

---

# ELECTRAMed: A NEW PRE-TRAINED LANGUAGE REPRESENTATION MODEL FOR BIOMEDICAL NLP

---

Giacomo Miolo\*, Giulio Mantoan\*, Carlotta Orsenigo

Department of Management, Economics and Industrial Engineering  
Politecnico di Milano, Italy

{giacomo.miolo,giulio.mantoan,orsenigo}@mip.polimi.it

## ABSTRACT

The overwhelming amount of biomedical scientific texts calls for the development of effective language models able to tackle a wide range of biomedical natural language processing (NLP) tasks. The most recent dominant approaches are domain-specific models, initialized with general-domain textual data and then trained on a variety of scientific corpora. However, it has been observed that for specialized domains in which large corpora exist, training a model from scratch with just in-domain knowledge may yield better results. Moreover, the increasing focus on the compute costs for pre-training recently led to the design of more efficient architectures, such as ELECTRA. In this paper, we propose a pre-trained domain-specific language model, called ELECTRAMed, suited for the biomedical field. The novel approach inherits the learning framework of the general-domain ELECTRA architecture, as well as its computational advantages. Experiments performed on benchmark datasets for several biomedical NLP tasks support the usefulness of ELECTRAMed, which sets the novel state-of-the-art result on the BC5CDR corpus for named entity recognition, and provides the best outcome in 2 over the 5 runs of the 7th BioASQ-factoid Challenge for the question answering task.

**Keywords** Pre-trained language models · ELECTRA · Biomedical NLP

## 1 Introduction

The immense body of biomedical scientific texts, which steadily grows at an exponential rate, makes it imperative to develop effective machine learning methods able to automatically extract the rich knowledge therein contained, that can be used to address several biomedical natural language processing (NLP) tasks.

Prominent NLP advancements in recent years have been mostly driven by the use of deep neural models, which require large corpora of annotated training data. Compared to the general domain, however, the collection of such data in the biomedical field is difficult and expensive, since it necessarily involves domain experts for accurate data labelling. For this reason, semi-supervised pre-trained language models, such as ELMo [1] and BERT [2], were developed and successfully applied in a wide range of NLP tasks.

ELMo and BERT leverage contextualized word embeddings for which the representation of a word depends on the context where it is used and, therefore, is a function of the entire input sequence. ELMo exploits a deep bidirectional language model pre-trained on a large corpus of texts. It was proposed to address both the syntactic and semantic complexities and ambiguities of words and has been proven to achieve notable results in a variety of NLP problems [1]. BERT [2] resorts to the transformer architecture to pre-train bidirectional language representations, instead of relying on recurrent neural networks. By embracing attention mechanisms, BERT showed to distinguish the sense of words at a very fine level and to grasp many of their syntactic and semantic properties. This ability made BERT the state-of-the-art contextualized word representation model for the most challenging natural language understanding problems.

---

\*Corresponding authors.

While ELMo and BERT architectures pre-trained on general-domain corpora are well-established top performers for general NLP tasks, they might yield poor results in case of scientific or specific domains, since the corpora used for pre-training, such as news articles and Wikipedia [3], might not include the same terminology adopted in the in-domain tasks. For specialized contexts past studies showed that general-domain language models can largely benefit from the use of in-domain textual data [4]. As a consequence, recent models for biomedical NLP relied on adapted versions of general-domain approaches. Among these, two of the most noteworthy and successful examples are represented by BioBERT [5] and BlueBERT [4], which are domain-specific language models initialized with the general-domain BERT, and then pre-trained on a wide range of biomedical and scientific corpora. In principle, these last methods rely on the assumption that initializing the pre-training from general-domain models might improve the overall performance for domain-specific purposes. However, it has been recently observed that for domains in which large corpora exist, like the biomedical field, pre-training language models from scratch yields better results than feeding the pre-training phase with general-domain knowledge [6].

To obtain contextualized word embeddings, BERT pre-training is based on masked language modelling (MLM) which aims at predicting a small random subset of masked input tokens, considering only the token context. This approach allows the model to learn bidirectional representations. A different input corruption procedure has been recently proposed, in which instead of masking, and therefore losing, some of the input tokens, these are replaced with plausible alternatives produced by a small generator network. By learning from all the input tokens this novel approach, called ELECTRA (Efficiently Learning an Encoder that Classifies Token Replacements Accurately), is computationally much more efficient than BERT, and has been shown to outperform the latter in several tasks [7].

