A Comprehensive Study of Gender Bias in Chemical Named Entity Recognition Models

Xingmeng Zhao, Ali Niazi, and Anthony Rios

Department of Information Systems and Cyber Security The University of Texas at San Antonio

{xingmeng.zhao, ali.niazi, anthony.rios}@utsa.edu

Abstract

Chemical named entity recognition (NER) models are used in many downstream tasks, from adverse drug reaction identification to pharmacoepidemiology. However, it is unknown whether these models work the same for everyone. Performance disparities can potentially cause harm rather than the intended good. This paper assesses gender-related performance disparities in chemical NER systems. We develop a framework for measuring gender bias in chemical NER models using synthetic data and a newly annotated corpus of over 92,405 words with self-identified gender information from Reddit. Our evaluation of multiple biomedical NER models reveals evident biases. For instance, synthetic data suggests that female names are frequently misclassified as chemicals, especially when it comes to brand name mentions. Additionally, we observe performance disparities between female- and maleassociated data in both datasets. Many systems fail to detect contraceptives such as birth control. Our findings emphasize the biases in chemical NER models, urging practitioners to account for these biases in downstream applications.

1 INTRODUCTION

Chemical named entity recognition (NER) is the extraction of chemical mentions (e.g., drug names) from the text. Chemical NER is essential in many downstream tasks, from pharmacovigilance (O'Connor et al., 2014) to facilitating drug discovery by mining biomedical research articles (Agarwal and Searls, 2008). For instance, Chemical NER systems are the first step in pipelines developed to mine adverse drug reactions (ADRs) (Farrugia and Abela, 2020; Mammì et al., 2013). However, it is unknown whether these systems perform the same for everyone. Who benefits from these systems, and who can be harmed? In this paper, we present a comprehensive analysis of

gender-related performance disparities of Chemical NER Systems.

Performance disparities have recently received substantial attention in the field of NLP. For example, there are differences in text classification models across sub-populations such as gender, race, and minority dialects (Dixon et al., 2018; Park et al., 2018; Badjatiya et al., 2019; Rios, 2020; Lwowski and Rios, 2021; Mozafari et al., 2020). Performance disparities can manifest in multiple parts of NLP systems, including the pre-trained models (e.g., word embeddings) and their downstream applications (Zhao et al., 2019; Goldfarb-Tarrant et al., 2021; Zhao et al., 2017). While previous research has explored these disparities for NER systems, the focus has been largely on synthetic data and non-biomedical NER applications (Mehrabi et al., 2020). Our study addresses this gap by providing a comprehensive examination of genderrelated performance disparities in Chemical NER, focusing on both synthetic and real-world data.

This paper is most similar to Mehrabi et al. (2020) with two primary distinctions. First, our focus is on Chemical NER, a less studied area in Biomedical NLP despite its having major bias implications. Second, while Mehrabi et al. (2020) uses synthetic data and templates (e.g., NAME in LOCATION) for bias analysis, we delve deeper into the potential including an analysis of the interaction of morphology patterns on bias. For instance, Lieven et al. (2015) highlighted a preference for linguistically feminine brand names in the market, leading drug companies to adopt such naming conventions. These patterns in training data can inadvertently cause models to misclassify female names as chemicals.

We also examine real-world data looking at the performance of chemical NER systems on groups that identify as male or female. For instance, Sundbom et al. (2017) shows that women are more frequently prescribed antidepressants than men. Other

studies, like Riley III et al. (1998), reveal gender differences in pain sensitivity and opioid prescriptions, with women receiving opioids twice as often. If chemical NER models struggle to detect the drugs often mentioned, then it may cause genderspecific biases in their performance. Our analysis identifies some of these patterns in real data.

Overall, this paper presents a dual approach: we explore template data but also assemble and annotate a novel real-world dataset with self-identified gender information. ¹ Synthetic data allows us to target specific biases in the models (e.g., morphological issues). Likewise, we believe exploring data from people who have self-identified their demographic information will provide a more realistic understanding of how these models will perform based on how people write and what they write about

Our main contributions are two-fold:

- We introduce a novel annotated Chemical NER dataset for social media data. Moreover, the dataset contains self-identified gender information to be used to measure gender bias in Chemical NER models. To the best of our knowledge, this is the first Reddit-based Chemical NER dataset, and it is the first Chemical NER dataset with self-identified gender information.
- 2. We provide a comprehensive testing framework for gender bias in Chemical NER using both synthetic and real-world data. To the best of our knowledge, our results are the first to conduct bias analysis for chemical NER models. This allows a better understanding of modern chemical NER techniques.

