

Developing A Deep Learning Natural Language Processing Algorithm For Automated Reporting Of Adverse Drug Reactions

Christopher McMaster^{1,3,*}, Julia Chan², David FL Liew^{1,2}, Elizabeth Su¹, Albert G Frauman^{1,2}, Wendy W Chapman³, Douglas EVPires³

1 Department of Clinical Pharmacology & Therapeutics, Austin Health, Melbourne, Victoria, Australia

2 Department of Medicine, University of Melbourne, Melbourne, Victoria, Australia

3 The Centre for Digital Transformation of Health, University of Melbourne, Melbourne, Victoria, Australia

* christopher.mcmaster@austin.org.au

Abstract

The detection of adverse drug reactions (ADRs) is critical to our understanding of the safety and risk-benefit profile of medications. With an incidence that has not changed over the last 30 years, ADRs are a significant source of patient morbidity, responsible for 5-10% of acute care hospital admissions worldwide. Spontaneous reporting of ADRs has long been the standard method of reporting, however this approach is known to have high rates of under-reporting, a problem that limits pharmacovigilance efforts. Automated ADR reporting presents an alternative pathway to increase reporting rates, although this may be limited by over-reporting of other drug-related adverse events.

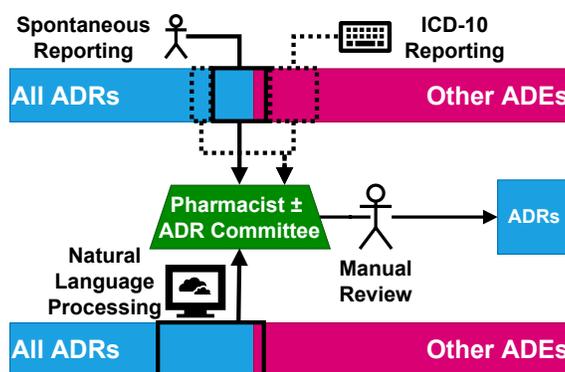
We developed a deep learning natural language processing algorithm to identify ADRs in discharge summaries at a single academic hospital centre. Our model was developed in two stages: first, a pre-trained model (DeBERTa) was further pre-trained on 150,000 unlabelled discharge summaries; secondly, this model was fine-tuned to detect ADR mentions in a corpus of 861 annotated discharge summaries. To ensure that our algorithm could differentiate ADRs from other drug-related adverse events, the annotated corpus was enriched for both validated ADR reports and confounding drug-related adverse events using. The final model demonstrated good performance with a ROC-AUC of 0.934 (95% CI 0.931 - 0.955) for the task of identifying discharge summaries containing ADR mentions.

1 Introduction

Adverse drug reactions (ADRs) are a subset of adverse drug events (ADEs), defined as injuries resulting from drugs administered at therapeutic doses [27]. With an incidence that has not changed over the last 30 years, ADRs are estimated to be responsible for 5-10% of acute care hospital admissions and one of the leading causes of death worldwide [8, 12, 25]. In Australia, ADRs result in approximately 2-4% of hospital admissions and are a significant source of patient morbidity [30].

Due to the challenges of clinical note processing, ADRs are currently reported spontaneously by clinical staff [9] or, less commonly, through the ICD-10 Y40.0-Y59.9 coding for ADEs [19] (see Figure 1). Spontaneous reporting is the predominant way in which ADRs are reported. Whilst it is a relatively accurate method, the under-reporting rate is estimated to range from 82-98% [9]. In Australia, clinician review and reporting of ICD-10 Y40.0-Y59.9 coded ADEs can capture some of these missed ADRs from spontaneous reporting, as ADEs are coded retrospectively by professional clinical coders who scan clinical data for patient diagnoses for the purposes of hospital funding. However, such utilisation for ADR detection is a relatively resource intensive process and not all hospitals have the capacity for this [22]. Further, not only is the ability of ICD-10 codes to capture ADRs yet to be validated [23], but it has also been demonstrated to capture less than half of the ADEs identified from medical record reviews [12]. With coding for ADRs performed by individuals not involved in clinical care, multiple opportunities for inconsistency in ADR reporting arise, including inter-coder variability [31], lack of context, knowledge, and time [12].

Figure 1. ADR reporting mechanisms. The current state of ADR reporting is primarily based on spontaneous reporting, with previous studies demonstrating a high rate of unreported ADRs. The addition of ICD-10 coding can increase this reporting rate, however many other ADEs are picked up in this process, placing strain on the review process. The aim of an NLP model is to match or exceed the ADR reporting rate from spontaneous reporting plus ICD-10 codes, whilst significantly reducing the burden of other ADE reports.



identify ADRs are required to rival the increasing number of drugs approved through expedited and provisional pathways [3, 16]. For these reasons, there is a need to create new systems to capture ADRs.

Considering that reviewing patient history is the most time-consuming process in maintaining a pharmacovigilance program, automated machine learning models have the potential to be trained to detect ADRs in real-time and flag potential high-risk cases for further review when required [19]. Furthermore, EMR data is generally inexpensive, accessible, can be obtained without interfering with patient care, and reflects real-world clinical outcomes [12]. In the same vein, machine learning models are currently being explored as methods for predicting ADRs in the preclinical stages of drug development [15].

