

Parsing Clinical Trial Eligibility Criteria for Cohort Query by a Multi-Input Multi-Output Sequence Labeling Model

Shubo Tian
Department of Statistics
Florida State University
Tallahassee, USA
stian2@fsu.edu

Pengfei Yin
Department of Health Outcomes
and Biomedical Informatics
University of Florida
Gainesville, USA
pengfeiyin@ufl.edu

Hansi Zhang
Department of Health Outcomes
and Biomedical Informatics
University of Florida
Gainesville, USA
hansi.zhang@ufl.edu

Arslan Erdengasileng
Department of Statistics
Florida State University
Tallahassee, USA
fel8b@my.fsu.edu

Jiang Bian
Department of Health Outcomes
and Biomedical Informatics
University of Florida
Gainesville, USA
bianjiang@ufl.edu

Zhe He, IEEE Senior Member
School of Information
Florida State University
Tallahassee, USA
zhe@fsu.edu

Abstract— To enable electronic screening of eligible patients for clinical trials, free-text clinical trial eligibility criteria should be translated to a computable format. Natural language processing (NLP) techniques have the potential to automate this process. In this study, we explored a supervised multi-input multi-output (MIMO) sequence labelling model to parse eligibility criteria into combinations of fact and condition tuples. Our experiments on a small manually annotated training dataset showed that the performance of the MIMO framework with a BERT-based encoder using all the input sequences achieved an overall lenient-level AUROC of 0.61. Although the performance is suboptimal, representing eligibility criteria into logical and semantically clear tuples can potentially make subsequent translation of these tuples into database queries more reliable.

Keywords— *Clinical trial, Eligibility criteria, Natural language processing*

I. INTRODUCTION

Randomized controlled trials are the gold-standard for evaluating the efficacy and safety of a treatment or intervention. Nevertheless, clinical trials often suffer from delayed patient accrual or insufficient participants, which may lead to early termination and cause significant financial loss for the sponsor. With the wide adoption of electronic health records (EHR), real-world EHR data allow us to evaluate the recruitment feasibility [1], perform electronic screening [2], and assess the generalizability of the trials before enrollment [3]. A necessary step to automate these analyses is to identify patients in the EHR data who satisfy the eligibility criteria of the trial, which are free-text sentences expressed in natural language and often with semantic ambiguities. It is thus important to extract key

elements from eligibility criteria and translate them into computable database queries. Natural language processing (NLP) is a key technology to facilitate such translation.

Typically, parsing eligibility criteria consists of 5 major tasks: (1) sentence chunking, (2) named-entity recognition (NER) and concept mapping, (3) relationship extraction, (4) temporal constraint detection, and (5) negation detection. Depending on the specific techniques, some tasks (e.g., NER and relation extraction) can be done in a single joint model. Manual annotation of eligibility criteria is required for building a robust criteria parser but it is expensive, labor intensive, and requires clinical domain knowledge. Therefore, an open question is “how to build a robust parser that can simultaneously perform multiple parsing tasks with limited annotated data of eligibility criteria?” In this work, we aim to investigate the use of a supervised multi-input multi-output (MIMO) sequence labelling model [4] to parse eligibility criteria. This architecture has two modules: a MIMO sequence labelling model, and a self-training method based on heuristic rule correction. In this architecture, multiple input sequences that can be generated automatically include: (1) word embeddings of the original text; (2) part-of-speech tags; (3) language model representation; and (4) concept, attribute, phrase (CAP) tagging. The tag sequences, which must be labelled manually, can be converted into fact and condition tuples jointly (i.e., multiple output). Expressing eligibility criteria in these tuples makes it possible to represent the named entities, temporal constraints (often as conditions), negations, and their relationships in a single universal framework. In this preliminary work, we demonstrate the feasibility of this approach for parsing eligibility criteria with a small labelled dataset.