2 RELATED WORK

Prior work extensively curated labeled data for chemical NER and developed domain-specific models. For example, the CHEMDNER corpus (Krallinger et al., 2015) was created for the 2014 BioCreative shared task on chemical extraction from text. Researchers recognize the importance of these systems and are working to make them as fair and accurate as possible. Likewise, the CDR (Li et al., 2016) dataset was developed to detect chemical-disease relations for the 2015 shared task. Similar to traditional NER tasks (Li et al., 2020), a broad range of approaches have

been proposed to detect chemicals (Rocktäschel et al., 2012; Chiu et al., 2021; Lee et al., 2020; Sun et al., 2021; López-Úbeda et al., 2021; Weber et al., 2021), from traditional conditional random fields to deep learning methods. Many recent neural network-based advances can be broken into three main groups of models, word, character, and contextual embedding-based models. For instance, Lee et al. (2020) trained a biomedical-specific BERT model that improved on many prior state-of-the-art results. HunFlair (Weber et al., 2021) introduced a method that matches the word, contextual, and character embeddings into a unified framework to achieve state-of-the-art performance. In this paper, we evaluate several state-of-the-art systems. Particularly, we focus on systems that use word embeddings, sub-word embeddings, and character embeddings, which allows us to understand the impact of morphological features of the chemical names on gender bias.

Several previous works have measured and highlighted bias in different NLP tasks. For instance, Sap et al. (2019) measures the bias of offensive language detection models on African American English. Likewise, Park et al. (2018) measures gender bias of abusive language detection models and evaluates various methods such as word embedding debiasing and data augmentation to improve biased methods. Davidson et al. (2019) shows racial and ethnic bias when identifying hate speech online and that tweets in the black-aligned corpus are more likely to be assigned hate speech. Gaut et al. (2020) creates the WikiGenderBias dataset to evaluate the gender bias in the relation extraction (RE) model, confirming that the RE system behaves differently when the target entities are of different genders. Cirillo et al. (2020) demonstrate that biases in biomedical applications can stem from various sources, such as skewed diagnoses resulting from clinical depression scales that measure symptoms more prevalent in women, potentially leading to a higher reported incidence of depression among this group (Martin et al., 2013). Other sources include the underrepresentation of minority populations such as pregnant women (Organization and for Women's Health in Society, 2009), non-representative samples in AI training data, and inherent algorithmic discrimination, all potentially contributing to inaccurate and unfair results.

Recent research has shown that although Large Language Models (LLMs) are now increasingly

¹The dataset and datasheet are available at https://zenodo.org/records/10905462

being used for tasks such as Named Entity Recognition (Ashok and Lipton, 2023; Wang et al., 2023) and relation classification (Wan et al., 2023), they also have the potential to reinforce or exacerbate gender biases, which emphasizes the importance of careful deployment to prevent the reinforcement of stereotypes (Kotek et al., 2023).

Overall, several metrics have been proposed to measure gender bias. One of the most commonly used metrics involves measuring bias by examining model performance disparities on male and female data points (Kiritchenko and Mohammad, 2018). Performance disparities have been observed across a wide array of NLP tasks such as detecting virus-related text (Lwowski and Rios, 2021), language generation (Sheng et al., 2019), coreference resolution (Zhao et al., 2018), named entity recognition (Mehrabi et al., 2020), and machine translation (Font and Costa-jussà, 2019). Most related to this study, researchers have shown that traditional NER systems (i.e., to detect people, locations, and organizations) are biased concerning gender (Mehrabi et al., 2020). Specifically, Mehrabi et al. (2020) demonstrates that female names are more likely to be misidentified as a location than male names. This stream of research underscores the importance of our investigation into performance disparities in NLP.

Finally, while not directly studied in prior NER experiments, it is important to discuss some background about morphological elements of chemical names. Morphological elements often representing masculinity or femininity are frequently used in chemical naming conventions. According to Lieven et al. (2015), consumers perceive linguistically feminine brand names as warmer and more likable. For instance, adding a diminutive suffix to the masculine form of the name usually feminizes it. The masculine names such as Robert, Julius, Antonio, and Carolus (more commonly Charles today) are feminized by adding the suffixes "a", "ia", "ina", or "ine" to generate Roberta, Julia, Antonia, and Caroline, respectively. The suffixes "ia" and "a" is commonly used for inorganic oxides such as magnesia, zirconia, silica, and titania (Hepler-Smith, 2015). Likewise, "ine" is used as the suffix in many organic bases and base substances such as quinine, morphine, guanidine, xanthine, pyrimidine, and pyridine. Hence, while these practices were not originally "biased" in their original usage, they can potentially impact model performance

	# of Chems.	# Sentences	# Words
CDR	4,409	14,306	346,001
CHEMDNER	84,355	87,125	2,431,247
CHEBI	24,121	12,913	423,577
AskDoc MALE	1,501	2,862	52,221
AskDoc FEMALE	1,774	2,151	40,184
AskDoc ALL	3,275	5,013	92,405
Synthetic MALE	2,800,000	2,800,000	25,760,000
Synthetic FEMALE	2,800,000	2,800,000	25,760,000
Synthetic ALL	5,600,000	5,600,000	51,520,000

Table 1: Dataset statistics.

