

**Fig. 2.** Overall NER Performance in 2018 MADE 1.0 Challenge.

was an increase in the F1-score to 0.93 with the integration of CNN and RNN developed by Wei et al. [13]. The performances of ML-based models-SVM, RF, and GB proposed by Yang et al. [14], and SVM developed by Wei et al. [13] were comparable or even outperformed the DL-based models derived from Bi-LSTM and BERT.

Fig. 4 illustrates the overall RC performance in the 2018 MADE 1.0 Challenge. RF proposed by Chapman et al. [19] was equivalent or even better than most of the DL-based models explored. However, the performance of the SVM and RFs model developed by Yang et al. [20] was much inferior to other models.

Fig. A.4 and Fig. A.5 show the F1-score for each relation type drawn from the models employed in the 2018 n2c2 and 2018 MADE 1.0 Challenge, respectively. Similar to the NER task, the extractions of “Duration-Drug”, “Reason-Drug”, and “ADE-Drug” were difficult for most of the teams. However, Fig. A.4 shows that BERT and BERT-related models (BERT, BioBERT etc.) improved the performance in these relation extractions. For the training set in the 2018 n2c2 shared task on ADE and medication extraction challenge, approximately 30 % of “ADE” and “Reason” were not in the same sentences as the related drugs [8,13]. Some models were trained on a single sentence without the context information in the NER task [13,14,21]. Similarly, the models trained on a single sentence in the RC task were not able to classify relations across several sentences [13]. This issue can be handled by utilizing an attention-based neural network architecture that can exploit more

context information [10]. Second, the co-location information was insufficient to detect the “Drug-ADE” and “Drug-Reason” relations correctly when several drugs were located in the same sentence [4].

### 3.3. End-to-End system

Table A.1 shows the end-to-end systems proposed by each team and their corresponding performance. The BERT model developed by Mahendran and McInnes [4] achieved the highest F-1 score of 0.94.

## 4. Discussion

Through studying the methods published by the different research teams and analyzing their submitted results, the models with state-of-the-art performance, as well as the current state of the literature reviewed, significant gaps or flaws in existing knowledge, future areas for study, and associations between our research and existing knowledge were identified.

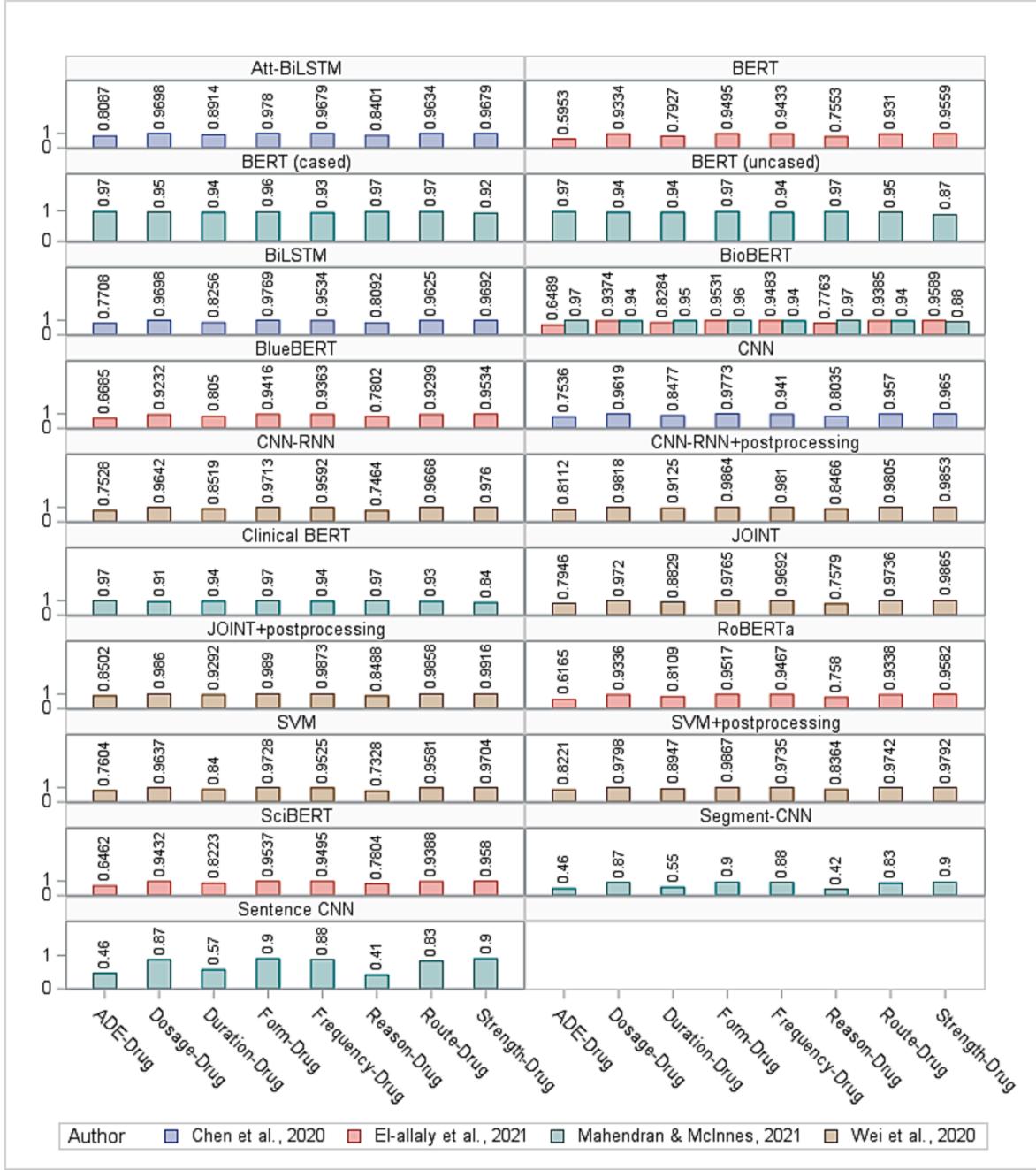
We provided the summary of advantages, limitations, and prospective avenues for further research for the methodologies under consideration in Table A.4. The methodologies considered for NER tasks, including CRF, BiLSTM-CRF, CNN-BiLSTM-CRF, LSTM-BiLSTM-CRF, CNN, RNN, LSTM-CRF, BERT, BioBERT, SciBERT, RoBERTa, and BlueBERT, each offer distinct strengths and weaknesses. CRF, while efficient











(a)

Fig. 4. Overall RC Performance in 2018 MADE 1.0 Challenge.

