts. The frequencies of all diverse entity types in our corpus are shown in Figure
135    2A; Figure 2B shows the frequencies of relations represented in the taxonomy. The
136    most frequent entity type is *GeneOrProtein,* which accounts for 14.7 percent of all
137    named entities in the corpus (see Figure 2A). The second most populous category is

138    *Process,* with nine percent tagged. *Process* has six sub-concepts and almost half of
139    *Process* instances (49.7 percent) are annotated as more specific sub-concepts; the
140    *BiologicalProcess* and the *MolecularProcess* are the fifth and seventh most frequent
141    entity types (see Figure 2). Entity type frequencies follow a heavy-tail distribution, with
142    the least frequent types being *Journal, Unit,* and *Citation* (see Figure 2). In addition to
143    190,679 named entities, we annotated 43,438 action terms, events connecting two or
144    more entities. The most annotated action term is *bind,* accounting for 28.4% of all
145    actions, see *Supplementary Figure 1.* When we normalize the action terms and
146    combine actions such as *bind, binds,* and *binding,* the normalized action *bind* accounts
147    for 31.8% of all actions, as shown in *Supplementary Figure 1.* We deployed a package
148    called NERO-nlp for researchers interested in diving deeper into our annotated corpus;
149    the installation guides and scripts are available online at https://pypi.org/project/NERO-
150    nlp and https://github.com/Bohdan-Khomtchouk/NERO-nlp respectively.
151        Below, we present two practical applications of our ontology and text annotations:
152    1) Machine learning experiments, which automatically identify named entities, and; (2)
153    Word embedding experiments, which leverage the automated discovery of semantic
154    relationships among real-world concepts referenced by a text's named entities.
155        *Machine learning experiments*: Using NERsuite [6], we conducted a ten-fold cross-
156    validation, dividing the corpus into training and test subsets. The classification results
157    are presented in *Supplemental Table 1*. The overall automated named entity recognition
158    performance is moderate, with 54.9% precision, 37.3% recall and a 43.4% $F_1$ score.
159    The best performance class, *GeneOrProtein,* had baseline results of 67.0% precision,
160    65.3% recall, and a 66.2% $F_1$ score. In addition to the default baseline implementation
161    of NERsuite, we added additional features in the training process to improve its
162    performance [7]. These are dictionary features derived from lookups in technical term
163    dictionaries. The classifier with dictionary features manifests 54.7% precision, 37.9%
164    recall and a 43.8% $F_1$ score. We observed a scant 0.35% increase in $F_1$ score from
165    adding dictionary features. We then implemented an ensemble method called stacking,
166    where we trained a higher-level model to learn how to best combine contributions from
167    each base model. The base model in this case is the baseline model from NERsuite.
168    Stacking yielded a 0.27% increase in $F_1$ score compared to baseline results. While
169    ensemble methods are commonly used to boost model accuracy by combining the
170    predictions of multiple machine learning models, choices of second-level and base
171    models can influence the amount of improvement in model accuracy. The overall
172    performance statistics are shown in *Supplementary Table 2*. As our corpus is made
173    public with this study's publication, we hope that other researchers will use this training
174    data to achieve core MR task performance that surpasses our initial experiments.
175        To examine how NERsuite performs in comparison to other popular open-source
176    Named-Entity recognition tools, we trained a custom NER model on our annotated
177    corpus using spaCy [8]. We evaluated the trained model on the test subset, which
178    consists of a random 10% sample from the corpus. Classification results are presented
179    in *Table 3*. Overall automated named entity recognition performance is low, with 30.9%
180    precision, 8.6% recall and a 13.4% $F_1$ score. The best performance class,
181    *GeneOrProtein,* had results of 45.1% precision, 36.4% recall, and a 40.3% $F_1$ score.
182    These statistics indicate a much poorer performance of spaCy compared to that of
183    NERsuite.

184   To help explain the huge difference in performance between NERsuite and
185 spaCy, we considered the set of input features used by each tool for insight. NERsuite's
186 baseline implementation uses an extra set of input features including the lemma, POS-
187 feature and chunk-feature, whereas our custom spaCy NER model only relies on
188 character offsets and entity labels. There is potential for further customizing spaCy's
189 processing pipelines by adding more components such as tagger and parser [8], but no
190 established approaches in this regard have been made available partly because
191 spaCy's model architecture is different from those of other popular NER tools. We also
192 observed that some entities classes, such as Gene and Protein, have zero values for
193 precisions, recalls and $F_1$ scores, which likely translate to no correct classifications
194 made for those entities. The zero values occur partly due to the relatively smaller
195 number of tokens for those entity classes in the training set, and as a result, the trained
196 NER model generalized poorly on the minority class entities in the test subset.
197   Due to spaCy's computational demands, we did not conduct 10-fold cross-
198 validation. NERsuite provides a well-integrated pipelined system where training a new
199 model consists of a few lines of code. In addition, NERsuite has a demonstrated record
200 [5] on two biomedical tasks, the BioCreative2 gene mention recognition task and the
201 NLPBA 2004 named entity recognition task. Therefore, one could argue that it offers an
202 advantage over spaCy for NLP tasks in specialized domains such as biomedicine.