1.1 Similar Work

In recent years there has been increasing interest in using natural language processing (NLP) for the detection of ADEs, a broader group of adverse events that encompasses ADRs, as well as poisonings and other medication-related harm. Several datasets have been developed for benchmarking algorithms, most notably the n2c2-2018 ADE [11] and the MADE 1.0 [14] datasets. Four distinct NLP paradigms have been used for ADE detection in these datasets, namely: 1) rule-based, 2) machine learning-based 3) deep learning-based, and 4) contextualized language model-based approaches. The best performance has been seen with deep learning methods [21], in particular using pre-trained large language models (LLMs) such as BERT [4] and BioBERT [20].

In contrast, ADR detection has not been as well studied. Distinguishing ADRs from other ADEs is a task that often requires extensive review of clinical notes to exclude both alternative causes and inappropriate medication usage. Attribution of causality hinges on scoring systems like the Naranjo Score [26]. One approach to ADR detection has been to predict the Naranjo Score from the clinical notes [29], however this assumes that the score components have been documented. When these components are not documented, that does not necessarily mean they were absent. This has the potential to introduce a further source of bias from differential documentation. An alternative approach to distinguishing ADEs from ADRs is to use a dataset enriched for both, so that our algorithms are forced to make this distinction.

1.2 Study Aim

The aim of this study was to develop an NLP algorithm to identify ADRs in discharge summaries and, in particular, to distinguish ADRs from other ADEs to augment ADR reporting (see Figure 1). In doing so, we developed a new corpus of ADR annotations, enriched for both ADRs and ADEs. Like the most successful ADE detection algorithms, we used a deep learning approach, in particular pre-trained LLMs. In order to adapt to the national (Australia) and institutional documentation practices and vocabulary, we further pre-trained our model on a large corpus of unannotated discharge summaries from our institution, prior to fine-tuning on the annotated corpus.

2 Methods

2.1 Institutional Setting

861 discharge summaries spanning a 5 year period (2015-2020) from the Electronic Medical Record (EMR) were collected and retrieved from the Clinical Research Data Warehouse of a 900-bed metropolitan tertiary teaching hospital network in Melbourne, Australia. All discharge summaries were from admissions coded with a Y40-59 ICD-10 code (ADE code), meaning that all admissions had been assessed by a clinical coder as containing an ADE – including a subset of 231 that had been previously reviewed and validated as true ADRs by our institutional ADR committee.

This model was governed by Austin Health Quality Improvement Number 41519.

2.2 Model Training and Evaluation

Figure 2 illustrates the full model training and evaluation process.

2.2.1 Annotation

The Prodigy annotation tool [1] was used to annotate drug names and known adverse drug reactions in the discharge summaries. We achieved rapid annotation of a moderate sized corpus by using active learning, an iterative process whereby annotation is performed in batches. On completion of each batch, an intermediate NLP model was retrained on the current annotations. This model was then used to suggest labels to the annotator, thus focusing the annotator’s attention on the most likely locations of interest within the text. The first set of 100 annotations were performed manually using the ner.manual Prodigy recipe to create an initial annotated corpus. The active learning process was then employed for all future batches (100-200 texts per batch), using the open-source Med7 model [17] as the initial intermediate NLP model. This ner.correct Prodigy recipe was used for all batches using the active learning process. This process was repeated until all documents were annotated. The two labels were ‘DRUG’ and