II. RELATED WORK

A number of NLP systems for clinical trial eligibility criteria parsing have been developed previously. These systems can be categorized into (1) rule-based, and (2) machine learning-based systems. Rule-based parsers (e.g., EliXR [5], ValX [6], rely on predefined rules, which may not be robust enough to handle complex criteria (e.g., unseen patterns). One the other hand, machine learning-based parsers (e.g., ELiE [7], Criteria2Query [8]) are robust, but require a large training corpus with annotated data to achieve satisfactory performance. Recently, two large manually annotated eligibility criteria datasets were released: the Chia data with 1000 trials [9] and the Facebook Research Data with 3314 trials [10]. Tian et al. [11] recently benchmarked 4 transformers-based NER models on these two datasets and RoBERTa pretrained with MIMIC-III clinical notes and eligibility criteria yielded the highest strict and relaxed F-scores in experiments with both datasets. Further, these existing methods often do not emphasize the representations of the parsing results, leading to difficulty of reusing the annotated training data or the parsing results.

III. METHODS

A. Data Source and Data Annotation

Eligibility criteria of Alzheimer's disease (AD) trials.

From the ClinicalTrials.gov, we obtained free-text eligibility criteria of 13 phase III and IV AD clinical trials for existing Food and Drug Administration (FDA) approved AD drugs.

Annotation process. We followed the tagging schema (i.e., "B/I-XYZ" and "O") in the original MIMO study [4] to annotate the eligibility criteria, where:

- B: beginning, I: inside;
- $X \in \{\text{fact, condition}\}$;
- $Y \in \{1: \text{subject}; 2: \text{relation}; 3: \text{object}\}$;
- $Z \in \{\text{concept, attribute, predicate}\}$.

We decomposed each eligibility criteria into a set of fact and condition tuples. The tags in "B/I-XYZ" format are used for tagging word tokens of each component in the fact and condition tuples, where "B" represents the start word of a tuple component, "I" represents words other than the start word of a tuple component; " $X \in \{f, c\}$ " represent the tuple types of fact (f) and condition (c); both fact and condition tuples are represented by 3 components (1) subject, (2) predicate, and object (3); and " $Z \in \{C, A, P\}$ " represent the component roles of concept (C), attribute (A) and predicate phrase (P). Using this format, each word of eligibility criteria can be annotated into 10 different tags as shown in TABLE I. Note that any words not in a component of fact and condition tuples are tagged as "O". An example is shown in Fig. 1.

Following this annotation schema, we developed an annotation guideline specially for annotating eligibility criteria. We completed the annotations in multiple rounds, and our annotation process is shown in Fig. 2. In each round of annotation, 2 trials were annotated by 2 annotators based on the annotation guideline and Kappa scores were calculated [12]. Conflicts

between the two annotators were resolved by a third annotator and discussed with the entire study team.

TABLE I. THE EXAMPLES OF "B/I-XYZ" TAGGING SCHEMA

	Fact	Condition
Subjects	B/I-f1C	B/I-c1C
Subject Attributes	B/I-f1A	B/I-c1A
Predicates	B/I-f2P	B/I-c2P
Objects	B/I-f3C	B/I-c3C
Object Attributes	B/I-f3A	B/I-c3A

===== NCT00097916 stmt3 =====							
WORD	Stable	dose	of	donepezil	for		3 months
POSTAG	JJ	NN	IN	NN	IN	CD	NNS
CAP	O	O	O	B-Drug	B-Temporal	I-Temporal	I-Temporal
f1	O	O	O	B-f3C	O	O	O
c1	B-c3C	I-c3C	O	B-c1C	O	O	O
c2	O	O	O	B-c1C	B-c2P	B-c3C	I-c3C