Inspired by previous research achievements which showed the effectiveness of language models built on domain-specific knowledge, and recognizing the importance of resorting to efficiently pre-trained methods while preserving downstream performances, in this study we describe a new ELECTRA-based model, called ELECTRAMed, suited for the biomedical domain.

In particular, the main contributions of the present work are as follows:

- We propose a novel ELECTRA-based language representation model (ELECTRAMed) pre-trained on biomedical corpora. To the best of our knowledge, this is the first ELECTRA-based model specifically developed for the biomedical domain, and the present study is the first which applies, to this domain, a transformer architecture different from BERT.
- We tested ELECTRAMed on several biomedical benchmark NLP tasks. The results achieved empirically show the effectiveness of the proposed approach, which performed at par, and sometimes better, than state-of-the-art models while leveraging the reduced computational effort required by ELECTRA-based architectures.
- We make publicly available the pre-processed datasets used in our study as well as the pre-trained weights of ELECTRAMed and the source code for fine-tuning the model<sup>1</sup>.

## 2 Materials and methods

As it is generally the case with pre-trained methods for language representation, the development of the new model encompassed two main phases represented, respectively, by pre-training and fine-tuning, as illustrated in the following sections. Since the proposed approach shares the same architecture of ELECTRA, a brief description of the latter is also required.

### 2.1 ELECTRA

ELECTRA (Efficiently Learning an Encoder that Classifies Token Replacements Accurately) is a pre-trained language model recently proposed to overcome the computational drawback of approaches based on masked language modelling (MLM), like BERT, where the input is first corrupted by masking some tokens in the sentence and a network is then trained to recover the original identity of the corrupted tokens. Since the network learns from a small percentage of masked-out tokens (typically 15%), these techniques usually require a huge amount of computations to be effective [7].

To address this inefficiency, ELECTRA implements an alternative pre-training procedure called Replaced Token Detection (RTD), in which a network is trained to distinguish real input tokens from synthetic but plausible replacements. Specifically, the model consists of two neural networks which are pre-trained jointly. The first is a generator that performs MLM by providing tokens substitutes, and that then learns to predict the original token from the masked form.

---

<sup>1</sup><https://github.com/gmpoli/electramed>

The second is a discriminator which is trained to detect the synthetic tokens, i.e. to distinguish real tokens from those replaced by the generator. After pre-training the generator is dropped and the discriminator is fine-tuned on labeled data for downstream tasks.

Compared to MLM, the innovative pre-training procedure embedded in ELECTRA is more computationally efficient since the discriminator is required to predict the class, real vs. replaced, of each of the input tokens, thereby learning from the entire input sequence instead of a small portion of it. This feature lets ELECTRA compete with state-of-the-arts models in several NLP problems while reducing dramatically the computational effort needed for pre-training, as described in the seminal paper [7]. Moreover, by resorting to plausible tokens alternatives for masking, RTD helps in mitigating a disadvantage of the traditional MLM approach, which introduces a potential mismatch in the learning framework between pre-training and fine-tuning due to the use, for input corruption, of the same conventional [MASK] token which does not appear in the fine-tuning phase. Its computational efficiency, together with the promising results achieved for different NLP tasks, motivated the adoption of ELECTRA in the present study.

## 2.2 ELECTRAMed pre-training

Differently from other prominent pre-trained models, such as BioBERT, ELECTRAMed wasn't initialized with the weights from ELECTRA. Indeed, pre-training was performed entirely on a biomedical domain corpus. The corpus at hand was published by [4] and consists of 28,714,373 PubMed abstracts (uncompressed, approximately 26GB), representing the whole amount of abstracts published on the free digital repository until September 2018. In more detail, the corpus contains  $\sim 181\text{M}$  sentences and  $\sim 4\text{B}$  words and was subject to the following common NLP pre-processing steps applied by the proponents of the dataset:

- Text lowercasing;
- Special characters ( $\backslash x00\text{-}\backslash x7F$ ) removal;
- Text tokenization using NLTK Treebank tokenizer <sup>2</sup>.

