
Biomedical Information Extraction for Disease Gene Prioritization

Jupinder Parmar^{*†}
Stanford University
jsparkmar@stanford.edu

William Koehler[†]
OccamzRazor
william@occamzrazor.com

Martin Bringmann
OccamzRazor
martin@occamzrazor.com

Katharina Sophia Volz
OccamzRazor
volz@occamzrazor.com

Berk Kapicioglu[†]
OccamzRazor
berk@occamzrazor.com

Abstract

We introduce a biomedical information extraction (IE) pipeline that extracts biological relationships from text and demonstrate that its components, such as named entity recognition (NER) and relation extraction (RE), outperform state-of-the-art in BioNLP. We apply it to tens of millions of PubMed abstracts to extract protein-protein interactions (PPIs) and augment these extractions to a biomedical knowledge graph that already contains PPIs extracted from STRING, the leading structured PPI database. We show that, despite already containing PPIs from an established structured source, augmenting our own IE-based extractions to the graph allows us to predict novel disease-gene associations with a 20% relative increase in hit@30, an important step towards developing drug targets for uncured diseases.

1 Introduction

Understanding diseases and developing curative therapies requires extracting and synthesizing relevant knowledge from vast swaths of biomedical information. However, with the exponential growth of scientific publications over the past several decades [1], it has become increasingly difficult for researchers to keep up with them. Moreover, most biomedical information is only disseminated via unstructured text, which is not amenable to most computational methods [2]. Thus, there is a growing need for scalable methods that can both extract relevant knowledge from unstructured text and synthesize it to infer novel biomedical discoveries.

To fill this need, we build an end-to-end biomedical IE pipeline [2, 3, 4] by leveraging SciSpacy [5], the most modern and actively developed open-source BioNLP library, and customizing its NER and RE components via transfer learning and BioBERT [6, 7]. We demonstrate that our pipeline outperforms the existing state-of-the-art (SOTA) for biomedical IE, such as PubTator Central [8], its RE extensions [9], and SciSpacy [5] itself.

We then run our pipeline on the PubMed [10] corpus, the largest repository of biomedical abstracts, and extract protein-protein interactions (PPI). Even though our pipeline can easily be trained to extract any relationship, we focus on PPIs because our understanding of them is only partially complete [11, 12, 13], they play an important role in identifying novel disease-gene associations [14], and there is already an established structured PPI database called STRING [15] that allows us to benchmark our extractions.

^{*}Work conducted while author was an intern at OccamzRazor.

[†]Equal contribution.

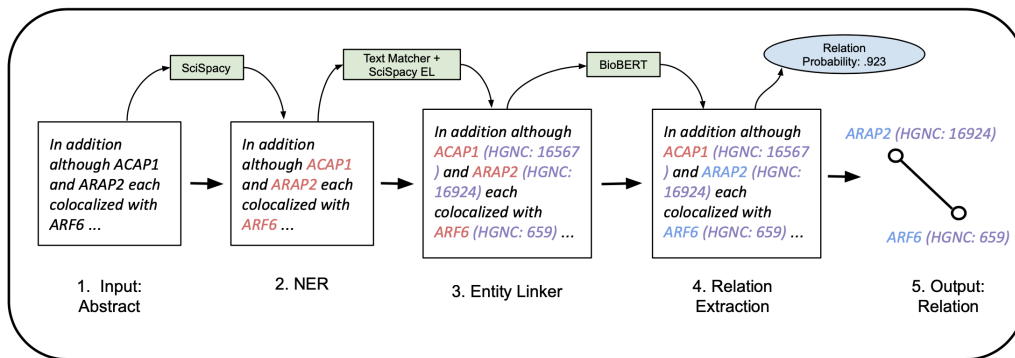


Figure 1: A high-level overview of our IE pipeline. We only display the single candidate relation (ARAP2, ARF6) for simplicity although three candidate relations are present.

