

Contextual embedding and model weighting by fusing domain knowledge on Biomedical Question Answering

Yuxuan Lu¹
Beijing University of Technology
Chaoyang Qu, Beijing Shi, China
luyuxuanleo@gmail.com

Jingya Yan¹
Beijing University of Technology
Chaoyang Qu, Beijing Shi, China
yanjy1998@163.com

Zhixuan Qi¹
Beijing University of Technology
Chaoyang Qu, Beijing Shi, China
zhixuanqi@outlook.com

Zhongzheng Ge¹
Beijing University of Technology
Chaoyang Qu, Beijing Shi, China
bubble327670@gmail.com

Yongping Du¹
Beijing University of Technology
Chaoyang Qu, Beijing Shi, China
ypdu@bjut.edu.cn

ABSTRACT

Biomedical Question Answering aims to obtain an answer to the given question from the biomedical domain. Due to its high requirement of biomedical domain knowledge, it is difficult for the model to learn domain knowledge from limited training data. We propose a contextual embedding method that combines open-domain QA model AoA Reader and BioBERT model pre-trained on biomedical domain data. We adopt unsupervised pre-training on large biomedical corpus and supervised fine-tuning on biomedical question answering dataset. Additionally, we adopt an MLP-based model weighting layer to automatically exploit the advantages of two models to provide the correct answer. The public dataset BIOMRC constructed from PubMed corpus is used to evaluate our method. Experimental results show that our model outperforms state-of-the-art system by a large margin.

KEYWORDS

biomedical question answering, contextual embedding, model weighting, domain knowledge.

1 INTRODUCTION

Question answering is a classic task in Natural Language Processing, requiring a model to understand natural languages. Cloze-style question answering problem has been a popular task because it is relatively easier to build cloze-style datasets. The cloze style question aims to select the best candidate answer regarding the specified context and fill in the blank of the question. Multiple cloze-style datasets have been published, such as CNN/Daily Mail [12], Children’s Book Test [13], etc. Models based on neural networks are proposed, such as AS READER [14], CAS READER [5], AoA Reader [4] and BERT [6].

These models have achieved good performance on several datasets. However, they do not perform well when facing domain-oriented problems. The main reason is that domain-oriented questions require more background knowledge to give an answer, and a large dataset is needed to allow the models to learn the required domain knowledge.

We make improvements to the existing model, AoA Reader, and validate our results on the BIOMRC dataset [20], the public biomedical dataset constructed from corpus from PubMed. We put forward the Contextual Word Embedding method and the MLP-based model

weighting strategy for the biomedical question answering task. By combining the open-domain QA model and domain-oriented contextual word embedding, the proposed method outperforms state-of-the-art system on biomedical domain question answering significantly, setting up a new state-of-the-art system.

The main contributions of this paper are listed as follows:

- Combining BioBERT and AoA Reader, which can take full advantage of contextual word embedding model pre-trained on large domain corpus and mining semantic and contextual information to choose the best answer. In particular, multiple aggregation methods are adopted and evaluated.
- An MLP-based model weighting strategy is proposed, which can automatically learn the preferences and biases of different models and exploit the advantages of both models to provide the correct answer.
- Our method is evaluated on the BIOMRC dataset, and the results show that it outperforms state-of-the-art system significantly. Our code is available at <https://github.com/leoleoas/MLP-based-weighting>.

2 RELATED WORK

The research of Question Answering has made rapid progress, which benefits from the publication of large-scale and high-quality datasets. Richardson et al. release MCTest [23], a multiple choice machine reading comprehension dataset that opened up research on statistical-based machine learning models. Hermann et al. release the CNN Daily Mail dataset [12], which includes over 1 million cloze-style data. More high-quality datasets have been released since then, such as SQuAD [21, 22], Facebook Children’s Book Test [13], etc.

Models based on deep learning technologies significantly outperform traditional models in extracting context information. Hermann et al. [12] propose an attention-based neural network and proves that the incorporation of attention mechanism is more effective than traditional, statistical-based baselines. Seo et al. propose the BiDAF [24] model, which uses different levels of encoding for linguistic representation and uses a bidirectional attention flow mechanism to obtain the query-aware context representation. Kadlec et al. [14] propose a simple model, the *Attention Sum Reader* (AS Reader),

which uses attention to directly pick answer. Fu et al. propose Ea-Reader [9], whose memory updating rule is able to maintain the understanding of document through read, update and write operations. Chen et al. propose McR² [3], which enables relational reasoning on candidates based on fusion representations of document, query and candidates. Fu et al. propose ATNet [8], which utilities both intra-attention and inter-attention to answer close-style questions over documents. Cui et al. propose *Attention-over-Attention Reader* (AoA READER) [4], which puts another level of document-to-query attention on top of query-to-document attention, achieving state-of-the-art performance on multiple datasets.