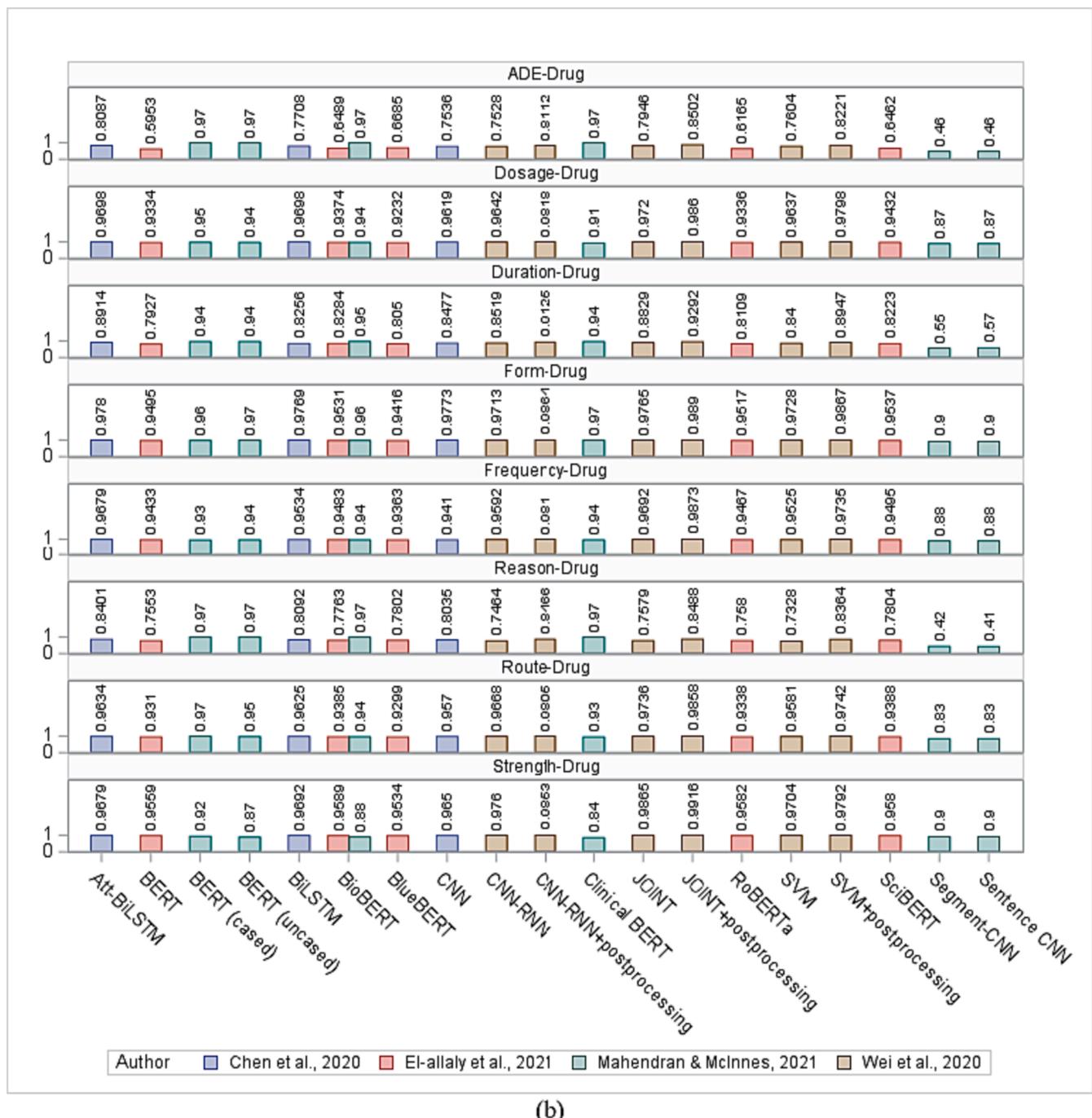
corpora and disease were included in the study [17].

Furthermore, while the literature predominantly focuses on machine learning and deep learning approaches, there is emerging interest in rule/ontology/knowledge-based methods [34,35] and semi-supervised learning methods [36,37]. Rule-based methods offer interpretability and explicit modeling of domain knowledge, but they can be limited by the complexity of rules and the need for manual intervention [35]. Semi-

supervised learning, on the other hand, leverages both labeled and unlabeled data, which can be beneficial when labeled data is scarce or costly to obtain [37].

#### 4.5. Future study

ADE is a critical and multifaceted aspect of biomedical research



(b)

Fig. 4. (continued).