In any NLP task, vocabularies are needed to encode tokens with numbers, and are generally built so to contain the most frequent words or subword units. In the present study we made use of SciVocab, a WordPiece vocabulary proposed by [8] and built on a scientific text corpus by means of the SentencePiece library <sup>3</sup>. Compared to the more commonly used vocabulary released with BERT, SciVocab is characterized by almost the same size ( $\sim 30\text{k}$ ) and 42% of tokens overlap, thereby showing a substantial difference in the words used frequently in scientific texts with respect to the general domain case. Notice that, the use of a more scientific-oriented vocabulary should reduce the incidence of out-of-vocabulary (OOV) tokens and, therefore, the loss of information in the text encoding phase. The hyperparameters applied for pre-training ELECTRAMed on the corpus mentioned above are reported in Table 1. These correspond to the same set of parameters used for the ELECTRA-base model, as described in [7]. For tokenization we instead resorted to the WordPiece scheme adopted by BERT.

It is well known that generating models which adhere to the pre-training and fine-tuning paradigm is a highly resource-intensive process. In the case of BERT, for example, the computational complexity of the self-attention layer increases quadratically with the length of the sequences after tokenization. As a result, limiting the maximum length of the input sequence to 128, at least for the first part of the training, and then increasing it to 512 at a later stage, is a commonly adopted strategy to gain efficiency. This is done despite the risk of reducing the ability of the model to capture long-distance dependencies and, therefore, to negatively affect its performance. Differently from other approaches, by leveraging the computational advantages provided by the ELECTRA framework, we were able to investigate the use of a maximum sequence length of 512, instead of 128, for the whole pre-training phase, while keeping the training time on par with other methods. Indeed, pre-training ELECTRAMed took  $\sim 10$  days by using one TPU v3 with 8 cores.

## 2.3 ELECTRAMed fine-tuning

After pre-training, ELECTRAMed was fine-tuned and tested on three biomedical NLP tasks, represented by named entity recognition (NER), relationship extraction (RE) and question answering (QA).

Named entity recognition (NER) is aimed at automatically finding and tagging in a text meaningful terms, called named entities. In a general domain these typically refer to sequences of words corresponding to specific entities in the real world, such as locations, persons, organizations. In the biomedical field named entities can represent the

---

<sup>2</sup><http://www.nltk.org/>

<sup>3</sup><https://github.com/google/sentencepiece>

Table 1: Hyperparameters used for ELECTRAMed pre-training

Hyperparameter	Value
Number of layers	12
Hidden size	768
FFN inner hidden size	3072
Attention heads	12
Attention head size	64
Embedding Size	768
Generator Size	1/3
Mask percent	15
Learning Rate Decay	Linear
Warmup steps	10000
Learning Rate	2e-4
Adam $\epsilon$	1e-6
Adam $\beta_1$	0.9
Adam $\beta_2$	0.999
Attention Dropout	0.1
Dropout	0.1
Weight Decay	0.01
Batch size	256
Train epochs	1M

name of genes, diseases, chemical compounds, and drugs, to name a few. The goal of NER is tagging each token in a sentence with one taken from a list of possible named entities, or with the "no entity" label. NER is often applied as a preliminary step for many applications, such as relationship extraction and knowledge base completion. It is worthwhile to notice that, compared to the general domain, NER in the biomedical field is considered to be more complex, since biomedical entities constantly grow in number with the scientific progress, may contain special characters and can be referred to using a wide variety of synonyms and abbreviations.

Among the available annotation schemes used to label multi-token named entities, in the present work we resorted to BIO tagging [9], for which three binary classifiers are trained to label each token in the text as *B* (the token is the beginning of a named entity), *I* (the token is inside a named entity but is not the beginning), and *O* (otherwise).

For fine-tuning and testing ELECTRAMed on biomedical NER, three publicly available corpora were used. The first, denoted here as NCBI-disease, is the disease corpus of the National Center for Biotechnology Information (NCBI) [10], which is a collection of 793 PubMed abstracts annotated at the mention and concept level. In detail, it contains 6892 disease mentions mapped into 790 unique disease concepts. The second dataset (BC5CDR) was used for the BioCreative V Chemical Disease Relation task and is composed by 1500 PubMed articles with 4409 annotated chemicals, 5818 diseases and 3116 chemical-disease interactions [11]. From the entire collection we extracted the subset with chemical and disease entities, along the lines of previous studies. The third dataset (JNLPBA) was released for an open challenge and derives from the GENIA corpus, which is a collection of abstracts retrieved by controlled search on MEDLINE. For the shared task, the authors simplified the original 36 classes into five super-classes represented, respectively, by proteins, DNA, RNA, cell types and cell lines [12]. These correspond to the biological named entities of interest.