203   We've also identified another package called scispaCy [9] that contains spaCy
204 models for processing biomedical, scientific or clinical text. SciSpaCy acts as an
205 extension to spaCy and provides a set of practical tools for text processing in the
206 biomedical domain [9]. In particular, scispaCy includes a set of spaCy NER models
207 trained on popular biomedical corpora, which covers entity types such as chemicals,
208 diseases, cell types, proteins and genes. As an extension to spaCy, it also has the
209 flexibility for users to train a custom NER model from scratch or update the existing
210 NER models with users' own training data. Since our NER ontology adopts a more
211 diverse and detailed annotation methodology for named entity types, it will be
212 challenging to update scispaCy's pretrained named entity recognizer with our annotated
213 corpora.

|  | p | r | f |
|---|---|---|---|
| Cell | 16.88 | 5.10 | 7.83 |
| CellComponent | 35.71 | 4.44 | 7.91 |
| GeneOrProtein | 45.11 | 36.44 | 40.31 |
| Organism | 16.99 | 5.53 | 8.34 |
| Disease | 11.79 | 8.12 | 9.61 |
| Drug | 11.11 | 1.20 | 2.17 |
| SmallMolecule | 0.00 | 0.00 | 0.00 |
| BiologicalProcess | 8.02 | 1.88 | 3.05 |
| MolecularProcess | 12.67 | 2.45 | 4.10 |
| Gene | 0.00 | 0.00 | 0.00 |
| Protein | 0.00 | 0.00 | 0.00 |
| BodyPart | 15.17 | 5.54 | 8.12 |
| AminoAcid | 12.50 | 0.88 | 1.64 |

Table 3: Experimental results using spaCy for NER evaluated on 10% of the corpus

*Word embedding experiments:* Semantic associations, automatically extracted from text using neural network embedding operations, can function as a kind of "digital double" of real-world phenomena embedded in text, facilitating inferences that were previously imagined only possible from the original experimental data. For example, word embeddings built from chemical and material science texts predict much of the subsequent decades' material discoveries[8], just as the corpus of molecules can recover the periodic table[9], and texts are able to recover the subtle, psychological and sociological biases of cultures that produced them [10,11]. We used word embedding models to evaluate the biomedical veracity of NERO and its text annotation. Embedding models like Google's $word2vec$ [12,13] initially received substantial attention based on their capacity to solve analogy problems and automatically capture deep semantic relationships among concepts. Building on these capacities [10,14,15], we proposed a general method for constructing meaningful dimensions by taking the arithmetic mean of word vectors representing antonyms along a dimension and using them to diagnose their meanings. This approach has been widely validated [15-19], and we employed it here to construct and compare the meanings embedded in NERO and our annotated corpus with ground truth data about drugs and diseases. In order to evaluate word embeddings based on NERO, we identified two disease properties —(1) severity and (2) gender specificity—and likewise two therapeutic drug properties —(1) toxicity and (2) expense—not directly present in text, but highly relevant to diagnosis and treatment, and on which text-independent ground truth data exists.