Finally, we augment our IE-based PPIs to a knowledge graph that already contains STRING-based PPIs and demonstrate that the augmentation yields a 20% relative increase in hit@30 for predicting novel disease-gene associations. Even though biomedical IE pipelines have previously been evaluated in downstream link prediction tasks when the IE-based extractions were the sole source of the graph [16, 17], to the best of our knowledge, we are the first to show a lift in a setting where the knowledge graph is already populated by an established structured database that contains the same relation type.

Increasing predictive accuracy in such a difficult setting demonstrates the quality of our biomedical IE pipeline, which is specifically designed to require only a small amount of training data to extract any biomedical relationship, and moves us one step closer towards developing drug targets for uncured diseases.

2 Biomedical Information Extraction

In Figure 1, we provide an overview of our biomedical IE pipeline that we train and evaluate on PPI data annotated by in-house biologists. In the following subsections we review how we configured the pipeline for biomedical text and show how each component outperforms its leading competitor in BioNLP.

2.1 Named Entity Recognition (NER)

We train our NER model using SpaCy [18], which we customize further via ScispaCy’s [5] word vectors pre-trained on biomedical text. Our training dataset consists of ~2000 PubMed abstracts tagged with proteins. We enforce strict annotation rules during the labeling process to help disambiguate unclear protein references, a task that we found other NER datasets do not accomplish effectively given the complex nature of biomedical literature. We then compare our model’s performance on the test set against two of the leading biomedical NER systems: PubTator Central [8], a web service that performs NER on PubMed abstracts, and ScispaCy [5], which provides its own protein NER model. As seen in Table 1, our model outperforms both of them.

System	Precision	Recall	F1
Our Model	78.41	73.87	76.08
PubTator	58.96	49.20	45.76
ScispaCy	37.81	57.96	53.64

Table 1: NER Test Results

2.2 Relation Extraction (RE)

For training and evaluating our RE model, we automatically annotate a separate set of ~2000 PubMed abstracts using our NER model, generate relation candidates between pairs of tagged proteins, and manually annotate whether a given candidate contains an interaction. Using our NER model for annotation ensures that our RE model is trained and evaluated based on the same data distribution it handles in production.

We then develop and evaluate a variety of RE models. First, we create models based on feature engineering that use GloVe embeddings [19] and various linguistic features known to perform well on BioNLP tasks [20]. Then, we develop models based on BERT [7], BioBERT [6], and SciBert [21]. We represent the task of relation extraction in these models using the entity start, mention pool, and masked input configurations discussed in [6, 22]. For BERT-based models, we experiment both with fine-tuning and feature extraction. In our feature extraction experiments we combine BERT-based features with our own engineered features.

System	Precision	Recall	F1
v1	43.24	45.71	44.44
v2	41.17	50.00	45.16
v3	31.37	68.57	43.04
Masked BioBERT	29.87	70.00	41.88

Table 2: RE Test Results.

SOTA model in terms of the F1 score. Since all of the models perform well on a different metric, we decide to run each of them on the entire PubMed corpus.

We compare each of our proposed configurations against the SOTA for biological RE [6], a masked input BioBERT model. We refer to our top three models as v1: BioBert feature extraction and feature engineering, v2: Fine-tuned SciBERT using mention pooling, and v3: Fine-tuned BioBERT using entity start. Table 2 reports the evaluation results for our top three models and the SOTA model. We note that each of our models outperforms the

3 Extracting Relations from PubMed

We run each of our pipeline configurations on PubMed [10], a repository of over 30 million biomedical abstracts that we filter down to 10 million based on their relevance to humans or mice.

After extracting PPIs from PubMed, we compare them to the ones in STRING [15], the leading structured PPI database, and ascertain to what extent our IE-based extractions are novel and in fact a segment of the siloed biomedical knowledge contained only in text. The results of the comparison are shown in Figure 2. We observe that IE-based PPIs do not significantly overlap with those in STRING as the highest proportion of extracted relations contained in STRING among the three pipelines is v1 at 24.32%. Additionally, we observe that each configuration behaves as we expect. Specifically, pipeline v3, whose relation extraction model has the highest recall, extracts the most relationships, whereas pipeline v1, whose relation extraction model has the highest precision, extracts the least relationships.

