

Abstract

Here we describe the application of BioBERT, a biomedical text-mining model to automate the extraction of information from public, scientific data sources in the construction of an internal knowledgebase of Myxoma virus research findings. This construct, in combination with other manually-curated data sets from the scientific literature and open access databases, will provide researchers with a highly explorable, graph-based database for new insight discovery about Myxoma virus biology. This research demonstrates the effectiveness of using available programming languages and deep learning techniques to rapidly construct and populate purpose-built biomedical research knowledgebases to enable new insight discoveries.

Introduction

BERT: Natural language processing (NLP) and text mining approaches have been used previously to identify valuable information from published literature. Bidirectional Encoder Representations from Transformers (BERT) showed powerful performance in a wide range of NLP research. Unlike other language models, BERT is designed to train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all learning layers[1]. Users of BERT can use the model which has been trained from immense resource by developers and fine-tune their model based on their target, rather than use their own hardware and time to train the model.

BioBERT: The performance of BERT leaves room for improvement because the word distribution of general corpuses of text, which prior NLP algorithms have been trained on, are not well suited to biomedical literature, creating challenges for the application of generalized text mining models in biomedical knowledge engineering uses. To address this, BioBERT was developed as a model specifically trained for biomedical text (Figure 1)[2].

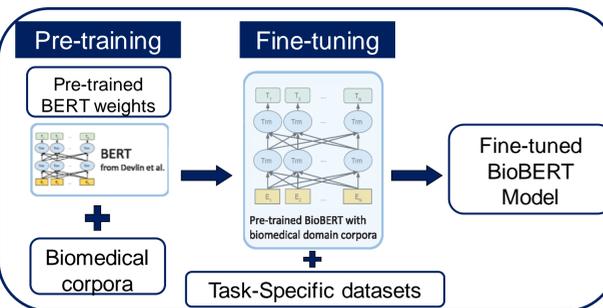


Figure 1. Overview of the pre-training and fine-tuning of BioBERT.

BERN: BERN is a neural Biomedical named Entity Recognition and multi-type Normalization tool [3]. BERN uses BioBERT named entity recognition(NER) models and recognize known entities and discover new entities. It provides Web-based service for researchers who want to discover and tag entities in raw text and even PubMed articles with PubMed ID.

Myxoma virus: We decided to focus on scientific knowledge on the Myxoma virus, due to the fact that the data would be of a manageable size and the types of knowledge, the data would present enough variability for us to determine the accuracy of BioBERT application.

Methods and Materials

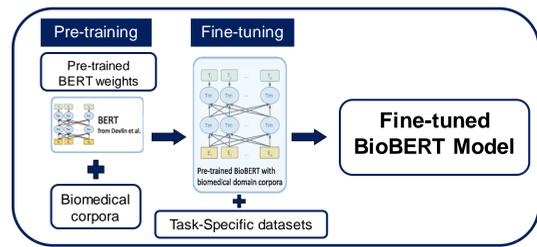
THE BIOMEDICAL TEXT-MINING PIPELINE

1) Pre-training of model

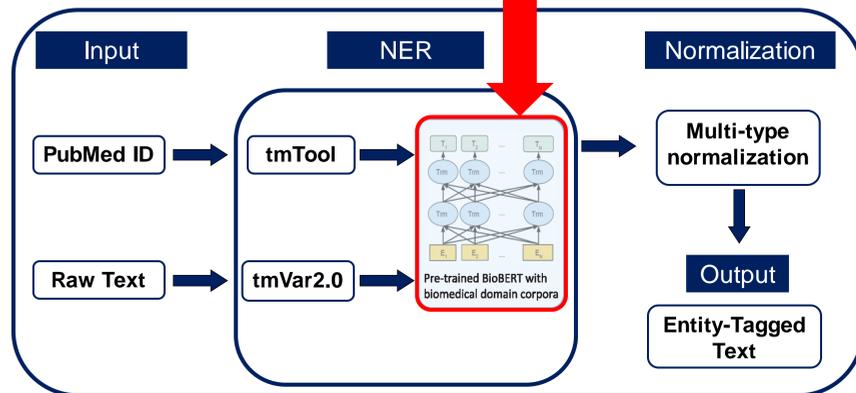
Using pre-trained weights of BERT as published in the BERT repository by Google which was pretrained on English Wikipedia and Books Corpus[4].

2) Fine-tuning of model

Using Speices-800 datasets to train and test model[5]. This dataset is focus on species names which can be used in NER.



3) BERN



MATERIALS

NVIDIA DGX-1 Workstation

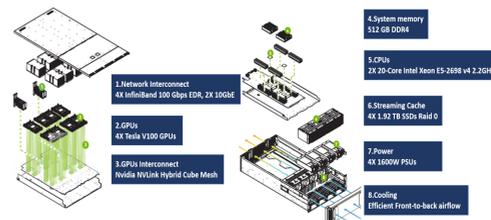


Figure 2. DGX-1 with V100 system components.

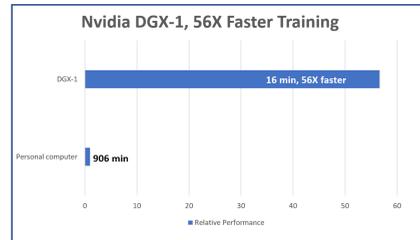


Figure 3. Relative performance based on pre-training time of BioBERT NER model.