239   We embedded named entities associated with diseases and drugs into a high-
240   dimensional space in which every NERO term was assigned a 300-dimensional vector,
241   (see Figure 3 for a three-dimensional projection of this embedding), along with a
242   selection of diseases and medications used to treat them. We then compared drug and
243   disease projections into the embedding dimensions for severity, gender, toxicity, and
244   expense with ground truth about these qualities. We constructed the severe-mild axis
245   with the following contrasting term pairs: (harmful, beneficial), (serious, benign), (life-
246   altering, common), (disruptive, undisruptive), (dying, recovering), (dangerous, safe),
247   (threatening, low-priority), (high mortality, harmless), (costly, cheap), (hospitalized, self-
248   administered ), (hospital, work), (debt, savings), (low quality of life, undisruptive), and
249   (hazard, routine). Then we compared disease projection in this dimension with World
250   Health Organization data on the burden of living with each of those diseases (DALYs [20])
251   and found a correlation of 0.329 ($p$=0.0614, $n$=33). We then constructed a gender
252   dimension with similarly contrasting pairs: (male, female), (prostate, ovary), (penile,
253   uterine), (penis, uterus), (man, woman), (men, women), (masculine, feminine), (he,
254   she), (him, her), (his, hers), (boy, girl), and (boys, girls). We compared the disease
255   projection in this gender dimension with the prevalence of those diseases for men and
256   women from a substantial sample of doctor-patient insurance records capturing
257   approximately 47% of all of U.S. doctor-patient visits between 2003 and 2011 and found
258   a correlation of 0.436 ($p$=1.46 $\times$ 10$^{-13}$, $n$=261).
259   Together, these patterns suggest that not only does NERO facilitate efficient and
260   accurate concept-by-concept annotation, but that the distribution of biomedical
261   properties underlying NERO-annotated texts have emergent validity and predict data
262   patterns not explicitly present in biomedical articles. Following the same pattern, we
263   projected medications onto a toxicity axis composed from: (harmful, beneficial), (toxic,
264   nontoxic), and (noxious, benign) and an expense dimension anchored by: (expensive,
265   inexpensive), (costly, cheap), (brand, generic), and (patented, off-patent). The
266   correlation of drug projections onto the toxicity dimension correlates at 0.32 ($p$=1.1 $\times$ 10$^{-4}$
267   ) with the median lethal dose, or dose required to kill 50% of subjects as documented
268   in the LD50 database [21]. Finally, the correlation of drug projections into an expense
269   dimension and the price of each drug as listed in the IBM MarketScan database [22] was
270   0.42 ($p$=1.5 $\times$ 10$^{-15}$) (see Figure 4). When a disease projects low in the *male – female*
271   dimension, it is much more likely to afflict women than men, such as ornithosis and
272   related infectious diseases. When a disease projects high in the *serious – benign*
273   dimension like leprosy, it is likely to incur substantial suffering. When a medication
274   projects high in the *toxic – nontoxic* dimension, such as Riluzole, a treatment for
275   amyotrophic lateral sclerosis with potential severe side effects ranging from unusual
276   bleeding to nausea and vomiting. Drug projections high in the *expensive – inexpensive*
277   dimension suggest a stiff medical bill, as in the case of Simvastatin, which is used to
278   reduce the risk of heart attack and stroke, and which, before it went off-patent, cost
279   hundreds of dollars per bottle. The robust accuracy of these associations suggest that
280   for qualities on which we do not have relevant or inexpensive data outside text,
281   associations from text represent a significant signal for biomedical research and can be
282   considered robust hypotheses meriting empirical study.
283   This study's main limitation is that, even though our NERO ontology aimed to
284   cover all entities contained in the biomedical research literature, we did not cover all

285 levels of granularity in classifying entities. Moreover, while the major concepts are well-
286 annotated, several concept types were not well-represented because of the heavy-tail
287 distribution of ontological class frequencies. In addition, we note that satisfactory results
288 of Named-entity Recognition (NER) rely heavily on a large quantity of hand-annotated
289 data, which is often costly in terms of time and resources spent. Therefore, adoption of
290 semi-supervised learning methods, which incorporates unlabeled data to improve
291 learning accuracy, could reduce the need for manual annotation [23].
292 While there is popular belief that pretraining on general-domain text can be
293 helpful for developing domain-specific language models, a recent study has shown that
294 for specialized domains, such as biomedicine, pretraining on in-domain text from
295 scratch offers noticeable improvements in model accuracy compared to continual
296 pretraining of general-domain language models [24]. Therefore, we trained on our
297 annotated corpus from scratch using in our machine learning experiments [25].
298 The resources offered in our study can be applied to a wide range of scientific
299 problems. First, the proposed NERO ontology can facilitate more robust and accurate
300 large-scale text mining of biomedical literature. As discussed above, NERO is the first
301 knowledge graph in this field, accounting for context-relevant levels of ambiguity. Graph
302 neural networks [26] can leverage such prior knowledge from human experts for learning
303 embedding of biomedical entities, which is likely to preserve both semantic meaning in
304 the original literature and domain knowledge from human experts. Second, researchers
305 can combine the curated corpus from this study with self-supervised learning [27]. Such a
306 learning scenario can utilize the unlabeled data in a supervised way by predicting part of
307 the sentence using the rest of the sentence. The annotated corpus from this study can
308 be used to fine-tune language models, orienting them for critical biomedical tasks.

309
310 **Competing interests**
311 The authors declare that they have no competing financial interests.
312
313 **Acknowledgments**