Fig. 1. An example of annotating a criterion

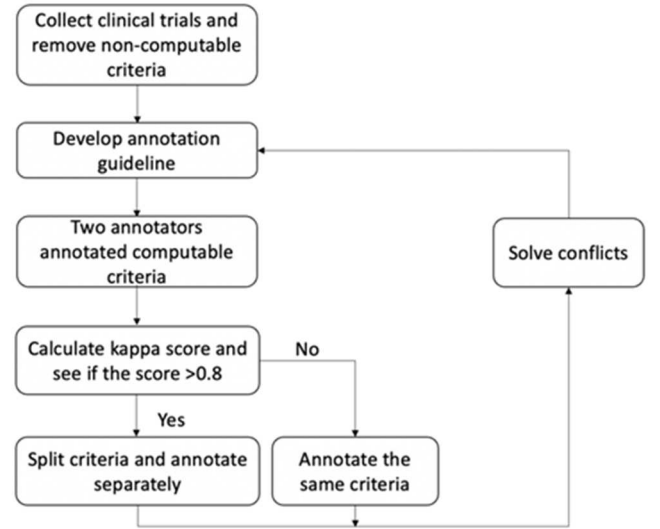


Fig. 2. The annotation process

B. The Multi-Input Multi-Output Sequence Labeling Model

The multi-input multi-output sequence labeling model, named as MIMO, was proposed by Jiang et al. as a framework for extracting fact and condition tuples from scientific text [4]. The advantage of MIMO is that it not only extracts the factual statements (i.e., fact tuples), but also considers the conditions when the fact tuples are true. The MIMO framework has two modules: (1) a multi-input module that takes four input sequences including pre-trained word embeddings, pre-trained language model outputs, part-of-speech (POS) tags, and CAP (i.e., Concepts, Attributes, and Phrases) tags of a sentence and uses a multi-head encoder-decoder model to generate a sequence representation of the input sentence. The multi-input gates were implemented to control the use of different input sequences [13]; (2) a multi-output module that takes the sequence representation

output of the multi-input module as input and predicts multiple tuple tag sequences for the fact and condition tuples. The multi-output module consists of a tuple component tagging layer, which predicts the tag sequences for fact and condition tuple components, and a tuple completion tagging layer, which predicts multiple tag sequences for the fact tuples and condition tuples. Finally, the complete fact and condition tuples were extracted from the predicted fact and condition tuple tag sequences, respectively, using the matching function as in [14].

Readers who are interested in the framework can refer to the original paper for more details. The code of the MIMO framework is publicly available at: https://github.com/twjjiang/MIMO_CFE.

C. Evaluation Metrics

We use standard evaluation metrics of precision (P), recall (R) and f1 score (F1) at strict and lenient levels to evaluate performance of the MIMO framework for component tagging and tuple extraction of both fact and condition tuples.

For evaluation of component tagging, each fact or condition tuple component was considered as a named entity. Strict level evaluation requires exact match between the predicted component and the ground truth annotated component for each type of components of fact and condition tuples. Lenient level evaluation requires only overlap between the predicted component and the ground truth annotated component.

For evaluation of tuple extraction, strict level evaluation requires exact match between the extracted tuple and the ground truth annotated tuple for each fact and condition tuple, i.e., each component of the extracted tuple matches exactly each component of the annotated tuple for all 5 different components of subject, subject attribute, predicate, object, and object attribute in each tuple. At the lenient level, an extracted tuple was considered as approximately matched as long as the subject, predicate and object of the extracted tuple overlap with a ground truth annotated subject, predicate and object respectively. An example illustrating exact and approximate match of tuples is shown in TABLE II.

TABLE II. EXAMPLES OF EXACT MATCH AND APPROXIMATE MATCH

Example of Exact Match and Approximate Match	
A criterion of the trial NCT00428389 states: <i>"Have received continuous treatment with donepezil for at least 6 months prior to screening, and received a stable dose of 5 mg/day or 10 mg/day for at least the last 3 of these 6 months."</i>	
Two of the annotated condition tuples are:	
A1	{{('continuous treatment with donepezil', 2, 6), 'NIL', ('at least', 7, 9), ('6 months', 9, 11), 'NIL'}}
A2	{{('continuous treatment with donepezil', 2, 6), 'NIL', ('at least', 27, 29), ('last 3 of these 6 months', 30, 36), 'NIL'}}
Two of the predicted condition tuples are:	
P1	{{('continuous treatment with donepezil', 2, 6), 'NIL', ('a least', 7, 9), ('6 months', 9, 11), 'NIL'}}
P2	{{('continuous', 2, 3), 'NIL', ('for at least', 26, 29), ('last 3', 30, 32), 'NIL'}}