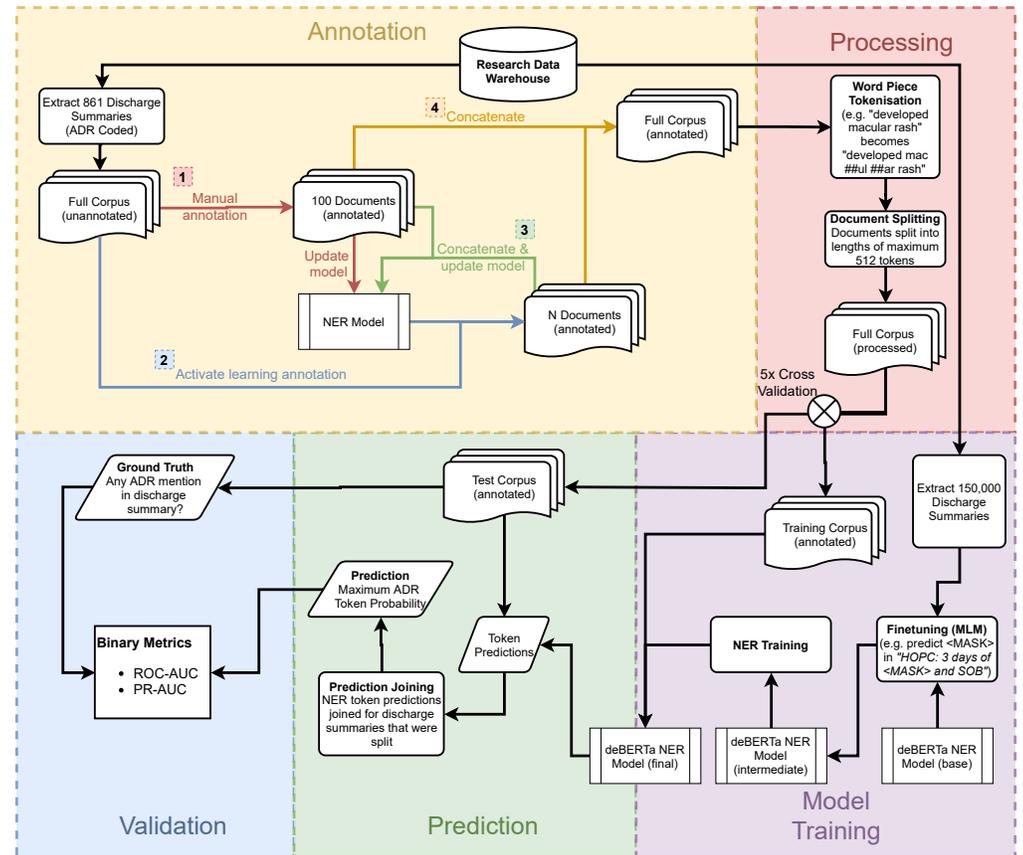


Figure 2. Full NLP Pipeline. Annotation is split into 4 steps performed in order – steps 2 & 3 are repeated in sequence with N discharge summaries each time, until either all data is annotated or the updated NER model shows minimal improvement when more annotated data is added.

‘ADR’, following a standardized label scheme (Appendix A) to create Named Entity Recognition (NER) training data for the model. A single annotator was trained to perform all annotations, which were then reviewed by a senior clinician.

2.2.2 Processing and Model Training

We used the DeBERTa model, a LLM neural network architecture that has been shown to outform BERT on many similar NLP tasks [10]. A pre-trained version is available from the HuggingFace Transformers library [33]. Pre-training with clinical texts has demonstrated improvement in other NLP tasks [2], so we therefore performed further pre-training of the DeBERTa model on masked-language modelling (MLM) [4], using a corpus of 150,000 unannotated discharge summaries from our EMR. In this step, 15% of the tokens in the unannotated corpus were masked, with the model trained to predict the mask from the surrounding text. We trained the model for 3 epochs on this task.

Annotated documents were tokenised using a word piece [34] tokeniser and then split into multiple chunks with a maximum token length of 512 tokens. These tokenised documents were then used for fine-tuning the model for the task of NER on drug and ADR annotations (see Figure 2). The maximum probability NER word labels were used to classify discharge summaries as ADR containing or not, and the final model was evaluated at the document level using the k-folds cross-validation method.

Model training steps were performed using Python version 3.9.6 [32] and the Transformers library [33] – using a PyTorch backend [28] – on 3 Nvidia 1080ti graphics processing units. Full details of model training, including hyperparameters, can be found in Appendix C. All code is available at <https://github.com/AustinMOS/adr-nlp>.

2.2.3 Validation and Testing

Using stratified k-fold cross validation, the annotated discharge summaries were shuffled randomly then split into 5 datasets. The model was then trained on 4 of the datasets and tested on the 5th (test) dataset. This allowed for the final model to be tested on previously unseen annotated EMR data not used in the training of the model. The process was repeated 5 times to calculate the mean and 95% confidence interval for each metric. The reported metrics are the area under the receiver operating characteristic curve (ROC-AUC) and the area under the precision-recall curve (PR-AUC).

2.2.4 Benchmarking

The performance of our model was compared to a previously published machine learning model from our institution [23]. This model was based on a dataset of Y40-59 ICD-10 coded admissions from December 2016 to November 2017. All of these admissions were flagged as possible ADRs and therefore assessed by an expert pharmacist using extensive chart review to determine the veracity of the reports. This model used only ICD-10 coding data and length of stay to discriminate between true and false ADR reports.

2.2.5 Ablation Study

We performed an ablation study to test the requirement for the intermediate pre-training step by comparing the final performance of the model with and without this step. Additionally, we examined the token-level F1-score during training to determine whether there was a trend towards “bridging” any gap in performance with additional training.