321
322
323

324
### References
326
327  1    Banko, M. & Brill, E. in *Proceedings of the 39th Annual Meeting on Association for*
328       *Computational Linguistics*   26-33 (Association for Computational Linguistics,
329       Toulouse, France, 2001).
330  2    Halevy, A., Norvig, P. & Pereira, F. The Unreasonable Effectiveness of Data. *Ieee*
331       *Intelligent Systems* **24**, 8-12, (2009).
332  3    Dogan, R. I., Leaman, R. & Lu, Z. NCBI disease corpus: a resource for disease name
333       recognition and concept normalization. *J Biomed Inform* **47**, 1-10, (2014).
334  4    Deng, J. *et al.* in *2009 IEEE Conference on Computer Vision and Pattern Recognition.*
335       248-255.
336  5    Natale, D. A. *et al.* The Protein Ontology: a structured representation of protein
337       forms and complexes. *Nucleic Acids Res* **39**, D539-545, (2011).
338  6    Wijffels, J. & Okazaki, N. *crfsuite: Conditional Random Fields for Labelling Sequential*
339       *Data in Natural Language Processing based on CRFsuite: a fast implementation of*
340       *Conditional Random Fields (CRFs)*, <https://github.com/bnosac/crfsuite> (2007-
341       2018).
342  7    Friedrich, C., Revillion, T., Hofmann-Apitius, M. & Fluck, J. Biomedical and chemical
343       named entity recognition with conditional random fields: The advantage of
344       dictionary features. (2006).
345  8    Honnibal, M. & Montani, I. *spaCy 2: Natural language understanding with Bloom*
346       *embeddings, convolutional neural networks and incremental parsing.*,
347       <https://spacy.io> (2017).
348  9    Neumann, M., King, D., Beltagy, I. & Ammar, W. *ScispaCy: Fast and Robust Models for*
349       *Biomedical Natural Language Processing.*  (2019).
350  10   Caliskan, A., Bryson, J. J. & Narayanan, A. Semantics derived automatically from
351       language corpora contain human-like biases. *Science* **356**, 183-186, (2017).
352  11   Garg, N., Schiebinger, L., Jurafsky, D. & Zou, J. Word embeddings quantify 100 years
353       of gender and ethnic stereotypes. *Proc Natl Acad Sci U S A* **115**, E3635-E3644,
354       (2018).
355  12   Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S. & Dean, J. Distributed
356       representations of words and phrases and their compositionality. *Advances in neural*
357       *information processing systems*, 3111-3119, (2013).
358  13   Mikolov, T., Chen, K., Corrado, G. & Dean, J. Efficient estimation of word
359       representations in vector space. *arXiv*, 1301.3781, (2013).
360  14   Austin, C. K., Taddy, M. & Evans, J. A. The Geometry of Culture: Analyzing the
361       Meanings of Class through Word Embeddings (vol 84, pg 905, 2019). *American*
362       *Sociological Review* **85**, 197-197, (2020).
363  15   Kozlowski, A. C., Taddy, M. & Evans, J. A. The Geometry of Culture: Analyzing the
364       Meanings of Class through Word Embeddings. *American Sociological Review* **84**,
365       905-949, (2019).
366  16   Kwak, H., An, J. & Ahn, Y.-Y. *FrameAxis: Characterizing Framing Bias and Intensity*
367       *with Word Embedding.*  (2020).

368  17  An, J., Kwak, H. & Ahn, Y.-Y. in *Proceedings of the 56th Annual Meeting of the*
369      *Association for Computational Linguistics (Volume 1: Long Papers).* 2450-2461
370      (Association for Computational Linguistics).
371  18  Bodell, M. H., Arvidsson, M. & Magnusson, M. Interpretable Word Embeddings via
372      Informative Priors. *ArXiv* **abs/1909.01459**, (2019).
373  19  Kang, D. & Evans, J. Against Method: Exploding the Boundary Between Qualitative
374      and Quantitative Studies of Science. **Quantitative Science Studies** (2020).
375  20  Mathers, C. D. History of global burden of disease assessment at the World Health
376      Organization. *Arch Public Health* **78**, 77, (2020).
377  21  US National Institutes of Health. *ChemIDplus*,
378      <https://chem.nlm.nih.gov/chemidplus/jsp/chemidheavy/help.jsp> (2020).
379  22  Hansen, L. The Truven Health MarketScan Databases for life sciences researchers.
380      *Truven Health Ananlytics IBM Watson Health*, (2017).
381  23  Liao, W. & Veeramachaneni, S. A Simple Semi-supervised Algorithm For Named
382      Entity Recognition. *Proceedings of the NAACL HLT Workshop on Semi-supervised*
383      *Learning for Natural Language Processing*, (2009).
384  24  Gu, Y. *et al.* Domain-Specific Language Model Pretraining for Biomedical Natural
385      Language Processing. *ArXiV* **abs/2007.15779**, (2020).
386  25  Ju, M., Nguyen, N. T. H., Miwa, M. & Ananiadou, S. An ensemble of neural models for
387      nested adverse drug events and medication extraction with subwords. *J Am Med*
388      *Inform Assn* **27**, 22-30, (2020).
389  26  Wu, Z. *et al. A Comprehensive Survey on Graph Neural Networks.* (2019).
390  27  Lan, Z. *et al. ALBERT: A Lite BERT for Self-supervised Learning of Language*
391      *Representations.* (2019).
392  28  Zipf, G. K. The meaning-frequency relationship of words. *J Gen Psychol* **33**, 251-256,
393      (1945).
394  29  Laherrere, J. & Sornette, D. Stretched exponential distributions in nature and
395      economy: "fat tails" with characteristic scales. *Eur Phys J B* **2**, 525-539, (1998).