[38,39]. Large language models (LLMs) have emerged as transformative tools in the domain of ADE extraction, promising unparalleled advancements in performance [40,41]. For example, Leas et al. found that ChatGPT accurately replicated human annotations for adverse event reports, achieving high agreement for both general (97.7 %; Kappa = 0.95) and serious adverse events (98.0 %; Kappa = 0.95) [42]. Li et al. applied and fine-tuned various LLMs, including GPT-2, GPT-3 variants,

GPT-4, and Llama 2 to evaluate the effectiveness of LLMs in extracting AEs related to the Influenza vaccine [40]. In the reviewed articles, the groundbreaking success of the BERT model, exemplified by its state-of-the-art F1-score of 0.94 in end-to-end tasks, positions it as a pivotal focus for future ADE extraction studies. However, only one drug-entity pair for the entire sentence was considered for the BERT model in the reviewed study, which did not represent all the actual cases [4]. Multiple drug-

entity pairs may exist in a sentence, which challenged the RC task in ADE extraction. Therefore, practical methods to represent multiple drug-entity pairs within a sentence should be constructed. In the work proposed by Mahendran and McInnes [4], an independent binary model for each class in the dataset was developed. In the future, a model for multiclass classification should be implemented and the effect of ensembles incorporated into multiple individual models should be studied. It has become a trend to employ character embeddings, word embeddings, and position embeddings as features. The coverage of selected embeddings directly affected the performance of ADE extraction [10]. In the study proposed by Dai et al. [10], the model with GloVe embedding greatly outperformed that with word2vec<sub>MIMIC</sub> and BioWordVec [10]. It remains to be studied which pre-trained embeddings will work best with BERT. A small dataset was used and it was challenging to extract ADE and reasons with limited training samples [13,14,17,20,21]. In the future study, a larger dataset with diverse medical domains will be included. As shown in Fig. 1 and Fig. 2, the integration of external resources into BiLSTM-CRF, LSTM-CRFs, and CRF did not enhance the performance in NER. It still needs to be explored whether the introduction of external resources works for BERT.

## 5. Conclusion

In conclusion, this study highlights the current state of ADE extraction research, emphasizing the potential of BERT models, effective embeddings, and the need for advanced strategies in entity extraction. The performance variations in NER and RC tasks suggest the importance of tailoring methods to specific challenges. Future research should focus on

optimizing the use of BERT models, tackling complex scenarios, leveraging pre-trained embeddings, and incorporating larger, more diverse datasets to further advance ADE extraction, ultimately enhancing patient safety and healthcare decision-making.

## CRediT authorship contribution statement

**Yiming Li:** Writing – original draft, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Wei Tao:** Visualization. **Zehan Li:** Writing – review & editing, Visualization. **Zenan Sun:** Data curation. **Fang Li:** Writing – review & editing. **Susan Fenton:** Writing – review & editing. **Hua Xu:** Writing – review & editing. **Cui Tao:** Writing – review & editing, Validation, Supervision, Project administration, Funding acquisition.

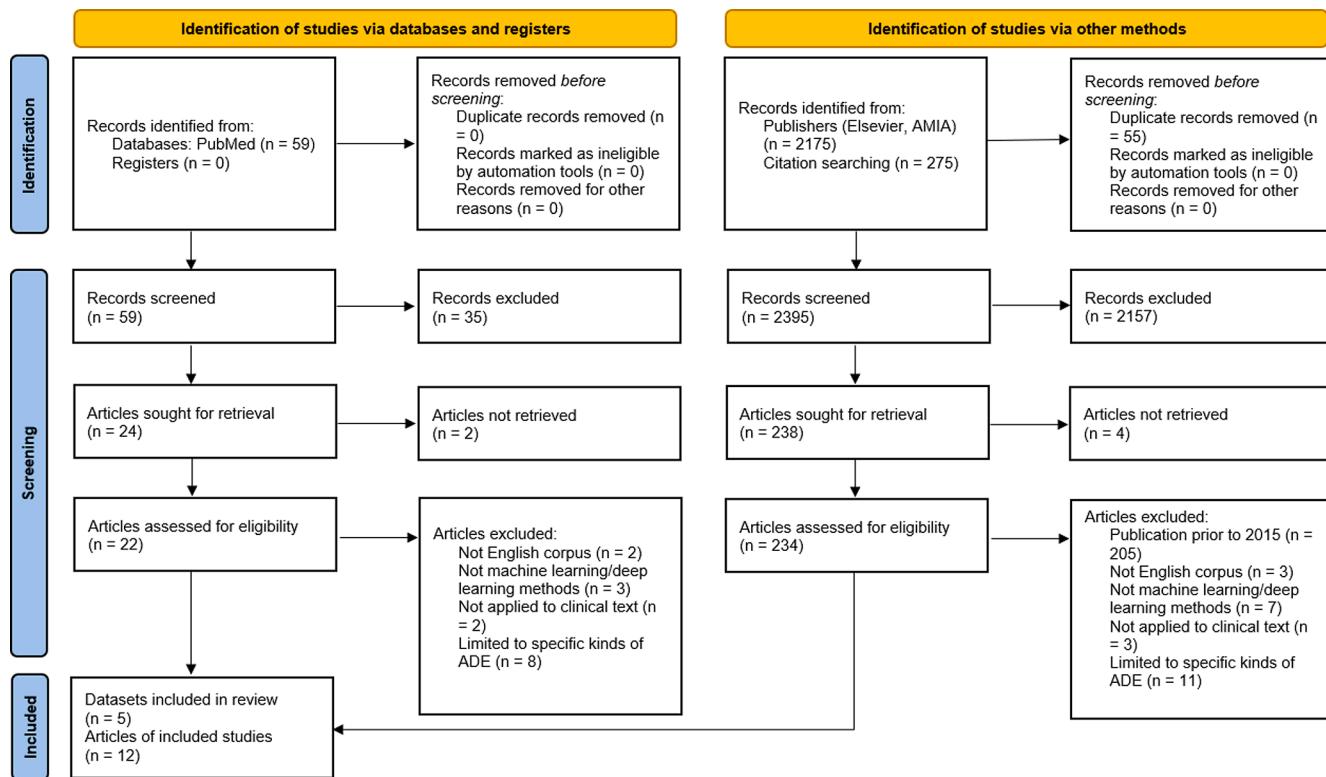
## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