3 Results

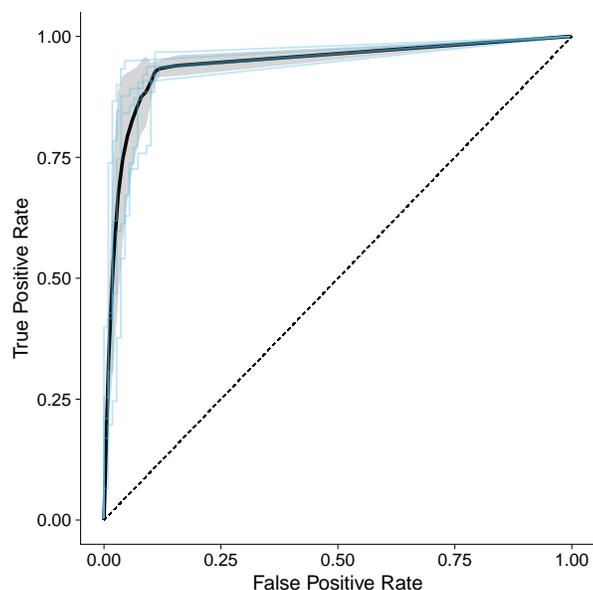
3.1 Data

In this study, 150,000 medical discharge summaries were used for intermediate pretraining of the model layered on top of the pre-existing DeBERTa transformer model. A further 861 annotated discharge summaries were used for the downstream NER task. All of the annotated discharge summaries came from admissions coded with an ICD-10 Y40.0-Y59.9, with 311 of the 861 containing an ADR mention. The admissions occurred across a broad range of specialities (23% general medicine, 10% geriatrics, 15% surgical).

3.2 Model

The model demonstrated good discriminative performance at the document level, with a ROC-AUC 0.934 (95% CI: 0.931 - 0.955) (see Figure 3) and PR-AUC 0.906 (95% CI: 0.885 - 0.926).

Figure 3. Receiver Operating Curves from 5-fold cross-validation. Mean performance is represented by the solid black line, with the grey shaded area representing the mean \pm standard deviation. Each individual fold is plotted in faint blue.



There was consistently good performance across all 5 folds with respect to the ROC (Figure 3). The final model did have the best performance – intermediate pretraining improved the PR-AUC by 0.12 and the ROC-AUC by 0.13, and the final pretrained NLP model had a higher ROC-AUC than the previously published ICD-10 model (see Table 1)). At the token level, the ablation study demonstrated an improvement in token-level F1-score with further pre-training that did not appear to improve with further model training (see Figure 5).

performance on one of the holdout folds demonstrates 93.1% classification accuracy (see Figure 4).

Using the Youden J-point threshold for classifying ADRs (30), the binary classification

Table 1. Precision-Recall Area Under Curve and Receiver Operating Characteristic Area Under Curve for ICD-10 model and Natural Language Processing model before and after pre-training.

Metric	NLP Model (no pre-training)	NLP Model (with pre-training)	ICD-10 Model
PR-AUC	0.894 (0.852 - 0.936)	0.906 (0.885-0.926)	-
ROC-AUC	0.921 (0.884 - 0.958)	0.934 (0.931-0.955)	0.803

Examining the 12 misclassified discharge summaries (see Appendix B) reveals distinct patterns. 5 out of the 7 false positives were adverse drug events, however they all fell short of our annotation criteria because of mild severity not warranting medication cessation. In contrast, 3 out of the 5 false negatives were anticoagulation-related adverse drug events where the anticoagulation was restarted at a lower dose – despite this, these events were annotated as ADRs because of their severity.

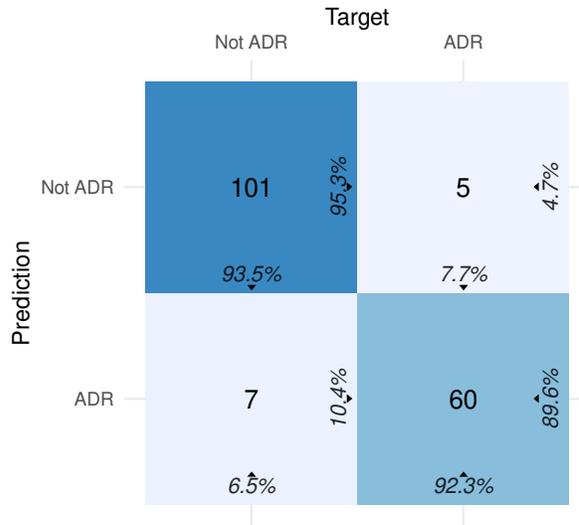


Figure 4. Confusion matrix from a single validation fold. This sample has the following metrics: sensitivity/recall 0.923, specificity 0.953, precision 0.896, F1 0.909, Matthew’s Correlation Coefficient 0.853.

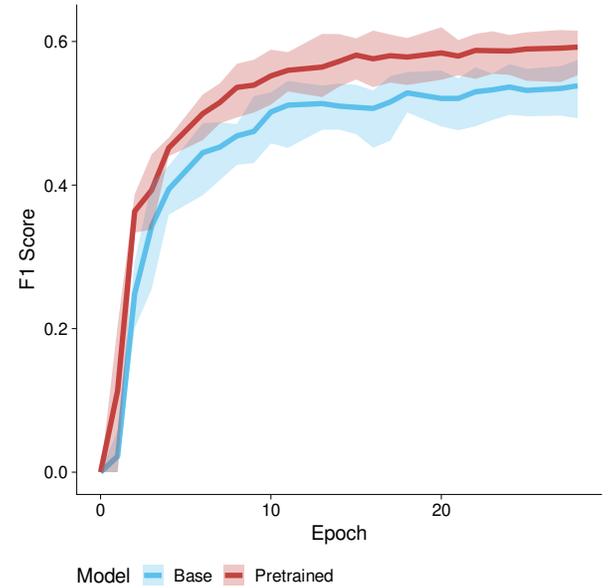


Figure 5. F1-score at the token level. The solid lines represent the mean F1 scores at each epoch (25 in total) and the shaded line shows the maximum and minimum values across all 5 folds.