Relationship extraction (RE) generally follows NER and is aimed at finding semantic relationships which may occur in a text between two or more entities. In the biomedical domain relationships are extracted between biological entities and, for example, can take the form of different kinds of relations among genes and diseases or of protein-protein interactions. RE is cast as a text classification problem in which, given a pair of entities, the goal is to assign the correct type of relationship, whether this exists. To validate ELECTRAMed on RE problems two datasets were analyzed. The first is CHEMPROT [13], which is a corpus containing annotated relationship types between chemicals and proteins. From the initial group of 10 different kinds of relationships, five of these (CPR:3, CPR:4, CPR:5, CPR:6 and CPR:9) were used for evaluation. The second corpus (DDI-2013) contains 792 texts extracted from the DrugBank database and other 233 MEDLINE abstracts describing drug-drug interactions (DDIs). The corpus was annotated by considering four different types of interactions (i.e. effect, mechanism, advice, and int), where the last was used when a DDI appeared in the text without providing any specific additional information [14].

Table 2: Description of the corpora used as benchmarks for biomedical NER

Dataset	Entity Type	N. entities	N. train	N. val	N. test
NCBI-disease	Disease	6,892 dis.	5,429	923	941
BC5CDR	Disease and Chemical	10,227 TOT 4,409 chemicals 5,818 diseases	3,951	3,957	4,145
JNLPBA	Gene and Protein	59,963 TOT 35,336 proteins 10,589 DNA 1,069 RNA 8,639 cell types 4,330 cell lines	16,845	1,743	3,869

Table 3: Description of the corpora used as benchmarks for biomedical RE

Dataset	Relationship	N. relations	N. train	N. val	N. test
CHEMPROT	Chemical-Protein	10,028 TOT 1,983 CPR:3 5,006 CPR:4 484 CPR:5 727 CPR:6 1,828 CPR:9	19,460	11,820	16,943
DDI-2013	Drug-Drug	2,937 TOT 1,212 effect 946 mechanism 633 advice 146 int	18,779	7,244	5,761

Question answering (QA) can be considered an extension of information retrieval and has the purpose of delivering direct and precise responses to questions asked in natural language. Focusing on the biomedical domain, and in particular on inquiries about the treatment of a disease, a QA tool would provide, as answers, specific drugs which are effective against that disease or short text passages containing the response. Indeed, there are three categories of questions which are typically asked to a QA system: confirmation questions, that can be dealt with by yes-or-no statements, list questions, which expect a list of entities or facts, and factoid questions, which require a single short phrase or sentence as response. In our work, we focused on this last type of questions, for which the goal is finding the correct answer inside a passage. This is achieved by training a model capable of predicting the starting and the ending tokens of the text segment containing the full expected response. For biomedical semantic QA we selected as benchmark the pre-processed version of the BioASQ Task 7b-factoid dataset [15], which contains both biomedical questions and gold standard answers, in the form of relevant concepts, articles, snippets, exact answers, summaries and the corresponding full PubMed abstracts as passages.

As for NER, question answering in the biological domain poses a major challenge in contrast to the open or other restricted domains, for the presence of a highly specialized terminology and for a potential larger gap in technicality between the questions made by non-expert users and the target documents. Another critical issue is the scarcity of labelled datasets. To tackle this problem, prior to addressing the biomedical QA task, we fine-tuned ELECTRAMed on the Stanford Question Answering Dataset (SQuAD v1.1), a large-scale general-domain reading comprehension dataset published by [16] and including 87,599 and 10,570 examples for training and testing, respectively. This additional fine-tuning was performed in line with previous studies, which showed its effectiveness in improving the performance on domain-specific activities [17], [15].