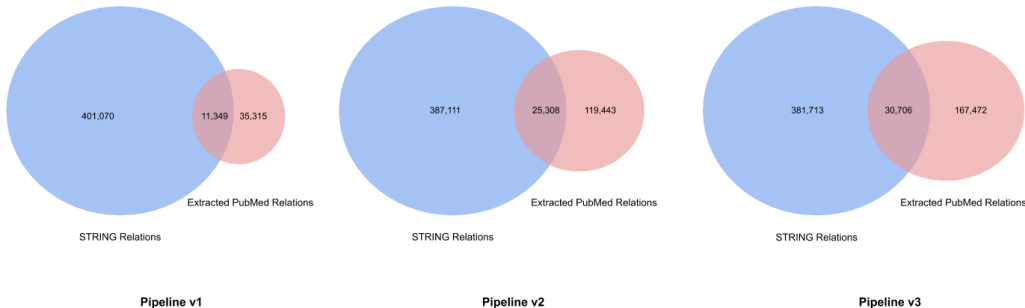


Figure 2: A comparison of different IE pipelines and STRING.

Finally, our pipeline extracts more PPIs than previous biomedical information extraction attempts. Most notably, Percha and Altman[9] extend PubTator [8] with RE functionality by using a dependency parser and clustering-based algorithms. They extract 41,418 PPIs, whereas each of our pipelines extract substantially more. In addition, we observe that the 198,178 PPIs pipeline v3 extracts is more in line with the biomedical expectation since researchers determined there to be roughly 650k PPIs in the human body of which only around 200k have been validated experimentally [11, 23].

4 Disease Gene Prioritization

The reason we developed our biomedical information extraction pipeline is to extract biomedical knowledge from unstructured text, construct a biomedical knowledge graph, and leverage this graph to infer novel biomedical discoveries. In previous sections we demonstrated that the components of our biomedical IE pipeline outperforms the leading NER and RE models in BioNLP. In this section, we demonstrate that our biomedical IE pipeline goes further and also enables novel biomedical discoveries.

Specifically, we focus on the problem of identifying disease genes, a set of genes associated with a particular disease. We formulate this task as a link prediction problem [24, 25] where we construct a biomedical knowledge graph and leverage the information in the graph to predict previously unknown links between genes and diseases. Identifying said links then helps in developing drug targets for uncured diseases.

Historically, biomedical IE pipelines have been evaluated in downstream link prediction tasks when the IE-based extractions were the sole source of the graph [16, 17]. In this paper, we attempt to ascertain whether a biomedical IE pipeline can also be used to complement an established structured database that provides edges of the same relation type.

To demonstrate this, we construct five different biomedical knowledge graphs. For evaluation, we use DisGeNET [26], the leading database for gene-disease associations. We split DisGeNET edges randomly into train (80%), valid (10%), and test sets (10%), and use the same valid and test sets for evaluating all five graphs. The only difference between the five graphs is the training data. The first graph only uses the train set of DisGeNET. The second graph augments the train set of DisGeNET with STRING. The remaining graphs augment the second graph, namely DisGeNET and STRING, with extractions from one of the three versions of our biomedical IE pipeline.

For each experiment, we train and evaluate a link prediction model using a graph embedding algorithm called RotatE [27] and use a library called Optuna [28] for hyper-parameter optimization. The results of the experiments are shown in Table 3. Note that MR is the mean of all gene-disease link ranks, MP is the mean of the rank divided by the pool for that disease, and hit@k describes the percentage of links we obtain in the top "k" ranks.

	MR	MP	hit@30	hit@3	hit@1
IE v3 + STRING + DisGeNET	1418.397	92.484	37.367%	15.302%	7.829%
IE v2 + STRING + DisGeNET	1441.802	92.262	35.409%	14.057%	7.473%
IE v1 + STRING + DisGeNET	1829.548	89.869	32.74%	13.701%	6.762%
STRING + DisGeNET	1952.084	89.362	31.139%	13.879%	7.651%
DisGeNET	7422.117	59.544	0.356%	0.178%	0.178%

Table 3: Link prediction results on various biomedical knowledge graphs.