In recent years, researchers have focused on combining the unsupervised pre-training on large corpus and supervised fine-tuning on the specific task. Vaswani et al. propose the Transformer [27] model, which uses attention mechanism to replace the CNN and RNN parts of the traditional model to improve the model while speeding up the training process. Devlin et al. build a large-scale unsupervised pre-training model BERT [6] on top of Transformer to pre-train the language representation model using the masked language model task and the next sentence prediction task, and innovatively propose a training strategy that separates pre-training and fine-tuning. While sharing the same pre-trained weights, BERT achieves state-of-the-art performance on many different downstream tasks and datasets. Other BERT-based model is proposed, for instance, Liu et al. propose RoBERTa [17], providing several techniques to robustly pre-train language models. Lan et al. propose ALBERT [15], providing two parameter-reduction techniques to lower memory consumption and increase the training speed of BERT.

In the biomedical domain, Tsatsaronis et al. launch the BioASQ [26] challenges. It contains multiple subtasks, including article / snippet retrieval, document classification and question answering. Pappas et al. construct two cloze-style datasets, BIOREAD [19] and BIOMRC [20], and compare the accuracy of experts, non-experts human and different baseline models and neural network models. The results show that the baseline methods fail to correctly answer questions of the BIOMRC dataset while neural MRC models perform well, indicating that the BIOMRC dataset is less noisy and has enough features for the model to learn. Tang et al. initiate the CORD-19 [25] dataset at the beginning of the global COVID-19 pandemic to help researchers in the biomedical field retrieve articles quickly.

Traditional neural network models require large-scale, high-quality supervised training data to obtain better results. However, it is challenging to build a large-scale and high-quality dataset for domain-oriented tasks because it requires domain experts' annotations. Recent researchers prove that combining the pre-training on large corpus and fine-tuning on supervised training data can achieve better performance on domain-oriented tasks. Gururangan et al. [11] proves that a second phase of pre-training in domain (domain-adaptive pre-training) leads to performance gains. Gu et al. propose BLURB (Biomedical Language Understanding & Reasoning Benchmark) [10] to test the biomedical language understanding ability of language models. Lee et al. propose BioBERT [16] and Cohan et al. propose SciBERT [1], both of which learn the corpus representation of papers on a large-scale corpus of biomedical papers or scientific papers and have obtained better results on natural language processing tasks in biomedical domains. We aim to further improve the performance by using BioBERT to obtain biomedical

contextual and semantic information and by using model weighting layers to combine different models.

3 METHOD

We propose a pre-training strategy based on the scientific pre-training model (SciBERT) and open-domain QA model (AoA Reader) to obtain the final answer to the question. In particular, different embedding and weighting strategies are used in the training process. fig. 1 shows the full structure of our model.

3.1 Formal Task Description

This model is aiming at tasks that comprise cloze-style questions, to which answers are closely related to the comprehension of the context documents included in the problems. A set of candidate answers are also provided alongside, and the model is supposed to choose an answer from the candidates. This task can be formalized as a triplet $\langle C, Q, \mathcal{A} \rangle$ that is inclusive of the given context $C = \{w_1, w_2, \dots, w_n\}$ made of words w_i , a query $Q = \{q_1, q_2, \dots, [MASK], \dots, q_m\}$ where the special token [MASK] marks the position where the answers are supposed to be placed, and answer candidates $\mathcal{A} = \{a_1, a_2, \dots, a_o\}$. A function F is expected to be learned by the model to predict the answer \mathcal{A} of question Q based on its comprehension of the proffered context C :

$$\forall a \in \mathcal{A}, P(a|C, Q) = \begin{cases} 1 & a \text{ is the correct answer} \\ 0 & a \text{ is not the correct answer} \end{cases} \quad (1)$$

$$F(C, Q, \mathcal{A}) = \max_{a \in \mathcal{A}} P(a|C, Q) \quad (2)$$

3.2 Training of SciBERT

3.2.1 Training data. We use the SciBERT [1], which has been pre-trained on the Semantic Scholar corpus. The corpus consists of 18% papers from the computer science domain and 82% from the broad biomedical domain. The unsupervised pre-training process using the large-scale corpus allows the model to obtain the semantic information of biomedical texts. Further, in order to let the model adapt to the cloze-style question answering task, BIOMRC dataset is used to fine-tune the model.