396
397  **Figure Legends**
398
399  **Figure 1.** Named Entity Recognition Ontology (NERO). The Ontology is shown here as
400  a multifurcating tree, with taxonomy nodes corresponding to ontology classes. Class
401  name and class mentions count in the corpus are shown in parentheses next to each
402  named entity class. Each taxonomy class is provided with a unique pictogram (black
403  and red shapes on yellow background) intended to simplify expert manual annotation of
404  the corpora. In total, we annotated 35,865 sentences. These sentences encapsulated
405  190,679 named entities and 43,438 events connecting two or more entities. In addition
406  to the almost two dozen, more sparsely-used branches (such as *ExperimentalFactor*
407  and *GeographicalLocation*) under the *NamedEntity* cluster, there are three heavily-
408  represented branches in our corpus: *AnatomicalPart*, *Chemical*, and *Process*. Slightly
409  more than half (51.6 percent) of all entities are from these three classes, with 26.6
410  percent of all entities originating from *Process* alone. We designed our ontology and its
411  annotations to capture the named entities associated with research activities and
412  facilities; these types of entities can be important for encoding methods used in

413    scientific experiments or patient treatment. The semantic classes *ResearchActivity* and
414    *MedicalProcedures* turn out to be the ninth and the tenth most frequent, respectively.
415    Other top concepts related to research include *Measurement*, *IntellectualProducts*,
416    *PublishedSourceOfInformation*, *Facility*, and *MentalProcess*.
417
418    **Figure 2. The relative abundance of annotated named entity classes in our**
419    **corpus.** As is typically the case with human languages, semantic classes are
420    represented unevenly in free texts, following a heavy-tail (Zipf's) distribution. (A) In
421    biomedical corpora, unsurprisingly, named entities associated with *genes* and *proteins*
422    are the most prevalent (15 percent), followed by *processes* (9 percent), *medical findings*
423    (8.8 percent), and *chemicals* (6.7 percent). At the low-frequency end of the named entity
424    spectrum, we find *journal names*, *units*, *citations*, and *languages*. (B) Events connecting
425    two or more entities are also approximately Zipf-law distributed. Event frequencies are
426    closely tracking corresponding named entity classes. For example, the most frequent
427    event, *bind*, is associated with the most frequently named entity, *GeneOrProtein*. We
428    tried fitting the rank-ordered frequency distribution of annotated named entities with a
429    Discrete Generalized Beta Distribution (DGBD). The result showed a significant
430    deviation from Zipf's law [28]: The observed distribution's tail was not heavy enough to
431    match Zipf's distribution, most likely due to the relatively small number of classes in our
432    ontology. [29] In other words, we expect that frequencies of semantic classes in a very