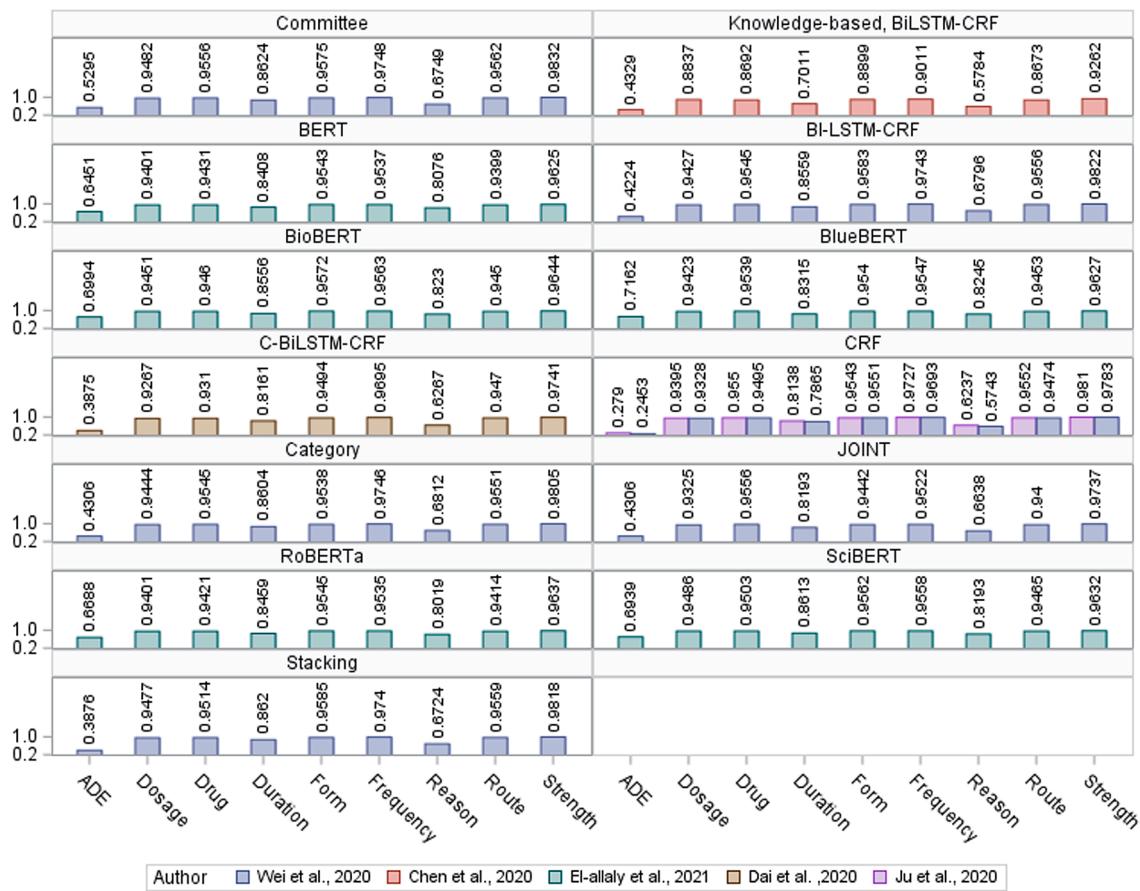
## Acknowledgments

This article was partially supported by the National Institute of Allergy And Infectious Diseases of the National Institutes of Health, United States under Award Numbers R01AI130460 and U24AI171008. I would like to express my sincere appreciation to Irmgard Willcockson for her invaluable contributions as the editor of this publication.

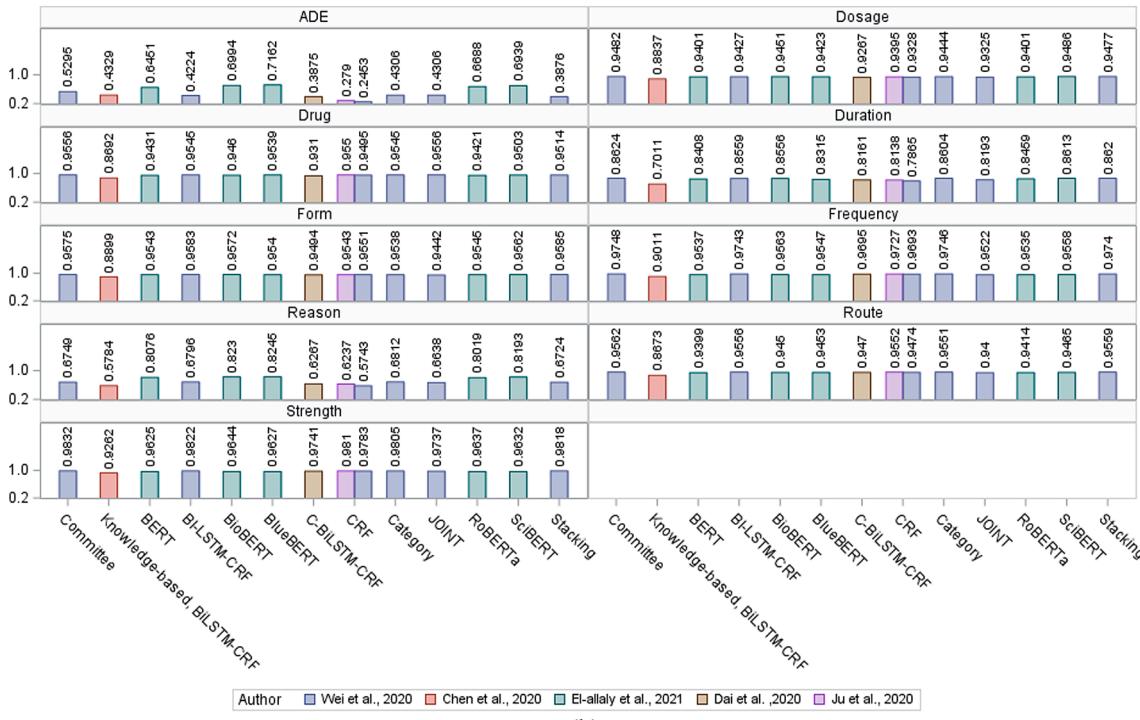
## Appendix A



**Fig. A1.** PRISMA 2020 flow diagram for new systematic reviews which included searches of databases, registers and other sources. Adapted from “The PRISMA 2020 statement: an updated guideline for reporting systematic reviews.” by Page et al. [43].