(e.g., feminine names can be detected as chemicals). Therefore, the patterns can cause biased models. As part of our approach to investigate this potential source of bias, we propose using synthetic data to quantify this phenomenon.

3 DATASETS

We use five main datasets used in our experiments: three are publicly-released datasets based on PubMed (CDR (Li et al., 2016), CHEMD-NER (Krallinger et al., 2015), and CHEBI (Shardlow et al., 2018)) and two are newly curated datasets, one using social media data and another based on templates. Table 1 provides their statistics. We selected the PubMed datasets for their prominence in chemical NER research. At the same time, the r/AskDocs subreddit was chosen for its large community, diverse health discussions, and consistent gender identification format, such as "I [25 M]". We provide complete descriptions of the publicly-released datasets in the Appendix. In this section, focus on the description of the newly collected and annotated data.

Synthetic (Template) Data We designed a new synthetic dataset to quantify the gender bias in the Chemical NER models. Intuitively, the purpose of the synthetic dataset is to measure two items. First, do gender-related names and pronouns get incorrectly classified as chemicals (i.e., cause false positives)? Second, does the appearance of gender-related names/pronouns impact the prediction of other words (i.e., cause false negatives)? Specifically, we create templates such as "[NAME] said they have been taking [CHEMICAL] for an illness." In the "[NAME]" column, we filled in the names associated with the male and female genders based on the 200 most popular baby names provided by the Social Security Administration ². Hence, we

²https://www.ssa.gov/oact/babynames/

Templates

[NAME] said they have been taking [CHEMICAL] for an illness. Did you hear that [NAME] has been using [CHEMICAL]. [CHEMICAL] has really been harming [NAME], I hope they stop. I think [NAME] is addicted to [CHEMICAL].

[NAME], please stop taking [CHEMICAL], it is bad for you.

Table 2: Templates used to create the synthetic dataset.

refer to these "gender-related" names in this paper. We recognize that gender is not binary and that names do not equal gender. We also recognize that the names do not accurately capture immigrants. This is a similar framework used by Mishra et al. (2020) and other gender bias papers (Kiritchenko and Mohammad, 2018). The "[CHEMICAL]" field is filled with the chemicals listed in the Unified Medical Language System (UMLS) (Bodenreider, 2004). For example, completed templates include "John said they have been taking citalopram for illness." and "Karen said they have been taking citalopram for illness." We created examples using five templates, 200 chemicals, and 200 names for each gender for each decade from 1880 to 2010, generating a total of 200,000 templates for each of the 14 decades. A list of additional templates is shown in Table 2. This dataset is only used for evaluation.

AskDocs We develop a new corpus using data from the Reddit community r/AskDocs. r/AskDocs provides a platform for peer-to-peer and patientprovider interactions on social media to ask medical-related questions. The providers are generally verified medical professionals. We collected all the posts from the community with self-identified gender mentions. To identify self-identified gender, we use a simple regular expression that looks for mentions of "I" or "My" followed by gender, and optionally age, e.g., "I [F34]", "My (23F)", "I [M]". Next, following general annotation recommendations for NLP (Pustejovsky and Stubbs, 2012), the annotation process was completed in two stages to increase the reliability of the labels. First, two graduate students annotated chemicals in the dataset resulting in an inter-annotator agreement of .874, achieving a similar agreement score as CDR and CHEMDNER. Second, a graduate student manually reviewed all disagreeing items to adjudicate the label and generate the gold standard. All students followed the same annotation guidelines developed for the CHEMDNER corpus. Contrary to the synthetic dataset, the actual data will allow users to measure biases arising from text content differences across posts with different self-identified gender mentions.

4 Methods

The goal of NER is to classify words into a sequence of labels. Formally, given an input sequence $\mathcal{X} = [x_1, x_2, \dots, x_N]$ with N tokens, the goal of NER is to output the corresponding label sequence $\mathcal{Y} = [y_1, y_2, \dots, y_N]$ with the same length, thus modeling the probabilities over a sequence $p(\mathcal{Y}|\mathcal{X})$. For this task, we conducted an experiment evaluating out-of-domain models on the AskDoc corpus. Specifically, models were trained and optimized on the CHEMDNER and CDR datasets and then applied to the AskDoc dataset. All models are evaluated using precision, recall, and F1. To measure bias, we use precision, recall, and F1 differences (Czarnowska et al., 2021). Specifically, let m be Males' performance metric (e.g., F1), and frepresent the Female metric. The bias is measured using the difference f - m.