The description of the corpora used for the three NLP tasks in terms of number of training, testing and validation examples, and number of observations for each class, are provided in Tables 2, 3 and 4. The hyperparameters applied for fine-tuning ELECTRAMed on all the tasks are instead indicated in Table 5. Notice that, the datasets at hand are gold standards for the evaluation of biomedical NLP tools. To enable a fair comparison with previous studies, therefore, we adopted the same partition into training, test and validation sets as originally proposed by the authors

Table 4: Description of the corpus (BioASQ 7b-factoid) used as benchmark for biomedical QA

Batch	N. questions	N. question-context pairs
Train	556	5537
Test batch 1	39	98
Test batch 2	25	56
Test batch 3	29	84
Test batch 4	34	90
Test batch 5	35	79

Table 5: Hyperparameters used for ELECTRAMed fine-tuning

Hyperparameter	NER	RE	QA_SQuAD	QA_BioASQ
Learning Rate	5e-5	5e-5	3e-5	5e-6
Adam $\epsilon$	1e-6	1e-6	1e-6	1e-6
Adam $\beta_1$	0.9	0.9	0.9	0.9
Adam $\beta_2$	0.999	0.999	0.999	0.999
Layerwise LR Decay	0.8	0.8	0.8	0.8
Learning Rate Decay	Linear	Linear	Linear	Linear
Attention Dropout	0.1	0.1	0.1	0.1
Dropout	0.1	0.1	0.1	0.1
Weight Decay	0	0	0	0
Batch size	32	32	16	16
Max Sequence Length	128	128	384	384
Document Stride	NA	NA	128	128

Table 6: Detected named entities and answers provided by ELECTRAMed (in bold) for samples taken from NER and QA corpora

Task	Dataset	Sample
NER	BC5CDR	Thus, <b>indomethacin</b> by inhibition of <b>prostaglandin</b> synthesis may diminish the blood pressure maintaining effect of the stimulated effect of the stimulated renin- <b>angiotensin</b> system in <b>sodium</b> and volume depletion.
NER	NCBI-disease	Occasional missense mutations in ATM were also found in <b>tumour</b> DNA from patients with <b>B-cell non-Hodgkins lymphomas (B-NHL)</b> and a <b>B-NHL</b> cell line..
QA	BioASQ-7b	Q: What is the cause of a STAG3 truncating variant? A: STAG3 truncating variant as the cause of <b>primary ovarian insufficiency</b> . <b>Primary ovarian insufficiency</b> (POI) is a distressing cause of infertility in young women ...
QA	BioASQ-7b	Q: Which receptor is targeted by Erenumab? A: ... Erenumab, a human monoclonal antibody that inhibits the <b>calcitonin gene-related peptide receptor</b> , is being evaluated for migraine prevention ...

of the corpora. Finally, by way of example, Table 6 shows the detected named entities and the answers provided by ELECTRAMed for samples extracted from the corpora used for NER and QA tasks.

### 3 Results

For named entity recognition the performance of ELECTRAMed was evaluated by means of the F<sub>1</sub>-score (F), defined as the harmonic mean between precision (P) and recall (R). Accordingly to [9], precision was computed as the percent-



Table 7: Precision (P), recall (R) and F<sub>1</sub>-score (F) for the NER task

Benchmark	Metrics	SOTA1	SOTA2	SOTA3	ELECTRAMed
NCBI-disease	P	88.22	?	?	85.87
	R	91.25	?	?	89.29
	F	<b>89.71</b>	89.13	88.85	87.54
BC5CDR	P	92.05	?	?	88.76
	R	87.91	?	?	91.34
	F	89.93	89.73	89.42	<b>90.03</b>
JNLPBA	P	?	72.24	?	69.33
	R	?	83.26	?	78.56
	F	<b>81.29</b>	77.59	77.03	73.65

State-of-the-art (SOTA) performance. For NCBI-disease, SOTA1 is BioBERT, [5], SOTA2 is Spark NLP, [18], SOTA3 is BioFLAIR, [19]. For BC5CDR, SOTA1 is RL+DS+PA, [20], SOTA2 is Spark NLP and SOTA3 is BioFLAIR. For JNLPBA, SOTA1 is Spark NLP, SOTA2 is BioBERT and SOTA3 is BioFLAIR.

age of named entities found by the model that were correct (i.e. were an exact match of the corresponding entity in the corpus), whereas recall was set as the percentage of entities included in the corpus and detected by the model. For the purpose of relationship extraction, the F<sub>1</sub>-score was applied as well. In this case, the predicted triplets entity-relation-entity were deemed correct if the relation and the two entities were the same as the ground truth. Finally, for question answering we resorted to three quality measures, commonly used for assessing tools which generate a list of possible responses to a given inquiry. The first is the mean reciprocal rank (MRR), given by the average of the reciprocal ranks of the results for a sample of queries. The other two are the strict accuracy (SACC) and the lenient accuracy (LACC), for which a question is correctly answered if the gold response is the first element of the list returned by the model, or it is included in the list, respectively.