Examining the NER output probabilities, the model correctly recognises appropriate entities like “eosinophilia” and “rash” as being common ADR features (see Figure 6). Attribution of these features to a drug increases the ADR class probabilities and this is robust to errors in spelling, even within important words (i.e., drug and ADR feature).

Exploring modifications to discharge summary text reveals that the algorithm is sensitive to medication mentions (see Table 2). Whilst it appears to be robust to non-Australian medication names like “acetaminophen”, trade names can result in misclassification. This likely reflects documentation practices in our institution, where generic names are much more commonly used. Removing the medication mention, or replacing it with an attribution to a disease or generic term like “a medication” greatly reduces the predicted ADR probability. In comparing these predictions to those produced by the model without intermediate pre-training, the pre-trained model produces robust predictions even when the ADR and causative drug mention are separated, as they may be in dot-point documentation. The model without pre-training requires these entities to be closely related in the text.

[CLS] Rash and eosinophilia likely secondary to new diagnosis of EG PA [SEP]
 [CLS] Rash and eosinophilia likely secondary to all op ur in ol [SEP]
 [CLS] Rash and oesinophillia likely secondary to al op ur in ol [SEP]

Figure 6. Examples of token-level ADR entity predictions. The color gradient goes from red (very low probability) to green (very high probability), representing the probability that each token belongs to the ADR class.

Description	Text	ADR Probability	
		With Pre-training	Without Pre-training
Baseline	# Pancreatitis - Lipase: 535 -> 154 -> 145 - Managed with NBM, IV fluids - CT AP and abdo USS: normal - Likely secondary to Azathioprine - ceased, never to be used again. - Resolved with conservative measures	0.837	0.425
Medication line removed	# Pancreatitis - Lipase: 535 -> 154 -> 145 - Managed with NBM, IV fluids - CT AP and abdo USS: normal - Resolved with conservative measures	0.003	0.428
Medication changed	# Pancreatitis - Lipase: 535 -> 154 -> 145 - Managed with NBM, IV fluids - CT AP and abdo USS: normal - Likely secondary to Methotrexate - ceased, never to be used again. - Resolved with conservative measures	0.967	0.424
No specific medication	# Pancreatitis - Lipase: 535 -> 154 -> 145 - Managed with NBM, IV fluids - CT AP and abdo USS: normal - Likely secondary to a medication - ceased, never to be used again. - Resolved with conservative measures	0.009	0.426
Non-Australian medication name	# Pancreatitis - Lipase: 535 -> 154 -> 145 - Managed with NBM, IV fluids - CT AP and abdo USS: normal - Likely secondary to Acetaminophen - ceased, never to be used again. - Resolved with conservative measures	0.842	0.425
Trade name	# Pancreatitis - Lipase: 535 -> 154 -> 145 - Managed with NBM, IV fluids - CT AP and abdo USS: normal - Likely secondary to Imuran - ceased, never to be used again. - Resolved with conservative measures	0.109	0.426
Non-medication cause	# Pancreatitis - Lipase: 535 -> 154 -> 145 - Managed with NBM, IV fluids - CT AP and abdo USS: normal - Likely secondary to alcohol - ceased, never to be used again. - Resolved with conservative measures	0.015	0.426
Single sentence	# Pancreatitis secondary to azathioprine	0.999	0.997

Table 2. Text-level ADR Predictions. Comparison with and without the additional MLM pre-training step. Output probabilities are scaled from 0 to 1.

4 Discussion

We trained a machine learning model on EMR discharge summaries using natural language processing parameters to detect drugs and adverse drug reactions. Compared with a previously published machine learning model [23] using ICD-10 coding data, our study demonstrates an improvement in differentiating true from false ADR reports using NLP. NLP has the benefit of bypassing the intermediate step of generating clinical codes from clinical notes, instead deriving predictions directly from the primary source. This can reduce the inconsistencies, missed ADRs, and human errors which may arise with the use of clinical coding for ADRs. By focusing on the language of clinical data for the identification of ADRs, we return to the fundamentals of ADR detection – the identification of a putatively causal relationship between the administration or ingestion of a drug and a subsequent adverse response. Although not explored in this study, there is a large potential for future studies to build upon this framework for identifying drugs and ADRs in novel clinical data and train models to recognise temporal relationships between these pairs.