3.2.2 Answer extraction. We adopt the answer extraction strategy Pappas et al. used in their BIOMRC dataset [20]. For each context-question pair, we first divide the context into sentences using NLTK [2]. Each sentence is concatenated to the question by [SEP] token, and they are fed to SciBERT respectively. In this way, we obtain the top-level embedding of the candidate entities and the placeholder in the question. The embeddings of each entity in the sentence are connected to the placeholder's embedding and are sent to a multilayer perceptron to obtain the score for the particular entity. If an entity appears multiple times in the paragraph, we choose the maximum value of its score.

3.3 Training of AoA Reader

In order to make AoA Reader [4] achieve better performance on domain-oriented tasks, we adopt different contextualized word embedding and attention aggregation strategies.

3.3.1 Contextualized word embedding. The AoA Reader uses direct word embedding. This approach converts each token in the context

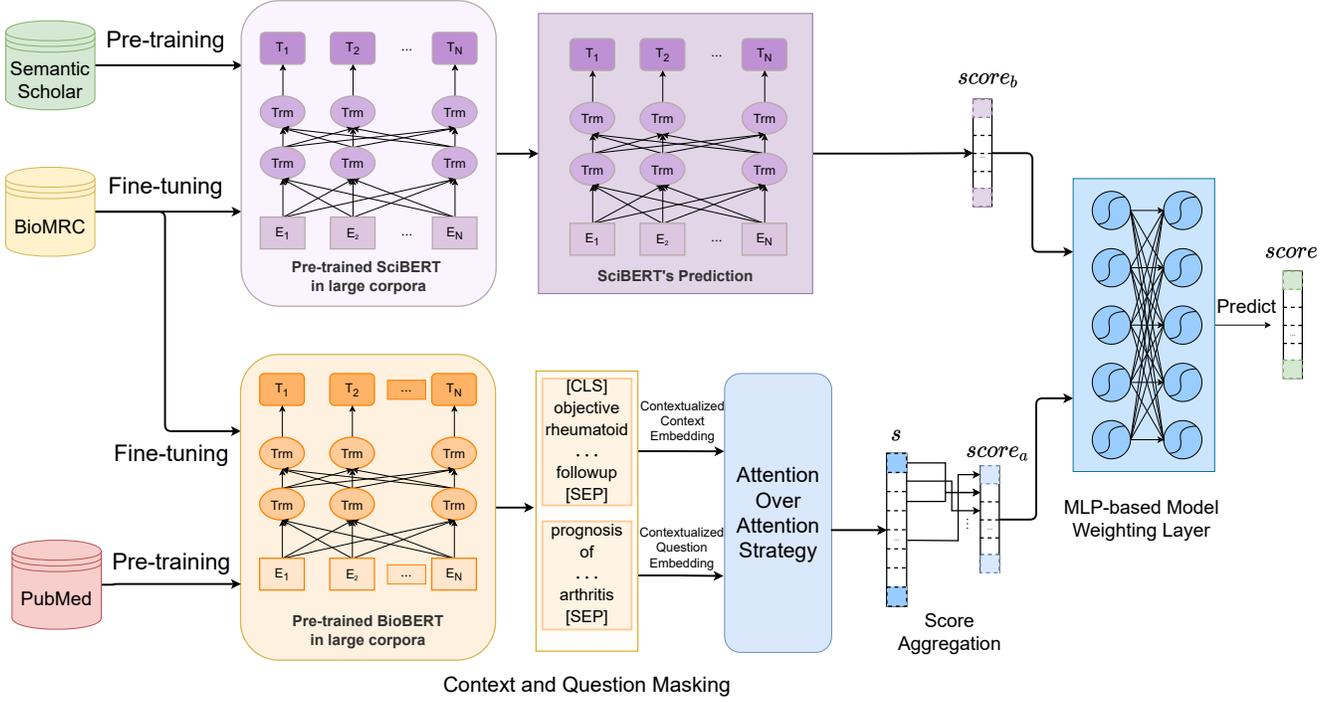


Figure 1: Model structure based on Pre-training and Weighting Strategy

C and the question Q into a one-hot vector. It then converts them into a continuous representation using the shared embedding matrix W_e . Using this method in biomedical question answering may result in severe Out-Of-Vocabulary issues, given that most terms in the biomedical domain are made through several word-formation methods. Meanwhile, the model cannot learn necessary knowledge and the meanings of terms due to the lack of domain-oriented training data.