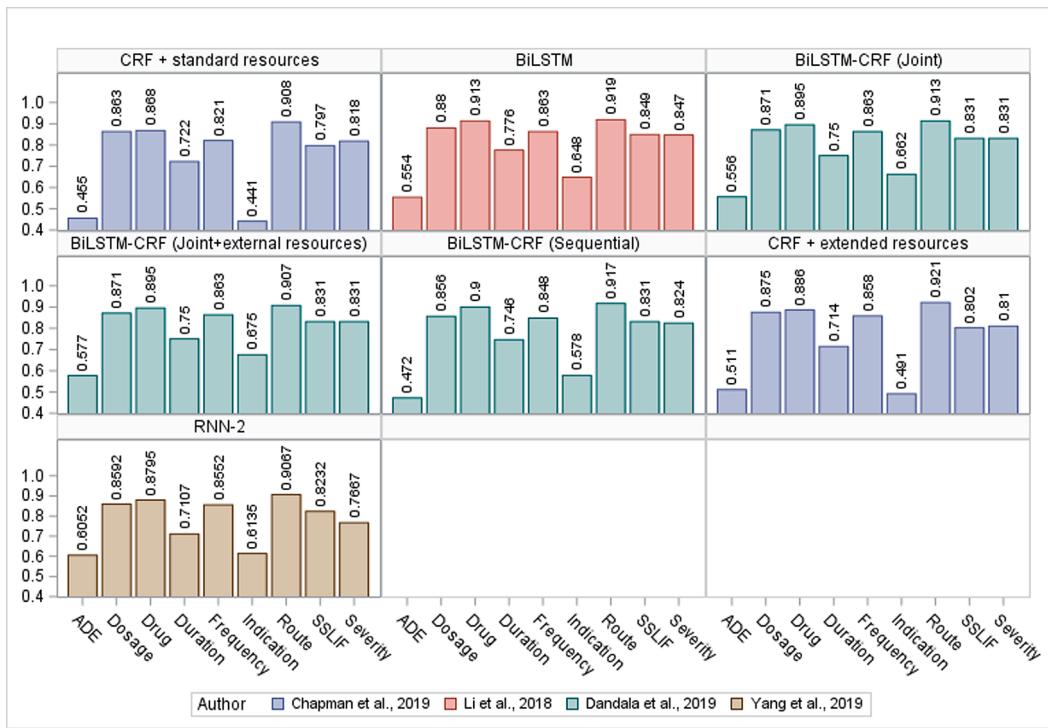


(a)

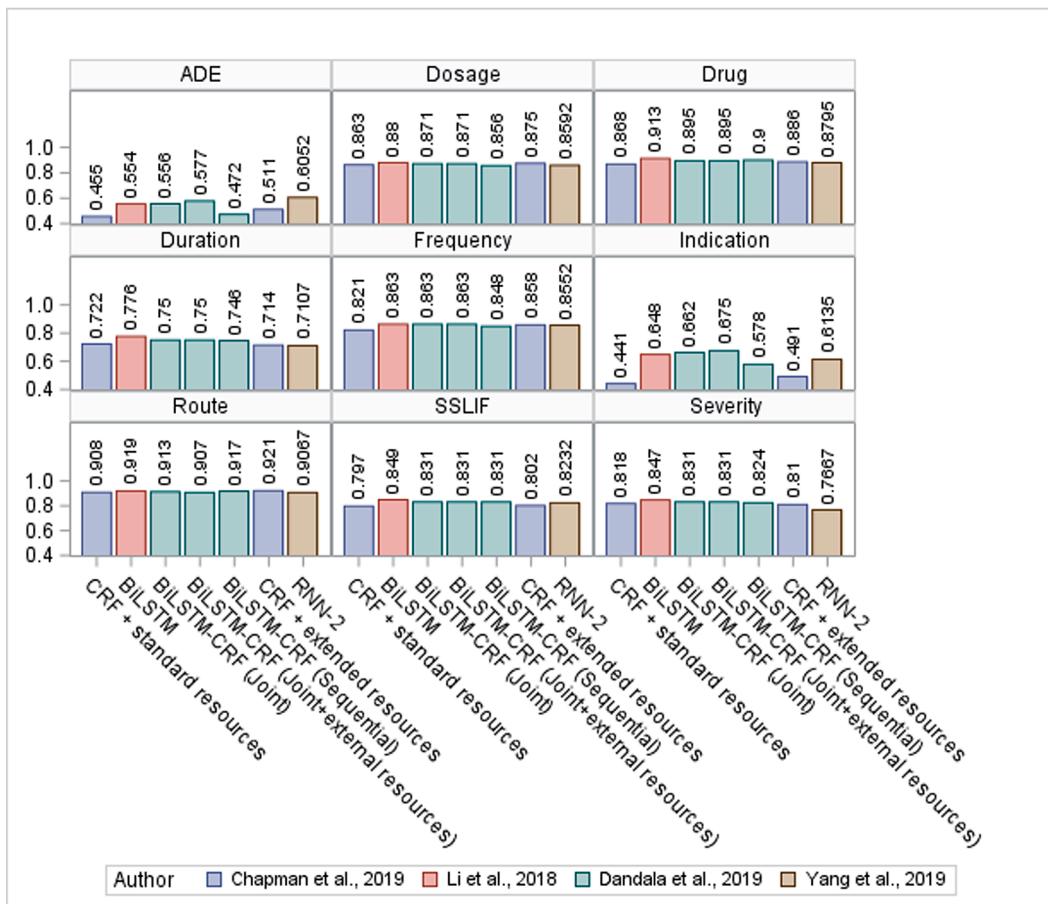


(b)