For each of the NLP problems ELECTRAMed was compared with the current best state-of-the-art (SOTA) models, including those developed by the participants of the 7th BioASQ Challenge runs (QA activity). The models comprised within each SOTA group are specified in the captions of Tables 7, 8 and 9, for every task and corpus. These tables also contain missing values denoted as “?”, to indicate the unavailability of the corresponding measure for the given pair model-dataset. Missing outcomes are associated to precision and recall for both NER and RE tasks. For the sake of completeness, we deemed relevant to provide the results for ELECTRAMed also in terms of these two quality measures. It is also worthwhile to observe that, all the values referred to ELECTRAMed in the tables are obtained by averaging the respective metrics over five runs with different seeds. This is in line with the experimental settings adopted for SciBERT [8], for which the outcomes were computed as average over multiple runs. For the other SOTA models, instead, we were unable to establish whether the results reported by the authors were averaged or corresponded to the best outcomes achieved.

The results obtained by ELECTRAMed for the task of named entity recognition are shown in Table 7. The proposed model reached the highest F<sub>1</sub>-score (90.03) on the BC5CDR dataset. By performing better than the current SOTA approaches, ELECTRAMed sets a novel state-of-the-art performance on this corpus for NER purposes in terms of F<sub>1</sub>-score. For NCBI-disease and JNLPBA datasets, ELECTRAMed was not among the top three SOTA models, but still reached comparable results on the first corpus.

The outcomes for the task of relationship extraction are indicated in Table 8. These results clearly demonstrate the superior effectiveness of SciBERT over the other models, but also show the promising performance of ELECTRAMed. Indeed, on the DDI-2013 dataset the proposed approach provided results which are close to the second-best model (SOTA2) and ranked at the third position by reaching a fairly higher F<sub>1</sub>-score compared to the current SOTA3 (79.13 vs. 72.90).

Finally, the results achieved for the question answering task are depicted in Table 9. ELECTRAMed was able to outperform all the competitors in two (batches 1 and 4) out of the five runs, and to provide comparable results with the best approach (KU-DMIS-5) for the second batch. For the remaining runs (batches 3 and 5), ELECTRAMed ranked at the seventh and sixth position, respectively, among all the participants. To investigate the ability of ELECTRAMed of providing high-quality responses across the whole challenge, besides the BioASQ baseline method we selected the models that competed in all the runs and, for each of them, we computed a score given by the ratio between their MRR and the highest MRR reached in a given batch. For each model, the sum of the scores over all the runs can be seen as a measure of its ability of providing responses which are close, if not equal, to the best ones in terms of MRR. The results of this analysis are reported in Table 10 and support the effectiveness of ELECTRAMed for the QA

Table 8: Precision (P), recall (R) and F<sub>1</sub>-score (F) for the RE task

Benchmark	Metrics	SOTA1	SOTA2	SOTA3	ELECTRAMed
CHEMPROT	P	?	77.02	?	75.47
	R	?	75.90	?	70.67
	F	<b>83.64</b>	76.46	74.40	72.94
DDI-2013	P	?	?	74.10	80.07
	R	?	?	71.80	78.24
	F	<b>84.08</b>	79.90	72.90	79.13

State-of-the-art (SOTA) performance. For CHEMPROT, SOTA1 is SciBERT, [8], SOTA2 is BioBERT, SOTA3 is BlueBERT, [4]. For DDI-2013, SOTA1 is DESC+MOL+SciBERT, [21], SOTA2 is BlueBERT, SOTA3 is Hierarchy Bi-LSTMs +Att.+SDP, [22].

task at hand. The proposed model, indeed, achieved the highest total score being associated with ratios more densely distributed around 1. In particular, ELECTRAMed performed better than systems based on the BioBERT approach (KU-DMIS-5 model).