Instead, we adopt BioBERT to generate contextualized word embedding, which has been pre-trained on a large biomedical corpus from PubMed containing biomedical literature proposed by Devlin [6]. Thus, we can obtain domain-oriented knowledge and terms.

For tokenization, WordPiece [28] through which new words can be represented by known tokens is used. Further, it can solve the Out-Of-Vocabulary issue and allow the model to better understand domain terms made by word-formation methods.

After connecting the context and the question with [SEP] token, it is fed into BioBERT. If it is longer than length limit (512 tokens), we trim the back of the context and keep the original question unchanged:

$$E(C, Q) = BERT([\text{CLS}]; C; [\text{SEP}]; Q; [\text{SEP}]) \quad (3)$$

$E(C, Q)$ in eq. (3) is the contextual embedding of the context C and the query Q .

To obtain the embeddings of context C and the query Q , we've applied a masking operation on $E(C, Q)$ for segmentation. This would conceal the representation of the other segment by zero

vectors, leaving the desired half acquired:

$$E(C)_i = \begin{cases} E(C, Q)_i & i \text{ is a context token} \\ 0 & i \text{ is a question token} \end{cases} \quad (4)$$

$$E(Q)_i = \begin{cases} E(C, Q)_i & i \text{ is a question token} \\ 0 & i \text{ is a context token.} \end{cases} \quad (5)$$

We adopt bi-directional RNN to further obtain the contextual representations $h_{context}(C) \in \mathbb{R}^{(|C|+|Q|+3)*2d}$ of the context C :

$$\overrightarrow{h_{context}(C)} = \overrightarrow{GRU}(E(C)) \quad (6)$$

$$\overleftarrow{h_{context}(C)} = \overleftarrow{GRU}(E(C)) \quad (7)$$

$$h_{context}(C) = [\overrightarrow{h_{context}(C)}; \overleftarrow{h_{context}(C)}] \quad (8)$$

and similarly we obtain $h_{question}(Q) \in \mathbb{R}^{(|C|+|Q|+3)*2d}$ for the question Q .

3.3.2 Pair-wise Matching Score. After obtaining the contextual embedding of the context $h_{context}$ and the question $h_{question}$, we calculate a pairwise matching matrix, which indicates the relevance between a token in the context and question, by calculating their dot product:

$$M(i, j) = h_{context}(i)^T \cdot h_{question}(j) \quad i \in C, j \in Q \quad (9)$$

3.3.3 Attentions over Attention Mechanism. After getting the pairwise matching matrix M , column-wise softmax is used to get the

Context	Because of reports of anaplastic transformation following irradiation, this study examines the incidence of anaplastic transformation and local control of these lesions. This review of seven @entity1 who had @entity189 of the @entity135 that was treated with irradiation shows local control in 71% of cases. There were no cases of anaplastic transformation. This report adds to the literature two cases of "de-differentiation" to less differentiated @entity957 ; one such case occurred after surgery alone. The literature is reviewed. Overall, anaplastic transformation is reported in 7% of @entity1 who had irradiation. De-differentiation occurs after surgery as well. The rate of local control with irradiation is less than 50%; with surgery it is 85%. It is concluded that surgery should be used if the procedure has acceptable morbidity. Otherwise, irradiation can be used. Failures can be salvaged surgically. "Anaplastic transformation" should not affect treatment approach.
Candidate Entities	@entity1: ['patients'] @entity135: ['head and neck'] @entity957: ['squamous carcinomas'] @entity189: ['verrucous carcinoma']
Question	Radiotherapy in the treatment of XXXX of the @entity135 .
Answer	@entity189: ['verrucous carcinoma']

Figure 2: A example of the BIOMRC dataset.

context-level attention regarding each token in the question:

$$\alpha(t) = \text{softmax}(M(1, t), \dots, M(|C|, t)) \quad (10)$$

$$\alpha = [\alpha(1), \alpha(2), \dots, \alpha(|Q|)] \quad (11)$$

We calculate a reversed attention using row-wise softmax to obtain the "importance" of each token in the question regarding each token in the context:

$$\beta(t) = \text{softmax}(M(t, 1), M(t, 2), \dots, M(t, |Q|)). \quad (12)$$

We average all $\beta(t)$ to get a question-level attention β :

$$\beta = \frac{1}{n} \sum_{t=1}^{|C|} \beta(t) \quad (13)$$

Then we adopt the attention-over-attention mechanism, by merging these two attentions to get the "attended context-level attention":

$$s = \alpha^T \cdot \beta \quad (14)$$

where s denotes the importance of each token in the context.