Results

MODEL TRAINING RESULT

Model	Total Tokens	Total Phrases	Phrases found	Phrases correct	Precision	Recall	F-1 score
BioBERT trained model	42298	767	827	592	71.58%	77.18%	74.28
Non-trained model	42298	767	5346	25	0.47%	3.26%	0.82

Figure 4. Performance evaluation result of the trained model and the non-trained model applied to scientific data

Results

NAMED ENTITY RECOGNITION RESULT

PubMed ID	Location	Category	Entity ID
30940649	1	Title	
30940649	1	Abstract	
30940649	25	Entity	Myxoma virus M013 protein gene CUI-less
30940649	6	Entity	Myxoma disease MESH:D009232 BERN:106924701
30940649	12	Entity	Myxoma virus species NCBI:txid10273
30940649	38	Entity	NF-kappaB gene MIM:164011 HGNC:7794 Ensembl:ENSG00000109320 BERN:324174102
30940649	167	Entity	myxoma virus species NCBI:txid10273
30940649	181	Entity	M013 gene CUI-less
30940649	214	Entity	viral pyrin domain-only protein gene CUI-less
30940649	247	Entity	vPOP gene CUI-less
30940649	279	Entity	M013 drug CUI-less
30940649	279	Entity	M013 protein gene CUI-less
30940649	344	Entity	NF-kappaB gene MIM:164011 HGNC:7794 Ensembl:ENSG00000109320 BERN:324174102
30940649	394	Entity	ASC1 gene MIM:604501 HGNC:12310 Ensembl:ENSG00000103671 BERN:325476802
30940649	403	Entity	NF-kappaB1 gene CUI-less
30940649	486	Entity	M013 pyrin domain gene CUI-less
30940649	505	Entity	PYD gene CUI-less
30940649	551	Entity	PYD gene CUI-less
30940649	637	Entity	M013 drug CUI-less
30940649	642	Entity	PYD gene CUI-less
30940649	708	Entity	aspartate drug CHEBI:72314 BERN:283828103
30940649	722	Entity	glutamate drug CHEBI:14321 BERN:295829803
30940649	786	Entity	aspartate drug CHEBI:72314 BERN:283828103
30940649	800	Entity	glutamate drug CHEBI:14321 BERN:295829803
30940649	858	Entity	ASC-1 gene MIM:604501 HGNC:12310 Ensembl:ENSG00000103671 BERN:325476802
30940649	926	Entity	M013 drug CUI-less
30940649	931	Entity	ASC-1 gene MIM:604501 HGNC:12310 Ensembl:ENSG00000103671 BERN:325476802
30940649	995	Entity	caspace-1 gene MIM:147678 HGNC:1499 Ensembl:ENSG00000103752 BERN:325264602
30940649	1055	Entity	M013 gene CUI-less
30940649	1089	Entity	C CHEBI:30039 BERN:305421203
30940649	1122	Entity	N drug CHEBI:29555 BERN:308520903
30940649	1133	Entity	PYD gene CUI-less
30940649	1213	Entity	M013 gene CUI-less
30940649	1267	Entity	M013 gene CUI-less
30940649	1276	Entity	NF-kappaB1 gene CUI-less
30940649	1427	Entity	vPOP M013 gene CUI-less
30940649	1432	Entity	M013 drug CUI-less
30940649	1511	Entity	NF-kappaB gene MIM:164011 HGNC:7794 Ensembl:ENSG00000109320 BERN:324174102

Figure 5. Result of NER mission by using BioBERT and BERN. There are four different entity categories, gene/protein, drug/chemical, disease and species. The targets are myxoma virus, related genes/proteins and drug/chemicals. In the title it only shows the article is related to M013 protein, NF-kappa B and inflammasome pathways. BioBERT and BERN discovered different genes which is related to myxoma virus in the abstract.

Conclusion

Text mining is a powerful and rapidly developing technology that can be useful for helping humans find relevant information in very large datasets. Text mining models pre-trained on biomedical corpora have been shown to be promising. This study provides solid evidence that text-mining could be effectively used on a specific subset of biomedical research to rapidly provide knowledgebases for new insight discovery. For virus research and research in other areas where the knowledge may be less well known or disseminated, text-mining can play a critical role in bringing together disparate sources into a single research resource. By combining appropriate tools, we are enabling the discovery of useful information more efficiently than current search methods.

References

- Devlin J, Chang M, Lee K, Toutanova K et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Minneapolis, MN, USA, pp. 4171-4186. Association for Computational Linguistics. <https://www.aclweb.org/anthology/N19-1423>.
- Lee J, Yoon W, Kim S, Kim D, Kim S, So C et al. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*. 2019.
- Kim D, Lee J, So C, Jeon H, Jeong M, Choi Y et al. A Neural Named Entity Recognition and Multi-Type Normalization Tool for Biomedical Text Mining. *IEEE Access*. 2019;7:73729-73740.
- naver/bioBERT-pretrained [Internet]. GitHub. 2020 [cited 3 April 2020]. Available from: <https://github.com/naver/bioBERT-pretrained/releases>
- HCMR S. SPECIES - Organism Name Identification in the Scientific Literature [Internet]. *Species.jensenlab.org*. 2020 [cited 3 April 2020]. Available from: <https://species.jensenlab.org/>.

Acknowledgements

This research was supported by Systems Imagination. I would like to acknowledge Chris Yoo for his mentorship and guidance throughout the project. I am also grateful to Pieter Derdeyn for his expertise and comments on the earlier stage of research.