**Fig. A2.** a) F1 score for each entity type in 2018 n2c2 Shared Task on ADEs and Medication Extraction (method based) b) F1 score for each entity type in 2018 n2c2 Shared Task on ADEs and Medication Extraction (entity type based)

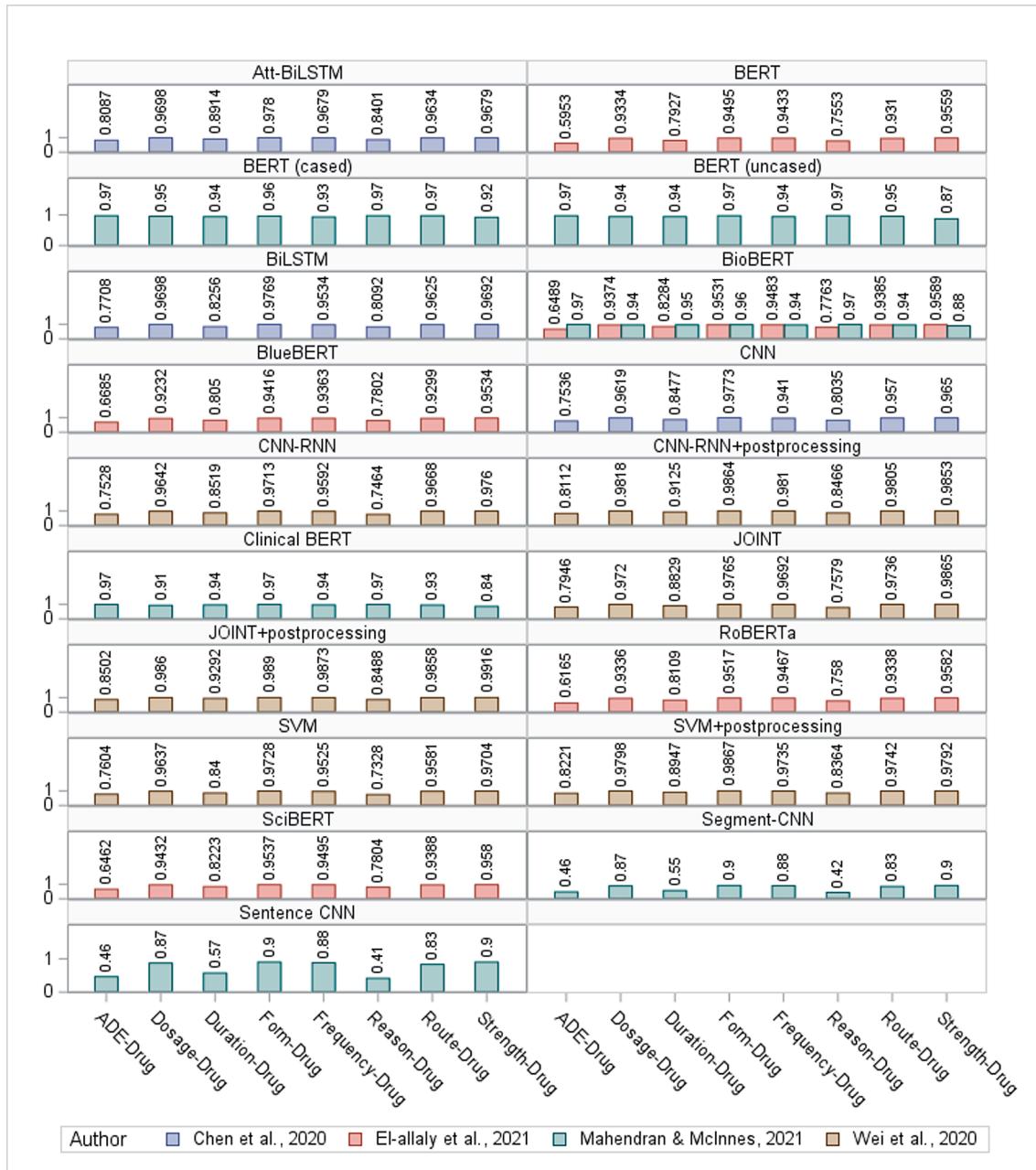


(a)



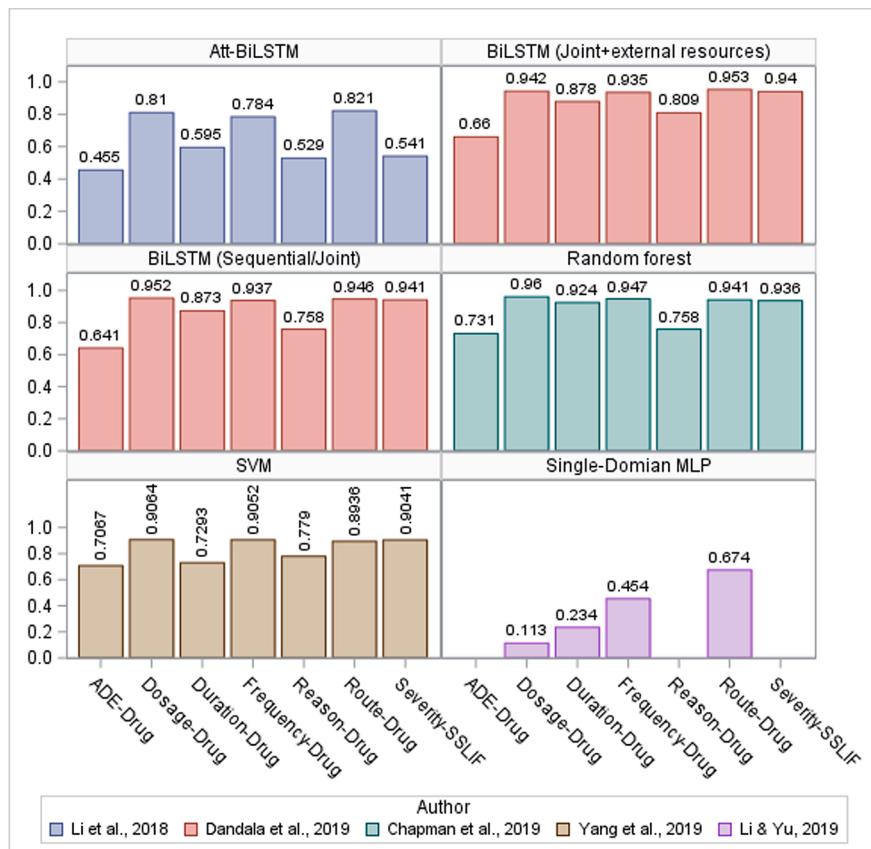
(b)