433    large corpus, annotated with classes from a hypothetical perfect named entity ontology,
434    would follow a Zipfian (discrete Pareto) distribution of named entity classes. Our action
435    annotations have moved beyond interactions between proteins and genes (*e.g.*, *bind*,
436    *inhibit*, *phosphorylate*, *encode*), into interactions involving genetic variants and
437    environmental factors (*e.g.*, *associated with*, *occur in presence of*, *trigger*, *lack*).
438    Ambiguity levels varied broadly across the named entities captured in our corpus. For
439    example, in the class *AnatomicalPart*, almost all (99.3 percent) are annotated at the
440    most specific levels, with the majority of entities belonging to *BodyPart*,
441    *CellularComponent*, and *Cell*. In contrast, the general (most vague) concept, *Chemical*,
442    turns out to be the most annotated within its cluster, although more specific subclasses,
443    such as *Protein*, *NucleicAcid,* and *Drug* are also well represented in the corpus. In the
444    *Process* concept cluster, about a third of all concept instances are annotated at a more
445    general *Process* level, and the rest of them are specific concepts, such as
446    *MedicalProcedure*, *MolecularProcess*, *ResearchActivity*, and *BiologicalProcess*. In
447    addition to these major clusters of concepts, several individual concepts are well
448    represented in the corpus. For example, *MedicalFinding* represents 7.3 percent of all
449    entities. Other well-represented concepts include *Duration*, *IntellectualProduct*,
450    *Measurement*, *Organism*, *PersonGroup*, *PublishedSource OfInformation*, and *Quantity*.
451    In total, about 70.4 percent of all entities are annotated at the most specific ontology
452    level. There are five concepts in the NERO ontology that allow the semantic flexibility
453    needed to avoid arbitrary concept assignment. Entities annotated as
454    *AminaoAcidOrPeptide*, *QuantityOrMeasurement*, *PublicationOrCitation*
455    *MedicalProcedureOrDevice,* and *GeneOrProtein* account for 17.8 percent of all entities,
456    while less than a quarter (23 percent) of entities representing either genes or proteins
457    are cleanly annotated with class *Gene* or class *Protein*. The remainder are annotated
458    with class *GeneOrProtein*. In addition to the action *bind,* actions indicating entities'
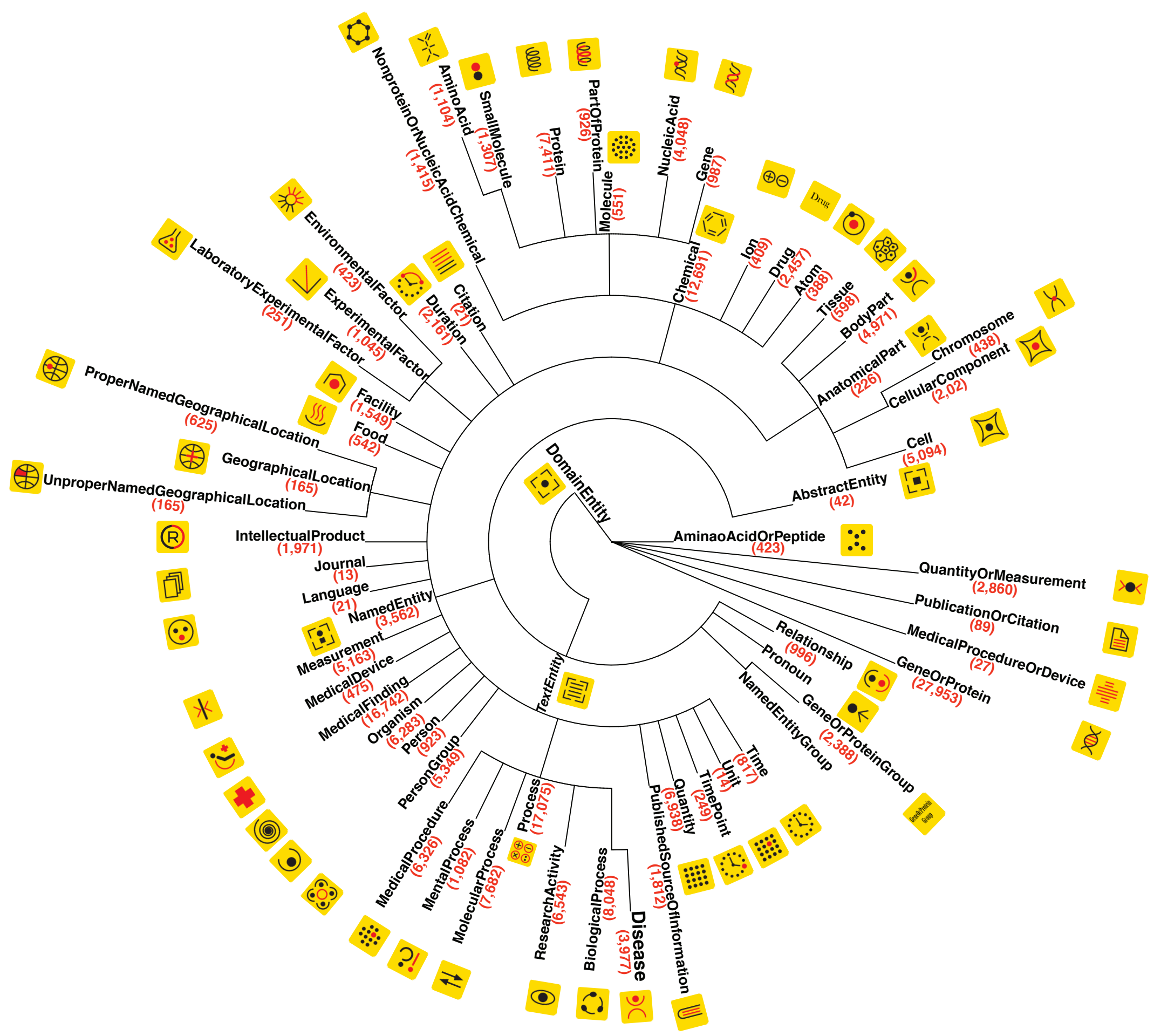
459  attributes are the next most frequent. Other biological relationships are also well-
460  represented in this annotation, such as *inhibit*, *activate*, *mediate*, *interact*, *contain*, and
461  *regulate*. The top 30 action categories account for 64.4 percent of all actions annotated
462  with the top ten action categories accounting for 52.2 percent. Interestingly, negations of
463  actions were also quite abundant in our annotated corpus. For example, *do not bind*
464  was the sixth most frequent normalized action. Other well-represented negations of
465  actions include *do not affect* and *do not inhibit* (see *Supplementary Figure 1*).