Table 9: Strict accuracy (SACC), lenient accuracy (LACC) and mean reciprocal rank (MRR) for the QA task

Batch	Competitor	SACC	LACC	MRR
1	(1) ELECTRAMed	44.62	51.28	<b>47.95</b>
	(2) KU-DMIS-1	41.03	53.85	46.37
	(3) BJUTNLPGroup	30.77	41.03	34.83
	(4) auth-qa-1	25.64	30.77	27.78
2	(1) KU-DMIS-5	52.00	64.00	<b>56.67</b>
	(2) ELECTRAMed	46.40	62.40	53.16
	(3) QA1	36.00	48.00	40.33
	(4) transfer-learning	24.00	44.00	32.67
3	(1) QA1	44.83	58.62	<b>51.15</b>
	(2) UNCC_QA_1	44.83	58.62	51.15
	(3)	41.38	65.52	50.23
	google-gold-input (7) ELECTRAMed	37.93	58.62	46.62
4	(1) ELECTRAMed	61.18	82.35	<b>69.55</b>
	(2) KU-DMIS-1	58.82	82.35	69.12
	(3) FACTOIDS	52.94	73.53	61.03
	(4) UNCC_QA3	52.94	73.53	61.03
5	(1) KU-DMIS-5	28.57	51.43	<b>36.38</b>
	(2) BJUTNLPGroup	28.57	40.00	33.81
	(3) UNCC_QA_1	28.57	42.86	33.05
	(6) ELECTRAMed	24.57	44.00	31.42

Note: The number in round brackets beside each model indicates the ranking in the challenge run.

Table 10: Scores over the five runs of the 7h BioASQ-factoid Challenge

Competitor	Batch1	Batch2	Batch3	Batch4	Batch5	Total
BioASQ Baseline	0.323	0.241	0.258	0.364	0.238	1.424
auth-qa-1	0.579	0.541	0.669	0.536	0.412	2.737
KU-DMIS-1	0.967	0.771	0.924	0.994	0.886	4.542
LabZhu,FDU	0.120	0.441	0.775	0.699	0.615	2.650
ELECTRAMed	<b>1.000</b>	0.938	0.911	<b>1.000</b>	0.864	<b>4.713</b>

Note: Three models by Lab Zhu at Fudan University were proposed. The table includes the one that performed better across the runs.



## 4 Conclusions & Future Developments

The recent literature on biomedical NLP has been heavily influenced by BERT-based architectures and by the usage of domain-specific corpora for pre-training. Meanwhile, in the general-domain NLP literature, a plethora of transformer-based architectures have flourished, bringing significant improvements to the first wide-spread implementation. In this study we presented ELECTRAMed, a new pre-trained language model for biomedical NLP. ELECTRAMed was pre-trained from scratch on a biomedical corpus using a domain-specific vocabulary, and was shown to obtain valuable results for some of the most commonly addressed NLP tasks arising in the biomedical field. For ELECTRAMed pre-training we leveraged the efficiency of the ELECTRA architecture, and we were able to perform at par of the current state-of-the-art models while keeping the computational effort low, in terms of both time and cost.

The results achieved in the present work encourage future studies that can be undertaken along different directions. From one side, it would be worthwhile to investigate the performance of ELECTRAMed when the maximum sequence length of the input is reduced from 512 to 128 for the entire, or at least, the most part of the training phase, with the aim of further reducing the computational requirements. From the other, it would be useful to explore the impact on the model performance of using a vocabulary built upon the selected pre-training corpus. With this study we hope to spark a new wave of transformer-based architectures in the biomedical domain. Consequently, as a future research line it would be also interesting to investigate potential improvements for tasks related to biomedical information extraction by combining existing biomedical domain knowledge resources (e.g. knowledge bases) into novel transformer-based learning frameworks.

## Acknowledgements

We thank Dr. Sandra Coecke from the Joint Research Center at European Commission and Dr. Anna Beronius from Karolinska Institute for their valuable and fruitful discussions that fostered a positive and encouraging environment which greatly contributed to the development of our work.