4.1 MODELS

We evaluate three distinct models: Word Embedding models (Mikolov et al., 2013b), Flair embedding models (Akbik et al., 2018), and BERT-based models (Devlin et al., 2019a). While the embeddings for each model type vary, the sequence processing component is the same for each method. Specifically, following best practices for state-ofthe-art NER models (Akbik et al., 2019a), we use a Bidirectional long short-term memory network (Bi-LSTM) (Hochreiter and Schmidhuber, 1997) due to its sequential characteristics and capability to capture long-term dependencies. Recent research has shown that Bi-LSTM models can produce state-ofthe-art performance when combined with contextual embeddings and Conditional Random Fields (CRFs) (Mueller et al., 2020; Veyseh et al., 2022). Hence, in this paper, we use the Bi-LSTM+CRF implementation in the Flair NLP framework (Akbik et al., 2019a). The Bi-LSTM+CRF model is flexible because it can accept arbitrary embeddings as input. It is not constrained to traditional word embeddings (e.g., Word2Vec).

4.2 EMBEDDINGS

We explore three sets of embeddings: Word2Vec, Flair, and BERT. For all embeddings, we ex-

periment with domain-specific (e.g., trained on PubMed) and general embeddings (e.g., Google News corpus). We chose these three embedding types because they cover word, subword, and character-level embedding methods. Social media texts are brief and informal. Drugs and chemicals are typically described in descriptive, nontechnical language with spelling errors. These issues are challenging for chemical NER models trained on social media data. Moreover, some medications, like "all-trans-retinoic acid", contain morphologically difficult parts. Yet, similar-structured phrases still generally represent similar things (Zhang et al., 2021). How we represent words can directly impact performance and bias. We describe each embedding we use below:

Word2Vec. We use Word2Vec domain-specific embeddings pre-trained on PubMed and PubMed Central (Pyysalo et al., 2013) and general embeddings trained on the Google News corpus (Mikolov et al., 2013a). The embeddings are publicly released as part of the FLAIR package. It is important to state that word embeddings have a major limitation. Word embeddings use a distinct vector to represent each word and ignore words' internal structure (morphology). This can result in models not particularly good at learning rare or out-of-vocabulary (OOV) words in the data. The growing number of emerging chemicals/drugs with diverse morphological forms makes recognizing chemical entities on social media platforms particularly challenging. Another challenge posed by user-generated content is its unique characteristics and use of informal language, typically short context, noisy, sparse, and ambiguous content. Hence, we hypothesize that word embeddings would perform worse than other methods. However, it is unclear how these differences can impact bias.

Flair/HunFlair. Weber et al. (2021) and Akbik et al. (2019b) recently proposed a Flair contextual string embeddings (a character-level language model). Specifically, we use two versions of the embeddings in the HunFlair extension of the Flair package (Weber et al., 2021). The domain-specific embeddings are pre-trained on a corpus of three million full-text articles from the Pubmed Central BioC text mining collection (Comeau et al., 2019) and about twenty-five million abstracts from PubMed. The general embeddings are trained on a one billion word news corpus (Akbik et al., 2019b).

Unlike word embeddings mentioned above, Flair embeddings are a contextualized character-level representation. Flair embeddings are obtained from the hidden states of a bi-directional recurrent neural network (BiRNN). They are trained without any explicit notion of a word. Instead, Flair models a word as sequences of characters. Moreover, these embeddings are determined by the text surrounding them, i.e., the same word will have different embeddings depending on its contextual usage. The variant of the Flair embedding used in this study is the Pooled Flair embedding (Weber et al., 2021; Akbik et al., 2018). Furthermore, we use the forward and backward representations of Flair embeddings returned from the BiRNN. Intuitively, character-level embeddings can potentially help improve model predictions with better OOV handling.

(Bio)BERT. We also evaluate two transformer-based embeddings: BERT and BioBERT. Specifically, we use the BERT variant "bert-base-uncased" available Flair and HuggingFace (Wolf et al., 2020). BERT was pre-trained using the BooksCorpus (800M words) and English Wikipedia (2,500M words) (Devlin et al., 2019b). Likewise, BioBERT embeddings further fine-tuned BERT on PubMed (Lee et al., 2020).

BERT embeddings are based on subword tokenization, so BERT can potentially handle OOV better than word embeddings alone. Intuitively, it fits somewhere between Flair (generating word embeddings from character representations) and Word2Vec (which independently learns embeddings for each word). Likewise, each word representation is context-dependent. Hence, BERT is better at handling word polysemy by capturing word semantics in context.