3.3.4 Final Predictions. AoA Reader regards an entity as a word no matter how many words it has. It uses *sum attention* mechanism proposed by Kadlec et al. [14] to get the confidence score of each candidate entity. In our model which uses WordPiece to obtain contextualized word embeddings, an entity may be either segmented into multiple tokens or composed of multiple words, and each token of the entity may occur multiple times in the context. So the confidence score of each candidate answer a is calculated by aggregating all the occurrences of all its tokens in the context:

$$P(a|C, Q) = F_1 \left(F_2 (s_i) \right)_{t \in \mathcal{T}(a) \quad i \in I(t, C)} \quad (15)$$

where $\mathcal{T}(a)$ is the result of segmenting the candidate answer a using WordPiece; F_1 and F_2 are aggregating functions, which can be either **maximum** or **sum**, and $I(t, C)$ indicates the position that the token t appears in the context C .

3.4 Model Weighting strategy

After completing the training of AoA Reader and SciBERT, a model weighting strategy is used to obtain the final answer by combining the advantages of both models.

Our previous study [7] demonstrates that better performance can be achieved using a dual-model weighting strategy. The weighting process is performed by calculating a weighted average of the answer's confidence score and the similarity of the answer derived from two models. Further, considering that different models perform differently against data with different features, we use a simple MLP with one hidden layer to allow the model weighting layer to *automatically* learn this difference and to take advantage of both models.

$$\text{score} = \text{MLP}([\text{score}_a, \text{score}_b]) \quad (16)$$

Where $\text{score}_a, \text{score}_b \in \mathbb{R}^{|\mathcal{A}|}$ is the confidence score of each model. In this way, the weighting layer would be able to learn the predilection and biases of each candidate model and achieve better performance.

4 EXPERIMENT

4.1 Datasets

We conduct the experiments on the BIOMRC LITE dataset [20] to verify the effect of our method. BIOMRC is a biomedical cloze-style dataset for machine reading comprehension. The contexts in each sample are extracted from PUBTATOR, a repository containing 25 million abstracts and their corresponding titles on PubMed. Biomedical entities in the abstract are extracted to form the candidate entities. The contexts are the abstracts themselves, and the questions are construct by randomly replacing a biomedical entity in the title with a placeholder. fig. 2 gives a sample in the BIOMRC dataset.

4.2 Experiment Settings

Our experiments are carried out on the machine with Intel i9-10920X (24) @ 4.700GHz, GPU of GeForce GTX 3090 24G, using pytorch 1.9.1 as the deep learning framework. To avoid overfitting, all models are trained for a maximum of 40 epochs, using early stopping on the dev, with a patience of 3 epochs.

Table 1: THE RESULT OF DIFFERENT AGGREGATION FUNCTIONS, COMPARED TO THE STATE-OF-THE-ART MODEL AND HUMAN EXPERTS

METHOD	Occurrence Aggregation	Token Aggregation	Train Time ¹	BIOMRC LITE		BIOMRC TINY ²
				Dev Acc	Test Acc	Test Acc
AS-READER	-	-	16.56hr	62.29	62.38	66.67
AoA-READER	-	-	60.90hr	70.00	69.87	70.00
SCIBERT-MAX-READER	-	-	83.22hr	80.06	79.97	90.00
HUMAN EXPERTS	-	-	-	-	-	85.00
AoA-READER WITH BioBERT EMBEDDING	max	max	1.50hr	78.54	78.11	90.00
	max	sum	0.88hr	83.40	83.36	93.33
	sum	max	3.60hr	80.98	81.20	90.00
	sum	sum	1.76hr	87.22	86.74	93.33

1: We conduct some code optimizations on the AoA-Reader model, so the training time of our implementation with BioBERT can not be compared to their original implementation.

2: The test set of BIOMRC TINY dataset only contains 30 samples, and so the results on it may be unstable. On the other hand, the demonstrated accuracy of human experts comes from averaging the results of multiple experts, so it is a bit more stable than other results.

During the process of fine-tuning SciBERT, the batch size is set to 1 and the top layer of SciBERT is frozen; other layers are trained with the learning rate of 0.001.

During the process of fine-tuning BioBERT and training AoA Reader, the batch size is set to 30, the learning rate is set to 0.001, and the learning rate for BioBERT is set to 10^{-5} . To reduce GPU memory usage, we use the mixed-precision training technique [18], setting precision to 16 bits.