The predicted tuple P1 is considered as exact match with A1 while P2 is considered as approximate match with A2. Then the precision (P), recall (R) and f1 score (F1) can be calculated using standard formulas based on true positive (TP) (i.e., the number of exact or approximate match tuples), false positive (FP) (i.e., the number of unmatched extracted tuples), and false negative (FN) (i.e., the number of tuples not being extracted).

IV. EXPERIMENT AND RESULTS

We conduct the experiments by implementing the MIMO framework with our annotated data and reusing the code made publicly available on GitHub with minor changes to accommodate our workflow and report the results as follows.

A. Experiment

The MIMO framework includes models of different architectures. We selected the MIMO framework with a BERT-based encoder for our experiment because it outperformed frameworks with other architectures as reported by the authors [4].

TABLE III. TUPLES AND COMPONENTS IN TRAINING AND TEST DATA

	Fact	Condition	Total
Training Data			
Tuples	188	110	298
Subjects	16	79	95
Subject Attributes	0	1	1
Predicates	21	68	89
Objects	185	90	275
Object Attributes	1	0	1
Total Components	223	238	461
Test Data			
Tuples	121	102	223
Subjects	20	55	75
Subject Attributes	0	0	0
Predicates	14	60	74
Objects	112	80	192
Object Attributes	0	0	0
Total Components	146	195	341

We split our annotated data by randomly selecting 8 trials as training data and using the remaining 5 trials as test data. The number of tuples and components in the training and test data is given in TABLE III. Following the best practice in [15], we used the NLTK (Natural Language Toolkit) package for word tokenization and POS tagging of the input sentences. We obtained the word embeddings with dimension of 50 from the MIMO repository on GitHub. The MIMO framework with a BERT-based encoder uses BERT [16] as the pre-trained language model. For CAP tagging, we used the NER tags predicted by a NER model based on RoBERTa [17], a transformer-based model first pre-trained with general English corpora and further pretrained with MIM-IC-III clinical notes [18] and eligibility criteria extracted from more than 350,000 clinical trial summaries on ClinicalTrials.gov. The RoBERTa NER model was then fine-tuned with a dataset derived from Chia, a corpus containing more than 12,000 annotated eligibility

criteria from 1,000 Phase IV trials in ClinicalTrials.gov [11]. We included entities of 6 major types including Condition, Value, Procedure, Drug, Measurement and Temporal in the derived dataset for training the NER model.

We used the default hyperparameters set in the MIMO framework and trained the MIMO framework [4] with a BERT-based encoder using different combination of inputs. We experimented with the different sets of input sequences and evaluated the performance using the evaluation metrics described in Section 3.3.

B. Results

Our experiment results show that the MIMO framework with a BERT-based encoder using all inputs of pre-trained word embeddings, pre-trained language model outputs, POS tags, and CAP tags achieves the best performance in terms of all evaluation measures at strict level for extraction of both fact and condition tuples, and achieves the best performance in terms of precision and f1 score for extraction of fact tuples at the lenient level. Detailed experimental results of the MIMO framework with a BERT-based encoder using all the input sequences are given in TABLE IV. and TABLE V.