5 RESULTS

CDR, CHEMDNER, and CHEBI Results. Table 3 reports the recall, precision, and F1 scores for each embedding type for the CDR, CHEMDNER, and CHEBI datasets. The reported scores are for the best models-hyperparameter combinations on their original validation datasets. Overall, we find that the Flair and BERT-based methods outperform word embeddings. The BERT embeddings result in the best performance for the CDR dataset. While in the CHEMDNER corpus, the PubMed Flair embeddings outperform the BERT embeddings (.9018 vs. .8938). For CHEMBI, the BioBERT embeddings work the best (.7720 vs. .7322 and .6372).

	Prec.	Rec.	F1
CDR + PubMed Word	.8962	.8797	.8615
CDR + PubMed Flair	.9090	.8984	.8920
CDR + BioBERT	.9030	.8913	.8971
CDR + General Word	.8046	.8006	.8026
CDR + General Flair	.8794	.8580	.8686
CDR + BERT	.9181	.9174	.9100
CHEMDNER + PubMed Word	.8963	.8887	.8846
CHEMDNER + PubMed Flair	.9133	.9112	.9018
CHEMDNER + BioBERT	.9112	.8861	.8985
CHEMDNER + General Word	.8267	.7570	.7903
CHEMDNER + General Flair	.8985	.8696	.8838
CHEMDNER + BERT	.9122	.8840	.8938
CHEBI + PubMed Word	.7384	.7123	.7251
CHEBI + PubMed Flair	.8051	.7384	.7703
CHEBI + BioBERT	.7858	.7703	.7780
CHEBI + General Word	.5999	.6793	.6372
CHEBI + General Flair	.7454	.7196	.7322
CHEBI + BERT	.7740	.7700	.7720

Table 3: CDR, CHEMDNER, and CHEBI Results.

Synthetic (Template) Results. We evaluated several Named Entity Recognition (NER) models across multiple datasets and embeddings to assess gender bias, as summarized in Table 4. Specifically, the aggregate measures in the bottom section of Table 4 highlight the overall trends in bias across embedding training data sources (PubMed vs. General) and embedding types (Word, Flair, and BERT). The bias analysis reveals that models generally perform differently on male versus female templates. Particularly, PubMed-trained (including BioBERT) embeddings across all datasets show an average precision bias of .0242 against female names. The General embeddings exhibit substantially more bias, especially in precision with an average difference of .0407. Moreover, while the average scores for Word and (Bio)BERT embeddings show less bias, the General and Flair embeddings indicate more significant bias in precision and F1 scores. These aggregate measures underscore the pervasive nature of gender bias in NER systems and the importance of addressing it in future work.

Overall, the major source of bias is that female names are being classified as chemicals. Intuitively, the word embeddings are less biased than Flair and (Bio)BERT-based embeddings because gender-related names are treated independently using word embeddings, or better, do not appear in the embeddings at all. This is particularly evident in the differences in performance between general word

embeddings and the PubMed-based word embeddings. The PubMed embeddings do not generally have any direct mentions of named (e.g., John or Jane), hence they are generally less biased than the general domain.

This finding that female names are classified as chemicals is consistent with prior research on naming conventions for brands being gendered (Lieven et al., 2015). To further investigate this, we randomly sampled 100 chemicals from all three datasets and measured the number of brand name mentions. Overall, we found one brand name in the CHEMDNER dataset, 19 in the CDR dataset, and 32 in the ASKDOC dataset, which generally matches the bias performance differences in Table 4 (i.e., biases are *generally worse* in CDR and ASKDOC datasets than the CHEMDNER dataset).

AskDoc Results. The AskDoc results, as shown in Table 5, highlight various biases in chemical NER systems on real-world data. This table presents results from models trained on CDR, CHEMDNER, and CHEBI datasets, using different embeddings such as Word, Flair, and (Bio)BERT. Again, the embeddings are both trained on general and domain-specific corpora (e.g., PubMed).

For the fine-grained results, we note that bias and performance can vary depending on unique combinations of the dataset and embedding types. However, for the aggregate results, we have two major findings. First, we find that general domain embeddings are more biased when applied to the chemical NER task (e.g., .0056 vs. .0330 precision). This further emphasizes the results from the synthetic data study. Second, we find that word embeddings are generally less fair than Flair BERT/BioBERT embeddings for precision (.0071 vs. .0156 and .0352) and F1 (.0158 vs. .0242 and .0245).