**Fig. A3.** a) F1 score for each entity type in 2018 MADE 1.0 Challenge (method based) b) F1 score for each entity type in 2018 MADE 1.0 Challenge (entity type based)

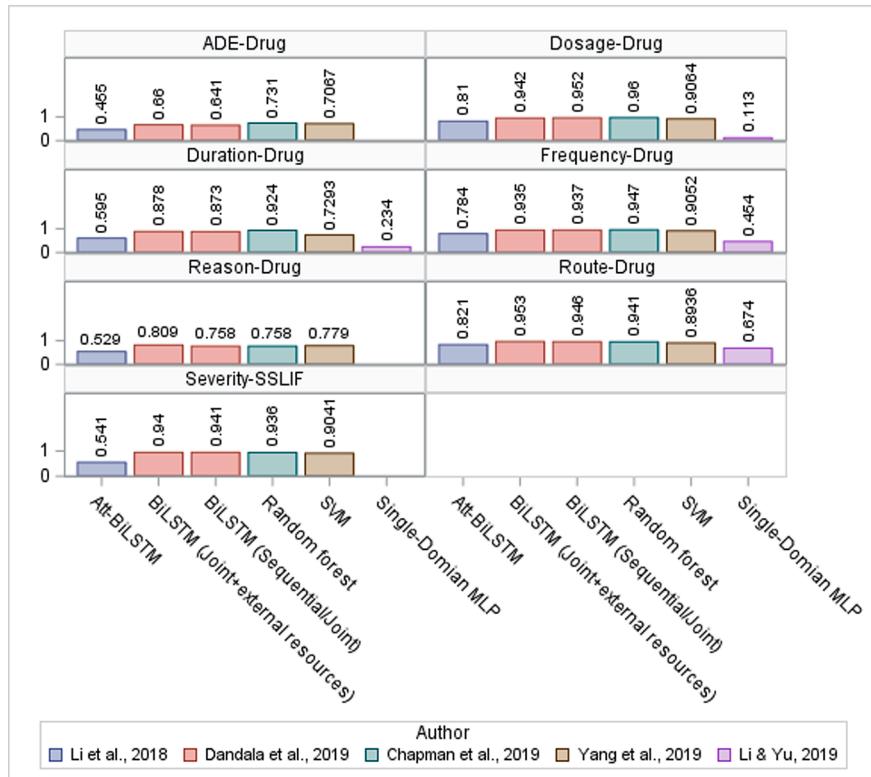


(a)

**Fig. A4.** a) F1 Score for Each Relation Type in 2018 n2c2 Shared Task on ADEs and Medication Extraction (method based) b) F1 Score for Each Relation Type in 2018 n2c2 Shared Task on ADEs and Medication Extraction (entity type based)