With the increasing transition to recording patient care in electronic formats, machine learning models have the potential to scan large amounts of data for ADRs and improve the accuracy and efficiency of pharmacovigilance within a health network. NLP is an effective method for processing electronic data into structured forms for clinical research. It has been demonstrated to have an accuracy comparable to professional clinical coders in the coding of radiographic reports and is superior to simple text searching methods [5, 13]. The broad applicability of NLP systems is demonstrated by their ability to be extended to recognize new patterns and types of information [6]. The MedLEE NLP system was initially trained for the automated processing of radiological reports, but it has also been successfully extended to detect adverse events in discharge summaries [7, 24]. In a similar vein, our model demonstrates how the open-source DeBERTa model may be extended to develop a machine learning model for ADR prediction. Furthermore, the additional step of pre-training on a larger corpus of discharge summaries from our institutional EMR lead to improved model accuracy. This may be because the model has learnt some of the linguistic features that are unique to our regional (Australia) and local (institution) context.

In recent years there has been increasing interest in using NLP models for the detection of ADEs. The open NLP challenge for ADE detection was held in 2018, using the MADE 1.0 corpus of annotated clinical notes related to drug safety surveillance [2]. However, these datasets have not been developed specifically to identify ADRs and distinguish them from other ADEs. This is especially important given the high rates of ADR under-reporting and their relatively rare occurrence in clinical notes. Therefore, a major strength of our model was the use of real-world clinical ADR reports in the annotation of our dataset. Our model was trained on 861 Y40.0-Y59.9 ICD-coded discharge summaries, of which 311 were labelled as ADR reports, with a significant subset (231 admissions) having been validated as true ADRs by our institutional ADR committee. This resulted in a smaller class imbalance than would have been observed had we used an unenriched sample of discharge summaries.

A key limitation of any NLP algorithm is the quality of the training data provided. For instance, a model will perform poorly where there are inconsistent annotations. A labelling scheme in our study was therefore used to reduce ambiguity and guide the manual annotation process for the training of our model. Although all annotations were reviewed by a senior clinician, the inherent biases of a single annotator are also carried forward into the biases of the final model and its prediction of ADR and drug labels. The misclassified labels we examined did demonstrate inconsistencies in labelling, where 8 of 12 misclassified discharge summaries were cases in which the offending medication was restarted – some labelled as ADRs and some not. Whilst these are not ADRs for

the purposes of a patient drug allergy history, when they fall on the severe end of the spectrum they are important to recognise for the purpose of population-level pharmacovigilance. More prescriptive annotation criteria might overcome this problem, however it might also mean that important severe events are missed. One way to overcome this might be by making event severity its own annotation and combining ADR and severity predictions to ensure severe events with less certainty about causality are captured for manual review.

Annotating large datasets for NLP can be time-consuming and costly. The ability to quickly and accurately annotate a corpus of this size was accelerated using active learning, allowing full review of each document whilst reducing the cost of producing further annotations with each subsequent batch. Although we did not reach a point in which we had saturated model improvement (4% improvement in token-level accuracy when using 100% vs. 75% of the corpus), increasing the corpus size beyond 861 is unlikely to result in significant improvements alone.

In terms of our dataset, our model was trained solely on ICD-10 coded Y40-Y59 discharge summaries rather than the full spectrum of EMR data. In order to be included in the training of the model, ICD-10 coding relies on the clinicians and clinical coders to have documented and identified an event that was thought to be an ADE. This may have missed some important confounding features in discharge summaries without this feature. Beyond discharge summaries, expanding the dataset to encompass other EMR notes (e.g. inpatient notes, pathology reports etc.) is another direction for future study. Additional improvements may be made to this model by annotating and training the model to identify linked concepts, such as drug indication and dose, separate annotations for non-ADR ADEs, and further training of the model to recognize drug-ADR and drug-dose-indication relations.

Whilst there is often consistency in the formatting of discharge summaries between health networks, it will require validation on these datasets to demonstrate generalizability. In particular, it is likely that we will observe a decrease in performance in institutions where trade name documentation is more prevalent. We anticipate that this problem can be prevented with synthetic training data generated from our current dataset, by simply replacing generic names with trade names from a thesaurus of generic-trade name pairs.

5 Conclusion

Our study demonstrates the potential for NLP models to be developed for automated ADR detection. This approach can address the under-reporting issues of current methods, bypass the resource limitations of current clinical workflows and increase the ADR reporting rates within the hospital. With additional pre-training on EMR data specific to our health network, the model was able to learn the patterns of discharge summary formatting, allowing correct classification even of distant relations when documented within the expected structure of a discharge summary. The unique construct of our corpus, particularly the presence of many validated ADRs alongside non-ADR ADEs, meant that our model had to differentiate ADRs from other incidents of medication-related harm. We plan to implement our model into the clinical workflow of ADR reporting. Specifically, we are working on an ADR dashboard to present ADR reports, derived from both spontaneous reporting and our NLP model. Reports will be presented with the relevant section of the discharge summary highlighted according to the NER outputs of the model. These annotations can be confirmed or rejected by the pharmacist reviewing the report, providing ongoing annotations to further train and refine our model. We expect to identify any dataset drift by this method, including novel ADRs, new medications and new documentation practices.