What does this mean in real-world terms? Considering a sample of 1,000,000 chemical mentions across male and female posts (a relatively small number in social media), a 4% recall difference results in an additional 40,000 false negatives for the female group. For example, there are well-known health disparities between men and women for depression, with absolute differences of less than 3% (Salk et al., 2017). Hence, a 4% recall difference can substantially impact findings if applied researchers or practitioners use out-of-domain models to understand medications for this disease. Such a considerable gap can markedly affect the utility and trustworthiness of these predictive outcomes

		Male		Female		I	Difference		
Dataset + Embeddings	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1
CDR + PubMed Word	1	.8230	.9029	1	.8230	.9029	.0000	.0000	.0000
CDR + PubMed Flair	.9711	.9486	.9597	.9344	.9494	.9418	.0367	0008	.0179
CDR + BioBERT	.8446	.9044	.8733	.7764	.9036	.8352	.0682	.0007	.0381
CDR + General Word	.9536	.6756	.7907	.8530	.6756	.7539	.1006	.0000	.0368
CDR + General Flair	.8325	.9400	.8827	.7610	.9397	.8408	.071 5	.0003	.0419
CDR + BERT	.9867	.8493	.9128	.9728	.8444	.9041	.0138	.0048	.0087
CHEMDNER + PubMed Word	.9990	.8625	.9257	.9968	.8622	.9246	.0021	.0003	.0011
CHEMDNER + PubMed Flair	.9982	.8836	.9374	.9885	.8852	.9340	.0097	007	.0034
CHEMDNER + BioBERT	.8847	.8968	.8907	.8625	.8963	.8790	.0222	.0005	.0116
CHEMDNER + General Word	.9614	.1966	.3264	.9311	.1957	.3233	.0302	.0009	.0030
CHEMDNER + General Flair	.9559	.8437	.8963	.9105	.8433	.8755	.0454	.0004	.0208
CHEMDNER + BERT	.9913	.8768	.9306	.9680	.8762	.9198	.0233	0006	.0107
ASKDOC + PubMed Word	.9739	.9330	.9530	.9739	.9330	.9530	.0000	.0000	.0000
ASKDOC + PubMed Flair	.8833	.9523	.9164	.8278	.9519	.8852	.0555	.0005	.0312
ASKDOC + BioBERT	.8026	.9444	.8677	.7703	.9443	.8483	.0323	.0001	.0194
ASKDOC + General Word	.9681	.6607	.7854	.9711	.6604	.7862	0030	.0003	0008
ASKDOC + General Flair	.8707	.9491	.9079	.8166	.9468	.8765	.0542	.0023	.0315
ASKDOC + BERT	.9394	.9288	.9340	.8967	.9282	.9121	.0427	.0006	.0220
CHEBI + PubMed Word	.9999	.8758	.9337	.9979	.8715	.9305	.0019	.0042	.0033
CHEBI + PubMed Flair	.9689	.9016	.9340	.9545	.9031	.9281	.0144	0015	.0060
CHEBI + PubMed BERT	.9170	.8673	.8914	.8690	.8689	.8690	.0480	0016	.0225
CHEBI + General Word	.9538	.5073	.6620	.9147	.4956	.6424	.0391	.0118	.0196
CHEBI + General Flair	.9832	.8720	.9242	.9677	.8701	.9163	.0155	.0019	.0079
CHEBI + BERT	.9779	.8892	.9314	.9223	.8882	.9048	.0556	.0011	.0266
Aggregate Measures									
AVERAGE PubMed/BioBERT	.9370	.8994	.9155	.9126	.8994	.9026	.0242	.0002	.0129
AVERAGE General	.9479	.7658	.8238	.9071	.7637	.8047	.0407	.0020	.0191
AVERAGE Word	.9763	.6919	.7850	.9548	.6897	.7771	.0214	.0022	.0079
AVERAGE Flair	.9329	.9114	.9199	.8951	.9112	.8998	.0378	0002	.0201
AVERAGE (Bio)BERT	.9181	.8946	.9040	.8797	.8938	.8840	.0382	.0011	.0199

Table 4: Synthetic (Template) Data Results. We **bold** the more biased aggregate measures and all differences greater than .01 to easily read the main findings.

in practical scenarios.

AskDoc Error Analysis. Our experiments show that Chemical NER systems are biased. However, what specifically is causing the errors? For the synthetic data, the answer is gender-related names. To understand the errors in the AskDoc data, we analyzed the errors made by the best NER models trained on the out-of-domain corpus (CHEMDNER and CDR) and tested the male and female splits of the AskDocs corpus. In Figure 1, we report the ratio of false negatives for different categories of drugs/chemicals. For every false negative made by the top models of each dataset-model combination, we manually categorized them into a general chemical class (e.g., Contraceptives, Analgesics/Pain Killers, and Stimulants). Formally, let FN_m^k repre-

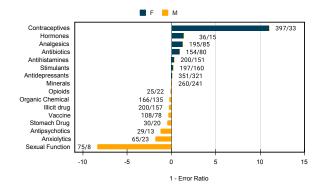


Figure 1: Ratio of false negatives for various drug categories. The ratio is represented next to each bar. For female-leaning errors, the female false negative count (FN_f^k) is in the numerator. For male-leaning errors, the male false negative count (FN_m^k) is in the numerator.