467  **Figure 3.** Properties of diseases and drugs visible in the first three principal
468  components of our multi-dimensional text embedding. The figure shows a projection of
469  text embedding into three-dimensional space, with named entities corresponding to
470  diseases and drugs shown with prisms and spheres, respectively. The figure represents
471  several projections of the same embedding, preserving spatial layout and projection,
472  with distinct elements of the embedding indicated by shape color. The central image
473  shows all disease systems and their corresponding medications together. More
474  specifically, the additional projections show: (A) Zollinger-Ellison syndrome and
475  associated medications; (B) cancers and associated therapies; (C) central nervous
476  system diseases and corresponding medications, and; (D) and (E) Viral and bacterial
477  infectious diseases, respectively, together with corresponding antiviral and antibiotic
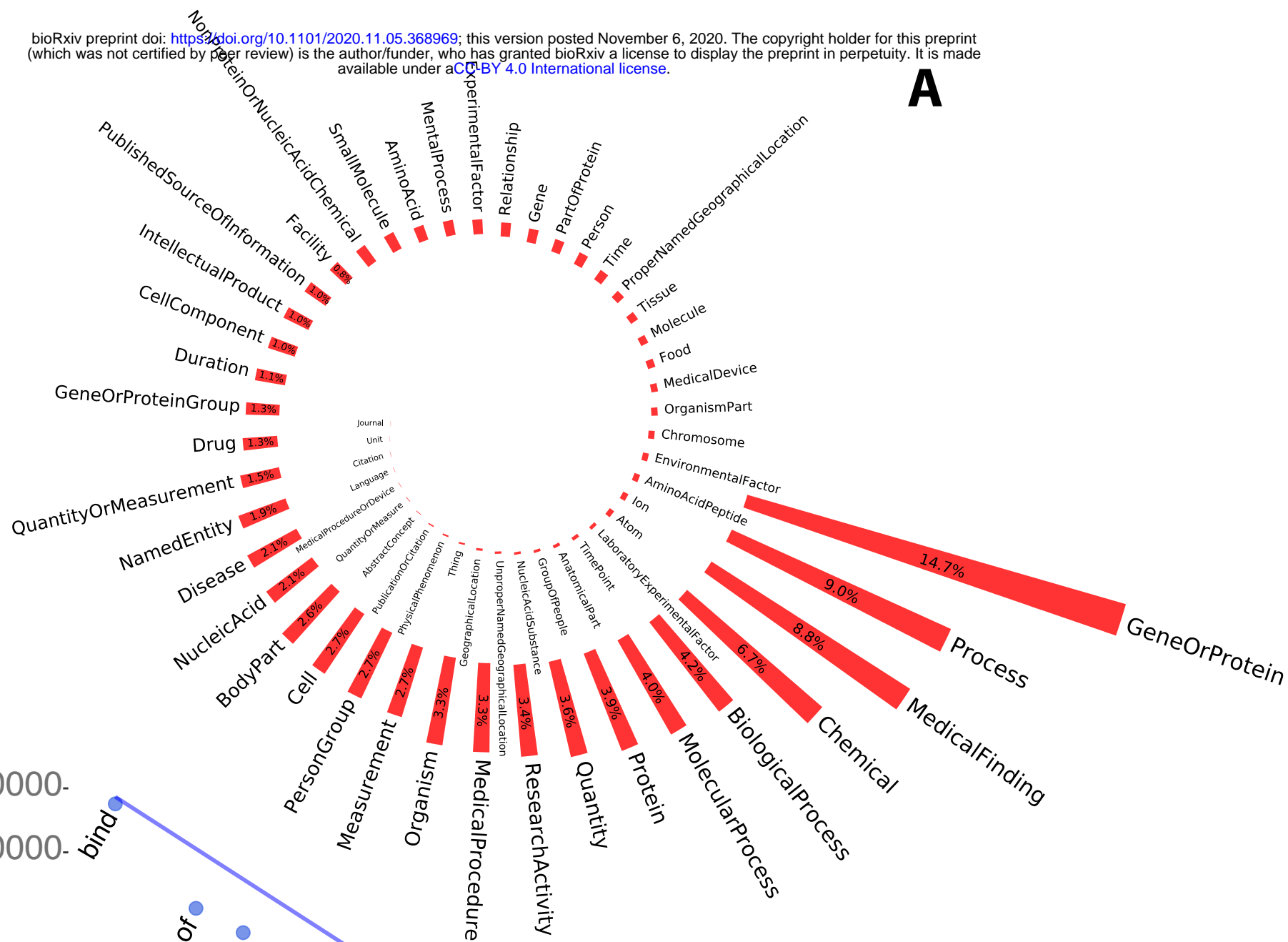478  agents. Another view of the same dataset is presented in Figure 4.