We observe that augmenting v3 of our IE extractions to the graph provided a lift across all metrics compared to the strong baseline of both STRING and DisGeNET. Specifically, MR had a relative reduction of 27.3%, hit@3 had a relative lift of 10.3%, and hit@30 had a relative lift of 20.0%.

This indicates that the large amount of relations extracted from PubMed contains high-quality edges and can be immediately helpful to a number of biomedical tasks. Additionally, by achieving better performance in disease gene identification when augmenting a knowledge graph that already contained PPIs from a structured resource with our extracted relations, we illustrate the tremendous representational power contained in our IE-based PPIs.

5 Conclusion

We have introduced a biomedical IE pipeline that can be configured to extract any biomedical relationship from unstructured text using a small amount of training data. We empirically demonstrated that its NER and RE components outperform their leading competitors such as PubTator Central [8], its RE extension [9], scispaCy [5], and BioBERT [6]. We then ran it on tens of millions of PubMed abstracts to extract hundreds of thousands of PPIs and show that these relations are novel

in comparison to the ones in leading structured databases. Finally, we evaluated our IE-based PPIs' ability to enable biomedical discoveries by augmenting them to a knowledge graph that already contains STRING-based PPIs and showed that the augmentation yielded a 20% relative increase in hit@30 for predicting novel disease-gene associations. We believe that increasing predictive accuracy in such a difficult setting demonstrates the quality of our biomedical IE pipeline, which we plan to use to uncover other biological relationships currently locked away in biomedical texts, and moves us one step closer to developing drug targets for uncured diseases.

References

- [1] Esther Landhuis. Scientific literature: Information overload. *Nature*, 2016.
- [2] Chung Chi Huang and Zhiyong Lu. Community challenges in biomedical text mining over 10 years: Success, failure and the future. *Briefings in Bioinformatics*, 2016.
- [3] Graciela H. Gonzalez, Tasnia Tahsin, Britton C. Goodale, Anna C. Greene, and Casey S. Greene. Recent advances and emerging applications in text and data mining for biomedical discovery. *Briefings in Bioinformatics*, 2016.
- [4] Mina Gachloo, Yuxing Wang, and Jingbo Xia. A review of drug knowledge discovery using BioNLP and tensor or matrix decomposition. *Genomics and Informatics*, 2019.
- [5] Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. ScispaCy: Fast and Robust Models for Biomedical Natural Language Processing. In *BioNLP Workshop and Shared Task*, 2019.
- [6] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: A pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 2020.
- [7] Jacob Devlin, Ming Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019.
- [8] Chih Hsuan Wei, Alexis Allot, Robert Leaman, and Zhiyong Lu. PubTator central: automated concept annotation for biomedical full text articles. *Nucleic Acids Research*, 2019.
- [9] Bethany Percha and Russ B. Altman. A global network of biomedical relationships derived from text. *Bioinformatics*, 2018.
- [10] Kathi Canese and Sarah Weis. PubMed: The bibliographic database. *The NCBI Handbook*, 2013.
- [11] Michael P.H. Stumpf, Thomas Thorne, Eric De Silva, Ronald Stewart, Jun An Hyeong, Michael Lappe, and Carsten Wiuf. Estimating the size of the human interactome. *Proceedings of the National Academy of Sciences (PNAS)*, 2008.
- [12] Kavitha Venkatesan, Jean François Rual, Alexei Vazquez, Ulrich Stelzl, Irma Lemmens, Tomoko Hirozane-Kishikawa, Tong Hao, Martina Zenkner, Xiaofeng Xin, Kwang Il Goh, Muhammed A. Yildirim, Nicolas Simonis, Kathrin Heinzmann, Fana Gebreab, Julie M. Sahalie, Sebiha Cevik, Christophe Simon, Anne Sophie de Smet, Elizabeth Dann, Alex Smolyar, Arunachalam Vinayagam, Haiyuan Yu, David Szeto, Heather Borick, Amélie Dricot, Niels Klitgord, Ryan R. Murray, Chenwei Lin, Maciej Lalowski, Jan Timm, Kirstin Rau, Charles Boone, Pascal Braun, Michael E. Cusick, Frederick P. Roth, David E. Hill, Jan Tavernier, Erich E. Wanker, Albert László Barabási, and Marc Vidal. An empirical framework for binary interactome mapping. *Nature Methods*, 2009.
- [13] Marc Vidal, Michael E. Cusick, and Albert László Barabási. Interactome networks and human disease. *Cell*, 2011.
- [14] Yves Moreau and Léon Charles Tranchevent. Computational tools for prioritizing candidate genes: Boosting disease gene discovery. *Nature Reviews Genetics*, 2012.
- [15] Damian Szklarczyk, Annika L. Gable, David Lyon, Alexander Junge, Stefan Wyder, Jaime Huerta-Cepas, Milan Simonovic, Nadezhda T. Doncheva, John H. Morris, Peer Bork, Lars J. Jensen, and Christian Von Mering. STRING v11: Protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Research*, 2019.
- [16] Julien Fauqueur, Ashok Thillaisundaram, and Theodosia Togia. Constructing large scale biomedical knowledge bases from scratch with rapid annotation of interpretable patterns. In *BioNLP Workshop and Shared Task*, 2019.
- [17] Meena Nagarajan, Angela D. Wilkins, Benjamin J. Bachman, Ilya B. Novikov, Shenghua Bao, Peter J. Haas, María E. Terrón-Díaz, Sumit Bhatia, Anbu K. Adikesavan, Jacques J. Labrie, Sam Regenbogen, Christie M. Buchovecky, Curtis R. Pickering, Linda Kato, Andreas M. Lisewski, Ana Lelescu, Houyin Zhang, Stephen Boyer, Griff Weber, Ying Chen, Lawrence Donehower, Scott Spangler, and Olivier