We train our model on the BIOMRC LITE dataset and evaluate it both on the BIOMRC LITE and TINY dataset, which have 100,000 and 30 samples, respectively. We use Setting A for BIOMRC, in which all pseudo-identifier like *@entity1* have a global scope, i.e., all biomedical entities have a unique pseudo-identifier in the whole dataset.

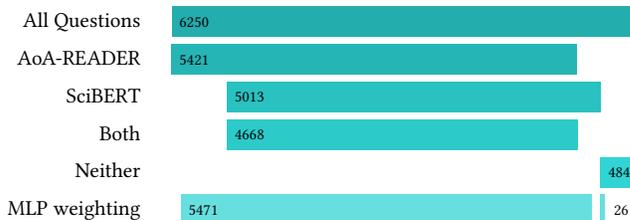
The model weighting layer is implemented after completing the training process of the two models, SciBERT and AoA Reader. The best weights evaluated by the *Dev Acc* are chosen to obtain the individual scores of each sample, which will later be used to train the model weighting layer.

4.3 Results

4.3.1 Performance of the Contextualized Word Embedding Strategy. Contextualized word embedding strategy based on BioBERT is used to obtain the final prediction answer. The selection of aggregating functions is crucial to the model performance. Therefore, multiple combinations of different aggregating functions are evaluated, and the results are shown in table 1.

It can be seen that choosing sum as both aggregation functions obtains better performance, and our model outperforms the state-of-the-art model significantly, which is about 6.77% absolute improvements on the BIOMRC LITE test sets.

Our model also shows an improvement on the BIOMRC TINY dataset, though the dataset contains only 30 samples, and this result may be unstable. Our performance on the larger BIOMRC LITE test

**Figure 3: The number of question answered correctly by different models on the BIOMRC LITE dataset.**

set still exceeds the average human expert performance on the BIOMRC TINY test set.

4.3.2 Performance of the Weighting Model. The MLP-based weighting model is used to further achieve better performance. We implement the SCIBERT-MAX-READER proposed by Pappas et al. [20], and use the model weighting strategy on top of AoA Reader with BioBERT embedding and SCIBERT-MAX-READER. The results of our experiments are shown in table 2. The result of SCIBERT-MAX-READER comes from our implementation of this model, which is used to train our MLP-based weighting layer. The results slightly differs to those in table 1.

It can be seen that our MLP-based weighting model improves the accuracy effectively. Especially, the accuracy on the test set is improved by 1.26% over the AoA Reader on the BIOMRC LITE test dataset.

In order to evaluate the effectiveness of the model weighting layer, we compare the result to the union accuracy of two single models, i.e., the percentage of the union of the questions answered correctly by the two models in the total number of questions.

The results when excluding data that both models failed to answer are shown in table 3. The number of questions answered correctly by different models on the BIOMRC LITE dataset is shown in

Table 2: THE RESULTS OF TWO SINGLE MODELS AND OUR MLP-BASED MODEL WEIGHTING MODEL, COMPARED TO THE UNION ACCURACY OF TWO SINGLE MODELS.

METHOD	BIOMRC LITE		BIOMRC TINY
	Dev Acc	Test Acc	Test Acc
AoA-READER WITH BioBERT EMBEDDING	87.22	86.74	93.33
SCIBERT-MAX-READER	79.74	80.21	86.67
MLP-BASED WEIGHTING MODEL (OURS)	88.76	88.00	96.66
THE UNION OF TWO SINGLE MODELS (IDEAL RESULT)	93.07	92.26	96.66

Table 3: THE RESULTS OF OUR MLP-BASED WEIGHTING MODEL, EXCLUDING DATA THAT BOTH MODELS FAILED TO ANSWER

METHOD	BIOMRC LITE		BIOMRC TINY ¹
	Dev Acc	Test Acc	Test Acc
AoA-READER WITH BioBERT EMBEDDING	93.71	93.21	96.55
SCIBERT-MAX-READER	86.67	86.94	89.66
MLP-BASED WEIGHTING MODEL (OURS)	95.36	95.38	100.00