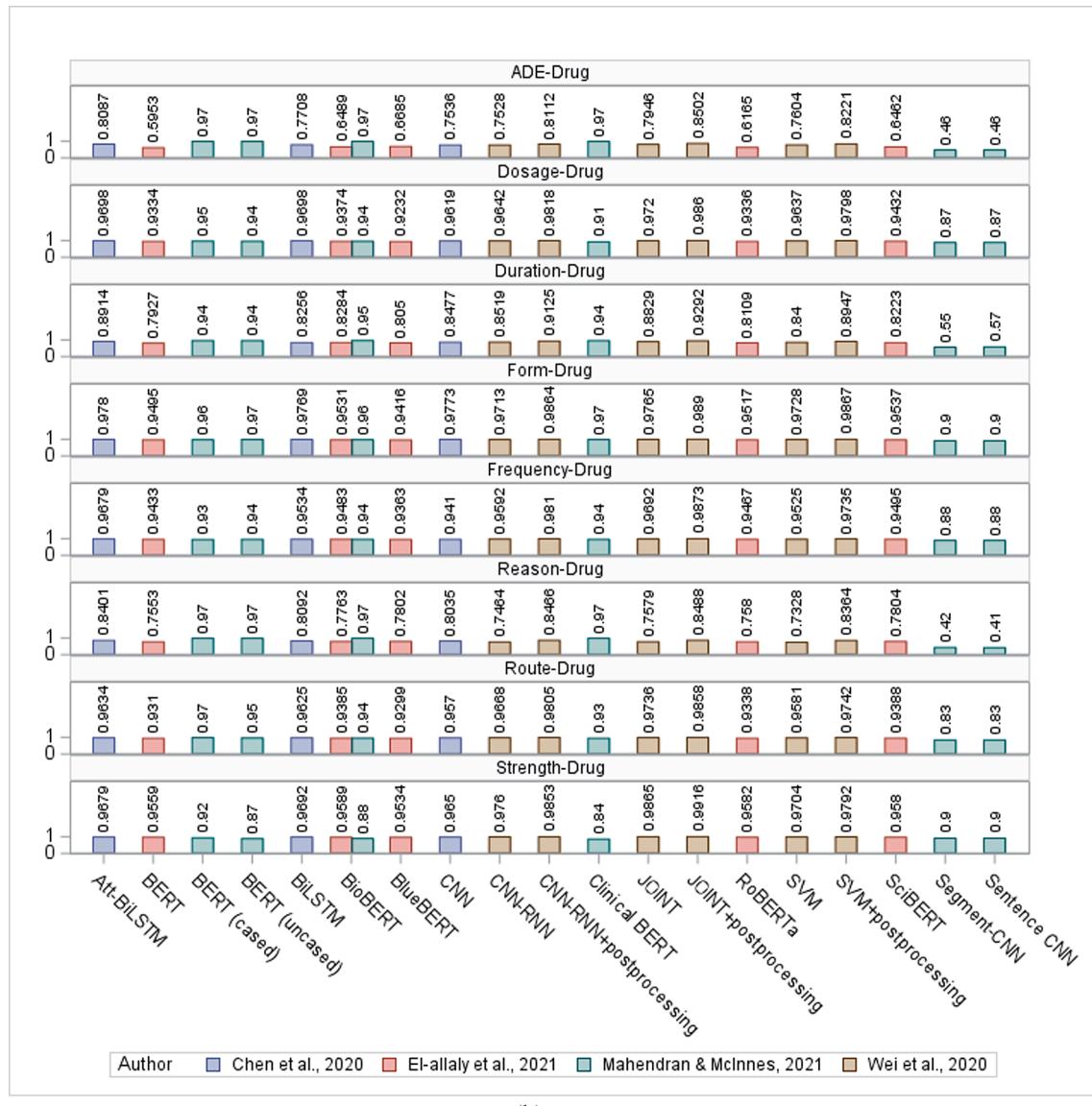


(a)



(b)

**Fig. A5.** a)F1 Score for Each Relation Type in 2018 MADE 1.0 Challenge (method based) b) F1 Score for Each Relation Type in 2018 MADE 1.0 Challenge (entity type based)



(b)

Fig. A4. (continued).

**Table A1**

Summary of included articles, related datasets, and their best end-to-end performances.

| Dataset                                  | Authors                   | Models used (NER model-RC model)   | Precision | Recall | F1                 |
|--|---------------------------|--|-----------|--------|--------------------|
| 2018 n2c2                                | Chen et al. [8]           | Knowledge-based- Attention-based Bidirectional long-short term memory (Att-BiLSTM)       | 0.8382    | 0.7539 | 0.7938             |
|  | Dai et al. [10]           | All Bidirectional long-short term memory (BiLSTM)- conditional random field (CRF) models | 0.939     | 0.900  | 0.919              |
|  | Ju et al. [15]            | BiLSTM-CRF (NER only)  | 0.9599    | 0.8979 | 0.9278             |
|  | Mahendran and McInnes [4] | Bidirectional Encoder Representations from Transformers (BERT)                           | 0.93      | 0.96   | 0.94               |
|  | Wei et al. [13]           | Convolutional neural network (CNN)- recurrent neural network (RNN)                       | NA        | NA     | 0.8905             |
|  | Yang et al. [14]          | LSTM-CRFs + gradient boosting (GB)-4   | 0.9187    | 0.8593 | 0.8880             |
| 2018 MADE 1.0 Challenge                  | Chapman et al. [19]       | CRF- random forest (RF)  | 0.721     | 0.534  | 0.612              |
|  | Dandala et al. [21]       | BiLSTM-CRF-BiLSTM (Joint + external resources)   | 0.696     | 0.632  | 0.662              |
|  | F. Li et al. [5]          | BiLSTM-CRF- BiLSTM-Attention   | NA        | NA     | 0.667              |
|  | Yang et al. [20]          | RNN-2-support vector machine (SVM)   | 0.5758    | 0.6542 | 0.6125             |
| 2018 MADE 1.0 Challenge<br>Cardio corpus | F. Li and Yu [17]         | Capsule network (CapNet)   | NA        | NA     | 0.827 <sup>a</sup> |

(continued on next page)





- Biomed. Health Inform. 24 (2020) 57–68, <https://doi.org/10.1109/JBHI2019.2932740>.
- [37] S. Gupta, S. Pawar, N. Ramrakhiyani, G.K. Palshikar, V. Varma, Semi-Supervised Recurrent Neural Network for Adverse Drug Reaction mention extraction, BMC Bioinf. 19 (2018) 212, <https://doi.org/10.1186/s12859-018-2192-4>.
- [38] Y. Li, J. Li, Y. Dang, Y. Chen, C. Tao, Temporal and Spatial Analysis of COVID-19 Vaccines Using Reports from Vaccine Adverse Event Reporting System, JMIR Prepr (2023), <https://doi.org/10.2196/preprints.51007>.
- [39] Y. Li, S.K. Lundin, J. Li, W. Tao, Y. Dang, Y. Chen, et al., Unpacking adverse events and associations post COVID-19 vaccination: a deep dive into vaccine adverse event reporting system data, Expert Rev. Vaccines 23 (2024) 53–59, <https://doi.org/10.1080/14760584.2023.2292203>.
- [40] Y. Li, J. Li, J. He, C. Tao, AE-GPT: Using Large Language Models to Extract Adverse Events from Surveillance Reports-A Use Case with Influenza Vaccine Adverse Events, 2023. ArXivorg. <https://doi.org/10.48550/arXiv.2309.16150>.
- [41] Y. Hu, I. Ameer, X. Zuo, X. Peng, Y. Zhou, Z. Li, Y. Li, J. Li, X. Jiang, H. Xu, Zero-shot Clinical Entity Recognition using ChatGPT, 2023. ArXivorg. <https://doi.org/10.48550/arXiv.2303.16416>.
- [42] E.C. Leas, J.W. Ayers, N. Desai, M. Dredze, M. Hogarth, D.M. Smith, Using Artificial Intelligence (ChatGPT) to Support Biomedical Text Analysis: A Case Study of Adverse Event Detection, JMIR Prepr (2023).
- [43] M.J. Page, D. Moher, P.M. Bossuyt, I. Boutron, T.C. Hoffmann, C.D. Mulrow, Prisma,, et al., explanation and elaboration: updated guidance and exemplars for reporting systematic reviews, BMJ 2021 (2020) n160, <https://doi.org/10.1136/bmj.n160>.