- Lichtarge. Predicting future scientific discoveries based on a networked analysis of the past literature. In *International Conference on Knowledge Discovery and Data Mining (KDD)*, 2015.
- [18] Matthew Honnibal and Ines Montani. spaCy2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. *To appear*, 2017.
- [19] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2014.
- [20] Jake Lever and Steven Jones. Painless Relation Extraction with Kindred. In *BioNLP Workshop and Shared Task*, 2017.
- [21] Iz Beltagy, Kyle Lo, and Arman Cohan. SCIBERT: A pretrained language model for scientific text. In *Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2020.
- [22] Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. Matching the blanks: Distributional similarity for relation learning. In *Association for Computational Linguistics (ACL)*, 2020.
- [23] Esti Yeger-Lotem and Roded Sharan. Human protein interaction networks across tissues and diseases. *Frontiers in Genetics*, 2015.
- [24] David Liben-Nowell and Jon Kleinberg. The link-prediction problem for social networks. *Journal of the American Society for Information Science and Technology*, 2007.
- [25] Víctor Martínez, Fernando Berzal, and Juan Carlos Cubero. A survey of link prediction in complex networks. *ACM Computing Surveys*, 2016.
- [26] Janet Piñero, Juan Manuel Ramírez-Anguita, Josep Saüch-Pitarch, Francesco Ronzano, Emilio Centeno, Ferran Sanz, and Laura I. Furlong. The DisGeNET knowledge platform for disease genomics: 2019 update. *Nucleic Acids Research*, 2020.
- [27] Zhiqing Sun, Zhi Hong Deng, Jian Yun Nie, and Jian Tang. RotatE: Knowledge graph embedding by relational rotation in complex space. In *International Conference on Learning Representations, (ICLR)*, 2019.
- [28] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A Next-generation Hyperparameter Optimization Framework. In *International Conference on Knowledge Discovery and Data Mining (KDD)*, 2019.