References

1. Prodigy · an annotation tool for AI, machine learning & NLP. <https://prodi.gy/>. Accessed: 2021-12-1.
2. E. Alsentzer, J. R. Murphy, W. Boag, W.-H. Weng, D. Jin, T. Naumann, and M. B. A. McDermott. Publicly available clinical BERT embeddings. Apr. 2019.
3. Australian Government Department of Health. Therapeutic Goods Administration. Provisional approval pathway: prescription medicines. <https://www.tga.gov.au/provisional-approval-pathway-prescription-medicines>. Accessed: 2021-12-1.
4. 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, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
5. M. Fiszman, W. W. Chapman, D. Aronsky, R. S. Evans, and P. J. Haug. Automatic detection of acute bacterial pneumonia from chest x-ray reports. *Journal of the American Medical Informatics Association: JAMIA*, 7(6):593–604, Nov. 2000.
6. C. Friedman. Towards a comprehensive medical language processing system: methods and issues. *Proceedings : a conference of the American Medical Informatics Association. AMIA Fall Symposium*, pages 595–599, 1997.
7. C. Friedman, G. Hripcsak, W. DuMouchel, S. B. Johnson, and P. D. Clayton. Natural language processing in an operational clinical information system. *Natural Language Engineering*, 1(1):83–108, Mar. 1995.
8. M. Hacker, W. S. Messer, and K. A. Bachmann. *Pharmacology: Principles and Practice*. Academic Press, June 2009.
9. L. Hazell and S. A. W. Shakir. Under-reporting of adverse drug reactions : a systematic review. *Drug safety: an international journal of medical toxicology and drug experience*, 29(5):385–396, 2006.
10. P. He, X. Liu, J. Gao, and W. Chen. DeBERTa: Decoding-enhanced BERT with disentangled attention. June 2020.
11. S. Henry, K. Buchan, M. Filannino, A. Stubbs, and O. Uzuner. 2018 n2c2 shared task on adverse drug events and medication extraction in electronic health records. *Journal of the American Medical Informatics Association: JAMIA*, 27(1):3–12, Jan. 2020.
12. C. M. Hohl, A. Karpov, L. Reddekopp, M. Doyle-Waters, and J. Stausberg. ICD-10 codes used to identify adverse drug events in administrative data: a systematic review. *Journal of the American Medical Informatics Association: JAMIA*, 21(3):547–557, May 2014.
13. G. Hripcsak, S. Bakken, P. D. Stetson, and V. L. Patel. Mining complex clinical data for patient safety research: a framework for event discovery. *Journal of biomedical informatics*, 36(1-2):120–130, Feb. 2003.

14. A. Jagannatha, F. Liu, W. Liu, and H. Yu. Overview of the first natural language processing challenge for extracting medication, indication, and adverse drug events from electronic health record notes (MADE 1.0). *Drug safety: an international journal of medical toxicology and drug experience*, 42(1):99–111, Jan. 2019.
15. S. Jamal, S. Goyal, A. Shanker, and A. Grover. Predicting neurological adverse drug reactions based on biological, chemical and phenotypic properties of drugs using machine learning models. *Scientific reports*, 7(1):872, Apr. 2017.
16. A. S. Kesselheim, B. Wang, J. M. Franklin, and J. J. Darrow. Trends in utilization of FDA expedited drug development and approval programs, 1987-2014: cohort study. *BMJ*, 351:h4633, Sept. 2015.
17. A. Kormilitzin, N. Vaci, Q. Liu, and A. Nevado-Holgado. Med7: a transferable clinical natural language processing model for electronic health records. Mar. 2020.
18. L. A. Ladewski, S. M. Belknap, J. R. Nebeker, O. Sartor, E. A. Lyons, T. C. Kuzel, M. S. Tallman, D. W. Raisch, A. R. Auerbach, G. T. Schumock, H. C. Kwaan, and C. L. Bennett. Dissemination of information on potentially fatal adverse drug reactions for cancer drugs from 2000 to 2002: first results from the research on adverse drug events and reports project. *Journal of clinical oncology: official journal of the American Society of Clinical Oncology*, 21(20):3859–3866, Oct. 2003.
19. A. Lavertu, B. Vora, K. M. Giacomini, R. Altman, and S. Rensi. A new era in pharmacovigilance: Toward real-world data and digital monitoring. *Clinical pharmacology and therapeutics*, 109(5):1197–1202, May 2021.
20. 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(4):1234–1240, Feb. 2020.
21. D. Mahendran and B. T. McInnes. Extracting adverse drug events from clinical notes. *AMIA Joint Summits on Translational Science proceedings. AMIA Joint Summits on Translational Science*, 2021:420–429, May 2021.
22. G. B. McLachlan, C. Keith, and C. Wood. The cost of pharmacovigilance: a time and motion study of an adverse drug reaction program. *The International journal of pharmacy practice*, 29(5):521–523, Oct. 2021.
23. C. McMaster, D. Liew, C. Keith, P. Aminian, and A. Frauman. A Machine-Learning algorithm to optimise automated adverse drug reaction detection from clinical coding. *Drug safety: an international journal of medical toxicology and drug experience*, 42(6):721–725, June 2019.
24. G. B. Melton and G. Hripcsak. Automated detection of adverse events using natural language processing of discharge summaries. *Journal of the American Medical Informatics Association: JAMIA*, 12(4):448–457, July 2005.
25. J.-L. Montastruc, M. Lafaurie, C. de Canecaude, G. Durrieu, A. Sommet, F. Montastruc, and H. Bagheri. Fatal adverse drug reactions: A worldwide perspective in the world health organization pharmacovigilance database. *British journal of clinical pharmacology*, 87(11):4334–4340, Nov. 2021.