	Male Female		Difference						
Dataset + Embeddings	Prec.	Rec.	F1	Prec.	Rec.	F1	Prec.	Rec.	F1
CDR + PubMed Word	.8375	.6023	.7007	.8206	.6249	.7095	0169	.0226	.0088
CDR + PubMed Flair	.8614	.6160	.7183	.8778	.6702	.7601	.0164	.0542	.0418
CDR + BioBERT	.8303	.6352	.7198	.8042	.6693	.7306	0261	.0341	.0108
CDR + General Word	.7538	.6724	.7108	.7489	.6986	.7229	0049	.0262	.0121
CDR + General Flair	.8479	.6501	.7359	.8542	.6707	.7514	.0063	.0206	.0155
CDR + BERT	.8742	.6453	.7425	.8638	.6589	.7475	0104	.0136	.0050
CHEMDNER + PubMed Word	.8057	.5966	.6855	.8158	.6049	.6947	.0101	.0083	.0092
CHEMDNER + PubMed Flair	.8891	.6155	.7274	.8871	.6282	.7356	0020	.0127	.0082
CHEMDNER + BioBERT	.8537	.6238	.7208	.8735	.6434	.7410	.0198	.0196	.0202
CHEMDNER + General Word	.7490	.5546	.6373	.7975	.5842	.6743	.0485	.0296	.0370
CHEMDNER + General Flair	.8159	.5678	.6696	.8821	.6021	.7157	.0662	.0343	.0461
CHEMDNER + BERT	.7165	.6315	.6713	.8309	.6349	.7198	.1144	.0034	.0485
CHEBI + PubMed Word	.7574	.5998	.6694	.7548	.6287	.6860	0026	.0289	.0166
CHEBI + PubMed Flair	.7540	.6415	.6932	.7571	.6740	.7131	.0031	.0325	.0199
CHEBI + BioBERT	.6896	.5969	.6399	.7380	.6148	.6708	.0484	.0179	.0309
CHEBI + General Word	.6047	.6541	.6284	.6132	.6687	.6397	.0085	.0146	.0113
CHEBI + General Flair	.6066	.5775	.5917	.6103	.6001	.6052	.0037	.0226	.0135
CHEBI + BERT	.6274	.6478	.6374	.6923	.6467	.6687	.0649	0011	.0313
		Aggre	gate Mo	easures					
AVERAGE PubMed/BioBERT	.8087	.6142	.6972	.8143	.6398	.7157	.0056	.0256	.0185
AVERAGE General	.7329	.6223	.6694	.7659	.6405	.6939	.0330	.0182	.0245
AVERAGE Word	.7514	.6133	.6720	.7585	.6350	.6879	.0071	.0217	.0158
AVERAGE Flair	.7958	.6114	.6894	.8114	.6409	.7135	.0156	.0295	.0242
AVERAGE (Bio)BERT	.7653	.6301	.6886	.8005	.6447	.7131	.0352	.0146	.0245

Table 5: AskDoc Results. We **bold** the more biased aggregate measures and all differences greater than .01 to easily read the main findings.

sent the total number of false negatives for chemical types k and male data m. Let FN_f^k represent the female false negatives. If FN_m^k is larger than FN_f^k , we define the ratio as $-(1-FN_m^k/FN_f^k)$. Likewise, if FN_f^k is greater than FN_m , then we define the ratio as $1-(FN_f^k/FN_m^k)$. Hence, when male ratios are higher, the score is negative; otherwise, it is positive.

Overall, we make several important findings. First, we find that the models make slightly more false negatives on the chemicals categories Contraceptives (e.g., birth control and Plan B One-Step), Hormones (e.g., Megace used to treat the symptoms of loss of appetite and wasting syndrome in people with illnesses such as breast cancer), Analgesics (i.e., Pain Killers such as Tylenol) and Antibiotics on the female dataset. In contrast, the models make slightly more errors in the chemical categories Anxiolytics (e.g., drugs used to treat anxiety), Antipsychotics (e.g., chemicals used to manage psychosis, principally in schizophrenia), and sexual function drugs (e.g., Viagra). Further-

more, while the ratio for the most male- and female-related errors (Contraceptives and Sexual Function) are similar, the absolute magnitudes are substantially different. For instance, there are 397 Contraceptive FNs in the female dataset, but only 75 Sexual Function FNs appear in the male dataset. This provides an explanation for the large differences in recall on the AskDoc corpus between the male and female datasets.