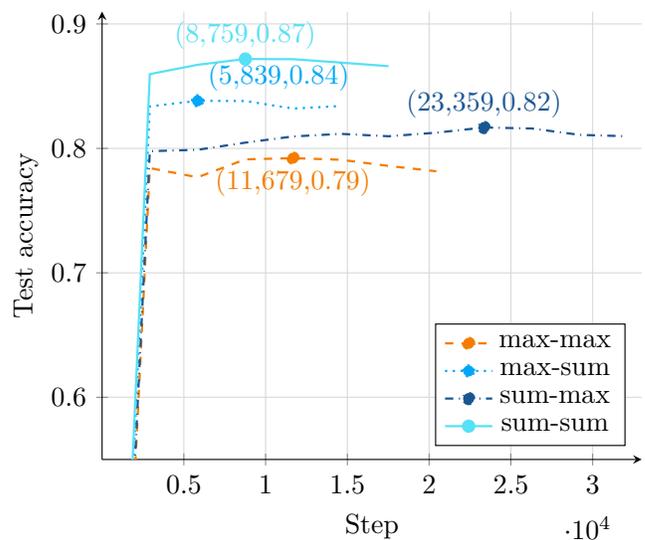
1: The test set of BIOMRC TINY dataset only contains 30 samples, and so the results on it may be unstable.

fig. 3. Here, *Both* in the figure refers to the the questions correctly answered by AoA-READER and SciBERT, and *Neither* in the figure refers to questions that cannot be answered correctly by any model.

As expected, both of the two models correctly answer some questions that the other model failed to answer. The proposed MLP-based weighting model not only gives the correct answer to the question that at least one model answers correctly, but also a small number of questions that both models fail to answer.

To further corroborate the model’s improvements in performance, we’ve applied the McNemar test to the results. Letting null hypothesis h_0 be our model has the same performance as SciBERT-MAX-READER, the alternate hypothesis h_1 would be there is a notable difference between the performance of our model and SciBERT-MAX-READER. The h_1 is accepted by the test where $n=2$, significance $\alpha = 0.025$. Since our model has a lower error rate in all tests, it sufficiently supports the hypothesis that our model has significantly outperformed the SciBERT-MAX-READER.

In general, our MLP-based weighting model improves the performance by 2.17% significantly compared to the original single model.

**Figure 4: The validation accuracy of AoA Reader after each epoch.**

Here, ‘max-max’ represent token aggregation and occurrence aggregation respectively.

These results demonstrate that the proposed method can automatically learn the biases and preferences and exploit the strengths of both models to achieve better performance.

4.3.3 Model Training Analysis. In our structure, SCIBERT-MAX-READER and AoA Reader are trained separately, the outputs of which are then gathered for the final weighting layer to learn.

The BioBERT embedded in the AoA Reader model is pre-trained on PubMed, and fine-tuned with the AoA Reader model on the BIOMRC dataset. It takes about 26 minutes for the AoA Reader to finish one epoch of learning based on the BIOMRC LITE dataset.

All models generated in epochs are saved, among which the one that has the best performance on the dev set will be used to train the weighting layer.

The training process for the AoA reader is illustrated in fig. 4. It can be seen that model using sum as both token and occurrence

aggregation converges faster compared to most models while giving best results.

5 CONCLUSIONS

We propose a contextual embedding and model weighting method, which can combine model pre-trained on a large corpus and open-domain QA model to mine semantic and contextual information in biomedical question answering. Especially, we adopt an MLP-based model weighting strategy which can automatically learn and utilize the preferences and biases of two models to combine their advantages. The results show that our method outperforms state-of-the-art system and has higher accuracy than experts. In future work, how to use the semantic similarity between entity tokens and context tokens in getting final predictions should be studied, i.e., a context token should contribute to the score of an entity if its semantic information is similar to that of entity token.