26. C. A. Naranjo, U. Busto, E. M. Sellers, P. Sandor, I. Ruiz, E. A. Roberts, E. Janecek, C. Domecq, and D. J. Greenblatt. A method for estimating the probability of adverse drug reactions. *Clinical pharmacology and therapeutics*, 30(2):239–245, Aug. 1981.
27. J. R. Nebeker, P. Barach, and M. H. Samore. Clarifying adverse drug events: a clinician’s guide to terminology, documentation, and reporting. *Annals of internal medicine*, 140(10):795–801, May 2004.
28. A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative style, High-Performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. Aché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
29. B. P. S. Rawat, A. Jagannatha, F. Liu, and H. Yu. Inferring ADR causality by predicting the naranjo score from clinical notes. *AMIA ... Annual Symposium proceedings / AMIA Symposium. AMIA Symposium*, 2020:1041–1049, 2020.
30. W. B. Runciman, E. E. Roughead, S. J. Semple, and R. J. Adams. Adverse drug events and medication errors in australia. *International Journal for Quality in Health Care*, 15 Suppl 1:i49–59, Dec. 2003.
31. E. J. Thomas, D. M. Studdert, W. B. Runciman, R. K. Webb, E. J. Sexton, R. M. Wilson, R. W. Gibberd, B. T. Harrison, and T. A. Brennan. A comparison of iatrogenic injury studies in australia and the USA. i: Context, methods, casemix, population, patient and hospital characteristics. *International journal for quality in health care: journal of the International Society for Quality in Health Care / ISQua*, 12(5):371–378, Oct. 2000.
32. G. van Rossum. Python tutorial. Technical report, Centrum voor Wiskunde en Informatica (CWI), Amsterdam, May 1995.
33. T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush. HuggingFace’s transformers: State-of-the-art natural language processing. Oct. 2019.
34. Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, and J. Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation. Sept. 2016.

Appendix

Appendix A: NER Annotation Label Scheme

1. Only label drugs that are possibly responsible for an ADR
2. ‘Drug’ labels are only given for drugs administered at therapeutic dose.

3. 'Drug' labels can be given if it is a clear cause of an ADR even if it is unclear what ADR occurred.
4. Preference labelling of specific drugs. Drug class can be labelled if specific causative drug is not mentioned (e.g. "angioedema after starting ACE-inhibitor")
5. ADR labels are only assigned to disease states, signs and symptoms, not pathological/radiological findings with no clinical consequence.
6. Abbreviations for drugs and ADRs are permitted
7. Preference drug and ADR labels which are in close proximity within the document
8. Preference labels with causal language – e.g. 'due to', 'secondary to', 'withheld', 'ceased', 'complicated by'.
9. Multiple ADRs may be labelled for a given drug
10. Multiple drugs may be labelled as (potentially) responsible for a given ADR
11. Do not label drugs and ADRs already recorded in the patient allergy list
12. Only label drugs documented in the body of the text, not in a drug list

Appendix B: Model Errors

False Negatives

- 1 instance of interstitial nephritis lacking drug attribution in the body of the text (likely association with NSAIDs only mentioned on the discharge plan at the very end of the document)
- 3 instances of anticoagulation being restarted after a bleed – all could conceivably be considered adverse drug events and not ADRs
- 1 instance of transaminitis from statin (only instance of statins being implicated in dataset, only instance of transaminitis)
 - Transaminitis still assigned probability of 41

False Positives

- Elevated transaminases in the setting of therapy for Mycobacterium tuberculosis
 - Could reasonably have been labelled as an ADR, but not done so because self-resolved without any changes to therapy
- Lithium toxicity already recognised prior to admission with plan to restart at a later date
 - Given plan to restart, not an ADR that could be placed on the patient's record, but important to note
- Past history of chemotherapy-induced pancytopenia, not an ADR relevant to current admission
- Multifactorial anaemia - partially attributed to anticoagulation (not ceased)
- Headache requiring lumbar puncture - anticoagulation withheld for the procedure

- Peripheral oedema secondary to amlodipine (not ceased) - adverse drug event that could conceivably be labelled as an ADR
- Possible medication-induced bradycardia vs. sick sinus syndrome - decision made to continue beta-blocker and tolerate slower rate – likely medication-induced, but not rising to the severity of an ADR

Appendix C: Model Training Details

Pre-training

- Learning rate: $5e-5$
- Weight decay: 0.01
- Training duration: 3 epochs (12 hours)

Fine-tuning

- Learning rate: $9e-6$
- Weight decay: 0.01
- Training duration: 25 epochs