## References

- [1] M.E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237. Association for Computational Linguistics, New Orleans, Louisiana, 2018.
- [2] J. Devlin, M.W. Chang, K. Lee, and K. Toutanova. 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)*, pages 4171–4186. Association for Computational Linguistics, Minneapolis, USA, 2019.
- [3] Y. Belinkov and J. Glass. Analysis methods in neural language processing: A survey. *Transactions of the Association for Computational Linguistics*, 7:49–72, 2019.
- [4] Y. Peng, S. Yan, and Z. Lu. Transfer learning in biomedical natural language processing: An evaluation of BERT and ELMo on ten benchmarking datasets. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 58–65. Association for Computational Linguistics, Florence, Italy, 2019.
- [5] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C.H. So, and J. Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36:1234–1240, 2020.
- [6] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, and H. Poon. Domain-specific language model pretraining for biomedical natural language processing. *arXiv preprint: 2007.15779*, 2020.
- [7] K. Clark, M.H. Luong, Q.V. Le, and C.D. Manning. ELECTRA: Pre-training text encoders as discriminators rather than generators. *arXiv preprint: 2003.10555*, 2020.
- [8] I. Beltagy, A. Cohan, and K. Lo. SciBERT: Pretrained contextualized embeddings for scientific text. *arXiv preprint: 1903.10676*, 2019.
- [9] E.F.T.K. Sang and S. Buchholz. Introduction to the CONLL-2000 shared task: Chunking. In *Proceedings of the fourth conference on computational natural language learning and of the second learning language in logic workshop (CONLL/LLL 2000)*, pages 127–132. Lissabon, Portugal, 2000.
- [10] R.I. Doğan, R. Leaman, and Z. Lu. NCBI disease corpus: a resource for disease name recognition and concept normalization. *Journal of Biomedical Informatics*, 47:1–10, 2014.

- [11] J. Li, Y. Sun, R.J. Johnson, D. Sciaky, C.-H. Wei, R. Leaman, A.P. Davis, C.J. Mattingly, T.C. Wieggers, and Z. Lu. BioCreative V CDR task corpus: a resource for chemical disease relation extraction. *Database: The Journal of Biological Databases and Curation*, 2016. 2016 baw068.
- [12] J.D. Kim, T. Ohta, Y. Tsuruoka, Y. Tateisi, and N. Collier. Introduction to the bio-entity recognition task at JNLPBA. In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications (NLPBA/BioNLP)*, pages 73–78. COLING, Geneva, Switzerland, 2004.
- [13] M. Krallinger, O. Rabal, S.A. Akhondi, M. Pérez, J. Santamaría, G. Pérez Rodríguez, G. Tsatsaronis, A. Intxaurreondo, J.A. López, U. Nandal, E.V. Buel, A. Chandrasekhar, M. Rodenburg, A. Lægreid, M.A. Doornenbal, J. Oyarzábal, A. Lourenço, and A. Valencia. Overview of the BioCreative VI chemical-protein interaction track. In *Proceedings of the sixth BioCreative challenge evaluation workshop*, pages 141–146. 2017.
- [14] M. Herrero-Zazo, I. Segura-Bedmar, P. Martínez, and T. Declerck. The DDI corpus: An annotated corpus with pharmacological substances and drug–drug interactions. *Journal of Biomedical Informatics*, 46(5):914–920, 2013.
- [15] W. Yoon, J. Lee, D. Kim, M. Jeong, and J. Kang. Pre-trained language model for biomedical question answering. In *Machine Learning and Knowledge Discovery in Databases*, pages 727–740. Springer International Publishing, 2020.
- [16] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. SQuAD: 100,000+ questions for machine comprehension of text. *arXiv preprint: 1606.05250*, 2016.
- [17] G. Wiese, D. Weissenborn, and M. Neves. Neural question answering at BioASQ 5B. *arXiv preprint: 1706.08568*, 2017.
- [18] V. Kocaman and D. Talby. Biomedical named entity recognition at scale. *arXiv preprint: 2011.06315*, 2020.
- [19] S. Sharma and R. Daniel. BioFLAIR: Pretrained pooled contextualized embeddings for biomedical sequence labeling tasks. *arXiv preprint: 1908.05760*, 2019.
- [20] F. Nooralahzadeh, J.T. Lønning, and L. Øvrelid. Reinforcement-based denoising of distantly supervised NER with partial annotation. In *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, pages 225–233. Association for Computational Linguistics, Hong Kong, China, 2019.
- [21] M. Asada, M. Miwa, and Y. Sasaki. Using drug descriptions and molecular structures for drug–drug interaction extraction from literature. *Bioinformatics*, 2020. btaa907.
- [22] Y. Zhang, W. Zheng, H. Lin, J. Wang, Z. Yang, and M. Dumontier. Drug–drug interaction extraction via hierarchical RNNs on sequence and shortest dependency paths. *Bioinformatics*, 34:828–835, 2018.