REFERENCES

- [1] Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A Pretrained Language Model for Scientific Text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3615–3620. <https://doi.org/10.18653/v1/D19-1371>
- [2] Steven Bird and Edward Loper. 2004. NLTK: The Natural Language Toolkit. In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*. Association for Computational Linguistics, Barcelona, Spain, 214–217. <https://doi.org/10.3115/1219044.1219075>
- [3] Wuya Chen, Xiaojun Quan, Chunyu Kit, Zhengcheng Min, and Jiahai Wang. 2020. Multi-Choice Relational Reasoning for Machine Reading Comprehension. In *Proceedings of the 28th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Barcelona, Spain (Online), 6448–6458. <https://doi.org/10.18653/v1/2020.coling-main.567>
- [4] Yiming Cui, Zhipeng Chen, Si Wei, Shijin Wang, Ting Liu, and Guoping Hu. 2017. Attention-over-Attention Neural Networks for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Vancouver, Canada, 593–602. <https://doi.org/10.18653/v1/P17-1055>
- [5] Yiming Cui, Ting Liu, Zhipeng Chen, Shijin Wang, and Guoping Hu. 2016. Consensus Attention-based Neural Networks for Chinese Reading Comprehension. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. The COLING 2016 Organizing Committee, Osaka, Japan, 1777–1786.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [7] Yongping Du, Jingya Yan, Yiliang Zhao, Yuxuan Lu, and Xingnan Jin. 2021. Dual Model Weighting Strategy and Data Augmentation in Biomedical Question Answering. In *Proceedings of the 2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*.
- [8] Chengzhen Fu, Yuntao Li, and Yan Zhang. 2019. ATNet: Answering Cloze-Style Questions via Intra-attention and Inter-attention. In *Advances in Knowledge Discovery and Data Mining (Lecture Notes in Computer Science)*, Qiang Yang, Zhi-Hua Zhou, Zhiguo Gong, Min-Ling Zhang, and Sheng-Jun Huang (Eds.). Springer International Publishing, Cham, 242–252. https://doi.org/10.1007/978-3-030-16145-3_19
- [9] Chengzhen Fu and Yan Zhang. 2019. EA Reader: Enhance Attentive Reader for Cloze-Style Question Answering via Multi-Space Context Fusion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 6375–6382. <https://doi.org/10.1609/aaai.v33i01.33016375>
- [10] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, et al. 2021. Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing. *ACM Transactions on Computing for Healthcare* 3, 1 (Oct. 2021), 2:1–2:23. <https://doi.org/10.1145/3458754>
- [11] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, et al. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 8342–8360. <https://doi.org/10.18653/v1/2020.acl-main.740>
- [12] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, et al. 2015. Teaching Machines to Read and Comprehend. In *Advances in Neural Information Processing Systems*, Vol. 28. Curran Associates, Inc.
- [13] Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2016. The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.).
- [14] Rudolf Kadlec, Martin Schmid, Ondrej Bajgar, and Jan Kleindienst. 2016. Text Understanding with the Attention Sum Reader Network. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Berlin, Germany, 908–918. <https://doi.org/10.18653/v1/P16-1086>
- [15] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *International Conference on Learning Representations*.
- [16] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, et al. 2019. BioBERT: A Pre-Trained Biomedical Language Representation Model for Biomedical Text Mining. *Bioinformatics* (Sept. 2019), btz682. <https://doi.org/10/ggh5qq>
- [17] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, et al. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv:1907.11692 [cs]* (July 2019). [arXiv:1907.11692 \[cs\]](https://arxiv.org/abs/1907.11692)
- [18] Paulius Mikićevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, et al. 2018. Mixed Precision Training. In *International Conference on Learning Representations*.
- [19] Dimitris Pappas, Ion Androutsopoulos, and Haris Papageorgiou. 2018. BioRead: A New Dataset for Biomedical Reading Comprehension. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. European Language Resources Association (ELRA), Miyazaki, Japan.
- [20] Dimitris Pappas, Petros Stavropoulos, Ion Androutsopoulos, and Ryan McDonald. 2020. BioMRC: A Dataset for Biomedical Machine Reading Comprehension. In *Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing*. Association for Computational Linguistics, Online, 140–149. <https://doi.org/10.18653/v1/2020.bionlp-1.15>
- [21] Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know What You Don't Know: Unanswerable Questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Melbourne, Australia, 784–789. <https://doi.org/10.18653/v1/P18-2124>
- [22] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, 2383–2392. <https://doi.org/10.18653/v1/D16-1264>
- [23] Matthew Richardson, Christopher J.C. Burges, and Erin Renshaw. 2013. MCTest: A Challenge Dataset for the Open-Domain Machine Comprehension of Text. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Seattle, Washington, USA, 193–203.
- [24] Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2017. Bidirectional Attention Flow for Machine Comprehension. In *The International Conference on Learning Representations*. [arXiv:1611.01603](https://arxiv.org/abs/1611.01603)
- [25] Raphael Tang, Rodrigo Nogueira, Edwin Zhang, Nikhil Gupta, Phuong Cam, Kyunghyun Cho, et al. 2020. Rapidly Bootstrapping a Question Answering Dataset for COVID-19. *arXiv:2004.11339 [cs]* (April 2020). [arXiv:2004.11339 \[cs\]](https://arxiv.org/abs/2004.11339)
- [26] George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R. Alvers, et al. 2015. An Overview of the BIOASQ Large-Scale Biomedical Semantic Indexing and Question Answering Competition. *BMC Bioinformatics* 16, 1 (April 2015), 138. <https://doi.org/10.1186/s12859-015-0564-6>
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, et al. 2017. Attention Is All You Need. In *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc.
- [28] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, et al. 2016. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *CoRR abs/1609.08144* (2016).