480  **Figure 4. Two-dimensional projections of diseases and medications. (A) We**
481  **projected diseases into two dimensions: female-male (X-axis) and severe-mild (Y-**
482  **axis).** We defined the "male-female" axis using the following pairs of terms: ('male,'
483  'female'), ('prostate,' 'ovary'), ('penile,' 'uterine'), ('penis,' 'uterus'), ('man,' 'woman'),
484  ('men,' 'women'), ('masculine,' 'feminine'), ('he,' 'she'), ('him,' 'her'), ('his,' 'hers'), ('boy,'
485  'girl'), and ('boys,' 'girls'). We defined the severe-mild axis with the following term pairs:
486  ('harmful,' 'beneficial'), ('serious,' 'benign'), ('life-altering,' 'common'), ('disruptive,'
487  'undisruptive'), ('dying,' 'recovering'), ('dangerous,' 'safe'), ('threatening,' 'low-priority'),
488  ('high mortality,' 'harmless'), ('costly,' 'cheap'), ('hospitalized,' 'self-administered'),
489  ('hospital,' 'work'), ('debt,' 'savings'), ('low quality of life,' 'undisruptive'), and ('hazard,'
490  'routine'). **(B)** We projected medications into "benign-toxic" (X-axis) and "cheap-costly"
491  (Y-axis). For the "benign-toxic" axis, we used the following pairs of antonym words:
492  ('harmful,' 'beneficial'), ('toxic,' 'nontoxic'), and ('noxious,' 'benign'). We defined the
493  "expensive–inexpensive" dimension using the following pairs of terms: ('expensive,'
494  'inexpensive'), ('costly,' 'cheap'), ('brand,' 'generic'), and ('patented,' 'off-patent').

496  *Table 1. **Inter-annotator Agreement Statistics.***

| Agreement Type | IAA (%) |
| --- | --- |
| Exact Match | 86.49 |
| Relaxed Match | 93.66 |
| Exact Match | 86.56 |
| Parent Match | 87.66 |
| Superclass Match | 86.72 |
| Ambiguity Match | 97.58 |

498

NonproteinOrNucleicAcidChemical (1,415)
AminoAcid (1,104)
SmallMolecule (1,307)
Protein (7,411)
PartOfProtein (926)
NucleicAcid (4,048)
Gene (987)
Molecule (551)
Chemical (12,691)
Ion (409)
Drug (2,457)
Atom (388)
Tissue (598)
BodyPart (4,971)
AnatomicalPart (226)
Chromosome (438)
CellularComponent (2,02)
Cell (5,094)
AbstractEntity (42)

EnvironmentalFactor (423)
LaboratoryExperimentalFactor (251)
ExperimentalFactor (1,045)
Citation (2)
Duration (2,161)

Facility (1,549)
Food (542)
ProperNamedGeographicalLocation (625)
GeographicalLocation (165)
UnproperNamedGeographicalLocation (165)

DomainEntity
AminaoAcidOrPeptide (423)
QuantityOrMeasurement (2,860)
PublicationOrCitation (89)
MedicalProcedureOrDevice (27)
GeneOrProtein (27,953)
Relationship (996)
Pronoun
GeneOrProteinGroup (2,388)
NamedEntityGroup

IntellectualProduct (1,971)
Journal (13)
Language (21)
NamedEntity (3,562)
Measurement (5,163)
MedicalDevice (475)
MedicalFinding (16,742)
Organism (6,283)
Person (923)
PersonGroup (5,349)
Process (17,075)
MedicalProcedure (6,326)
MentalProcess (1,082)
MolecularProcess (7,682)
ResearchActivity (6,543)
BiologicalProcess (8,048)
Disease (3,977)
PublishedSourceOfInformation (1,812)
Quantity (6,938)
TimePoint
Unit (249)
Time (817)
TextEntity

**A**



**B**

A

Zollinger-Ellison syndrome

ranitidine
rabeprazole
omeprazole

B

anthracyclines
folic acid antagonist

breast cancer

ovarian cancer
gastric malignancy
hepatoma
myeloproliferative disorders
cutaneous T-cell lymphomas

ifosfamide
hydroxyurea
epirubicin
doxorubicin

C

torticollis
cerebral palsy

Huntington's
vertigo

riluzole
propoxyphene        ergotamine
                    lamotrigine
analgesics
                    sumatriptan

ketamine

F

Drug
Disease
CNS/Psychatric
Digestive
Infectious/Immune
Neoplastic
Other

E

cryptococcosis
                        tularemia
blastomycosis           histoplasmosis
candidiasis
aspergillosis
coccidioidomycosis

Amphotericin B

D

poliomyelitis
influenza

oseltamivir

vaccines