TABLE IV. TABLE TYPE STYLES

	P	R	F1
Tuples			
Fact	0.347	0.554	0.427
Condition	0.067	0.078	0.072
Total	0.240	0.336	0.280
Fact Components			
Subject	0.550	0.550	0.550
Subject Attribute	0.000	0.000	0.000
Predicate	0.409	0.643	0.500
Object	0.506	0.723	0.596
Object Attribute	0.000	0.000	0.000
Total Components	0.500	0.692	0.581
Condition Components			
Subject	0.269	0.327	0.295
Subject Attribute	0.000	0.000	0.000
Predicate	0.531	0.433	0.477
Object	0.493	0.438	0.464
Object Attribute	0.000	0.000	0.000
Total Components	0.423	0.405	0.414
Total Components			
Subject	0.333	0.387	0.358
Subject Attribute	0.000	0.000	0.000
Predicate	0.493	0.473	0.483
Object	0.502	0.604	0.549
Object Attribute	0.000	0.000	0.000
Total Components	0.463	0.528	0.493

From our experiment, we observed that the MIMO framework tended to extract more tuples than the annotated gold-standard. This brings higher recall but lower precision in most of the experiments. Another observation is that the MIMO framework achieved better performance for fact tuple extraction than performance for condition tuple extraction. One of the reasons for this may be because condition tuples in eligibility criteria of clinical trial summaries are more complicated than

fact tuples. In addition, the small sample size of the annotated data from only 13 trial summaries may not be adequate for training a deep learning model to achieve a good performance.

TABLE V. TABLE TYPE STYLES

	P	R	F1
Tuples			
Fact	0.674	0.851	0.752
Condition	0.454	0.373	0.409
Total	0.590	0.632	0.610
Fact Components			
Subject	0.850	0.900	0.874
Subject Attribute	0.000	0.000	0.000
Predicate	0.500	0.714	0.588
Object	0.756	0.920	0.830
Object Attribute	0.000	0.000	0.000
Total Components	0.738	0.897	0.810
Condition Components			
Subject	0.582	0.582	0.582
Subject Attribute	0.000	0.000	0.000
Predicate	0.776	0.600	0.677
Object	0.747	0.563	0.642
Object Attribute	0.000	0.000	0.000
Total Components	0.695	0.580	0.632
Total Components			
Subject	0.644	0.667	0.655
Subject Attribute	0.000	0.000	0.000
Predicate	0.690	0.622	0.654
Object	0.753	0.771	0.762
Object Attribute	0.000	0.000	0.000
Total Components	0.717	0.716	0.716

V. DISCUSSION AND CONCLUSIONS

In this preliminary work, we evaluated the feasibility of using the MIMO framework to parse clinical trial eligibility criteria. Using 13 AD trials, we achieved a reasonable performance in terms of lenient-level F1 for recognizing components of fact (0.81) and condition tuples (0.72), respectively and then the entire tuples (0.61). The reason for the lower performance of condition tuples could be attributed to the small sample size. And the unsatisfactory performance of the strict-level evaluation is mainly due to inaccurate tuple components extraction. Nevertheless, representing eligibility criteria into logical and semantically clear fact and condition tuples can potentially make subsequent translation of these tuples into database queries more reliable. In future work, we will refine the annotation guideline and annotate more trials to increase the training samples. We will also integrate the results with the entity type recognition model (RoBERTa-MIMIC-Trial) that we previously built [11], which can potentially improve the model performance. We will explore ways of building database queries against real-world EHR data using the tuples and evaluate cohort identification performance.

ACKNOWLEDGMENT

This study was supported in part by the National Institutes of Health (NIH) under awards R21AG061431, R21AG068717, R21CA253394, and UL1TR001427.

REFERENCES

- [1] J. Doods, F. Botteri, M. Dugas, and F. Fritz, "A European inventory of common electronic health record data elements for clinical trial feasibility," *Trials*, vol. 15, no. 1, p. 18, Jan. 2014, doi: 10.1186/1745-6215-15-18.
- [2] S. R. Thadani, C. Weng, J. T. Bigger, J. F. Ennever, and D. Wajngurt, "Electronic Screening Improves Efficiency in Clinical Trial Recruitment," *J Am Med Inform Assoc*, vol. 16, no. 6, pp. 869–873, 2009, doi: 10.1197/jamia.M3119.
- [3] Z. He *et al.*, "Clinical Trial Generalizability Assessment in the Big Data Era: A Review," *Clin Transl Sci*, vol. 13, no. 4, pp. 675–684, Jul. 2020, doi: 10.1111/cts.12764.
- [4] T. Jiang, T. Zhao, B. Qin, T. Liu, N. Chawla, and M. Jiang, "Multi-Input Multi-Output Sequence Labeling for Joint Extraction of Fact and Condition Tuples from 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)*, Hong Kong, China: Association for Computational Linguistics, 2019, pp. 302–312. doi: 10.18653/v1/D19-1029.
- [5] C. Weng, X. Wu, Z. Luo, M. R. Boland, D. Theodoratos, and S. B. Johnson, "EliXR: an approach to eligibility criteria extraction and representation," *J Am Med Inform Assoc*, vol. 18, no. Suppl 1, pp. i116–i124, Dec. 2011, doi: 10.1136/amiajnl-2011-000321.
- [6] T. Hao, H. Liu, and C. Weng, "Valx: A system for extracting and structuring numeric lab test comparison statements from text," *Methods Inf Med*, vol. 55, no. 3, pp. 266–275, May 2016, doi: 10.3414/ME15-01-0112.
- [7] T. Kang *et al.*, "EliIE: An open-source information extraction system for clinical trial eligibility criteria," *J Am Med Inform Assoc*, vol. 24, no. 6, pp. 1062–1071, Nov. 2017, doi: 10.1093/jamia/ocx019.
- [8] C. Yuan *et al.*, "Criteria2Query: a natural language interface to clinical databases for cohort definition," *J Am Med Inform Assoc*, vol. 26, no. 4, pp. 294–305, Feb. 2019, doi: 10.1093/jamia/ocy178.
- [9] F. Kury *et al.*, "Chia, a large annotated corpus of clinical trial eligibility criteria," *Sci Data*, vol. 7, Aug. 2020, doi: 10.1038/s41597-020-00620-0.
- [10] Y. Tseo, M. I. Salkola, A. Mohamed, A. Kumar, and F. Abnoui, "Information Extraction of Clinical Trial Eligibility Criteria," *arXiv:2006.07296 [cs]*, Jul. 2020, Accessed: Jul. 20, 2021. [Online]. Available: <http://arxiv.org/abs/2006.07296>
- [11] S. Tian *et al.*, "Transformer-based named entity recognition for parsing clinical trial eligibility criteria," in *Proceedings of the 12th ACM Conference on Bioinformatics, Computational Biology, and Health Informatics*, in BCB '21. New York, NY, USA: Association for Computing Machinery, Aug. 2021, pp. 1–6. doi: 10.1145/3459930.3469560.
- [12] Stephanie, "Cohen's Kappa Statistic," *Statistics How To*, Dec. 09, 2014. <https://www.statisticshowto.com/cohens-kappa-statistic/> (accessed May 31, 2023).
- [13] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2016, pp. 770–778. doi: 10.1109/CVPR.2016.90.
- [14] G. Stanovsky, J. Michael, L. Zettlemoyer, and I. Dagan, "Supervised Open Information Extraction," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, New Orleans, Louisiana: Association for Computational Linguistics, Jun. 2018, pp. 885–895. doi: 10.18653/v1/N18-1081.
- [15] S. Bird, E. Klein, and E. Loper, *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*, 1st edition. Beijing ; Cambridge Mass.: O'Reilly Media, 2009.
- [16] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *arXiv:1810.04805 [cs]*, May 2019, Accessed: May 10, 2020. [Online]. Available: <http://arxiv.org/abs/1810.04805>
- [17] Y. Liu *et al.*, "RoBERTa: A Robustly Optimized BERT Pretraining Approach," *arXiv:1907.11692 [cs]*, Jul. 2019, Accessed: Dec. 13, 2019. [Online]. Available: <http://arxiv.org/abs/1907.11692>
- [18] A. E. W. Johnson *et al.*, "MIMIC-III, a freely accessible critical care database," *Scientific Data*, vol. 3, no. 1, Art. no. 1, May 2016, doi: 10.1038/sdata.2016.35.