6 LIMITATION

There were several limitations to our study. First, the adjudication of disagreeing items was dependent on the judgment of a single graduate student, potentially introducing human error and bias compared to a multi-adjudicator approach. Second, the vast volume of data from the active r/AskDoc subreddit community makes the feasibility of one person's comprehensive review debatable. Although our annotation method is in line with standard practices, a more multi-faceted approach involving numerous annotators and adjudicators might

offer improved accuracy and consistency in future datasets. Third, our study focuses on binary representations of gender (ignoring non-binary people). Moreover, the Social Security's Most Popular Baby Names (SSN) names may not adequately mention immigrant-related names. Hence, the results may be European-specific.

7 ETHICAL CONSIDERATIONS

In this study, we consider binary gender biases. While binary gender is a common area of study in NLP literature (Mehrabi et al., 2020), and we follow best practices of using self-identified gender (Larson, 2017), it leaves a large portion of individuals out of the study (i.e., not counted). Moreover, we also follow prior (Mehrabi et al., 2020) by relating names to gender. Nevertheless, names are not directly related to gender identity. Hence, in future work, we intend to explore data collection methods beyond binary gender. Specifically, we plan to collect data from other groups for detailed studies of model performance.

Additionally, using data from platforms like Reddit's r/AskDocs, where individuals share personal health experiences, raises ethical concerns about the potential exposure of personally identifiable information (PII) and sensitive personal health information (PHI). While our research aims to assess gender bias without examining personal details, the potential for identifiable information necessitates careful handling to protect privacy and confidentiality, following established ethical guidelines for internet research (Fiesler et al., 2024).

8 CONCLUSION

In this paper, we evaluate the gender bias of Chemical NER systems. Moreover, we compare bias measurements from synthetic data with real-world self-identified data. We make two major findings. First, Chemical NER systems are biased with regard to gender for synthetic data. Specifically, our study found that female name-like patterns feature prominently in chemical naming conventions. This characteristic leads to a notable bias in NER systems, where female names are disproportionately identified as chemicals, inadvertently escalating the gender bias in these systems. Second, we explored the performance of these models in real-world scenarios and found that most models perform better on male-related data than female-related data. A striking revelation was

the system's poor performance when identifying chemicals frequently found in female-related data, such as mentions of contraceptives.

In conclusion, the results of our study emphasize the urgent need for deliberate bias mitigation strategies in Chemical NER systems. Our findings spotlight the necessity for incorporating both synthetic and real-world data considerations to develop models that are both fair and reliable. There are two major paths for future research. First, while large language models are still behind in terms of performance for NER systems (Wang et al., 2023), they are becoming more common. Future work should explore biases in prompting-based NER solutions. Second, we plan to explore how the chemical NER biases impact downstream tasks such as relationship classification and question answering.

ACKNOWLEDGEMENTS

This material is based upon work supported by the National Science Foundation (NSF) under Grant No. 1947697 and NSF award No. 2145357.

References

Pankaj Agarwal and David B. Searls. 2008. Literature mining in support of drug discovery. *Briefings in bioinformatics*, 9 6:479–92.

Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019a. Flair: An easy-to-use framework for state-of-the-art nlp. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 54–59.

Alan Akbik, Tanja Bergmann, and Roland Vollgraf. 2019b. Pooled contextualized embeddings for named entity recognition. 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 724–728.

Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649.

Dhananjay Ashok and Zachary C Lipton. 2023. Promptner: Prompting for named entity recognition. *arXiv* preprint arXiv:2305.15444.

Pinkesh Badjatiya, Manish Gupta, and Vasudeva Varma. 2019. Stereotypical bias removal for hate speech detection task using knowledge-based generalizations. In *The World Wide Web Conference*, pages 49–59.

- James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *Journal of machine learning research*, 13(2).
- Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270.
- Elisa Chilet-Rosell. 2014. Gender bias in clinical research, pharmaceutical marketing, and the prescription of drugs. *Global Health Action*, 7(1):25484.
- Yu-Wen Chiu, Wen-Chao Yeh, Sheng-Jie Lin, and Yung-Chun Chang. 2021. Recognizing chemical entity in biomedical literature using a bert-based ensemble learning methods for the biocreative 2021 nlm-chem track. In *Proceedings of the seventh BioCreative challenge evaluation workshop*.
- Davide Cirillo, Silvina Catuara-Solarz, Czuee Morey, Emre Guney, Laia Subirats, Simona Mellino, Annalisa Gigante, Alfonso Valencia, María José Rementeria, Antonella Santuccione Chadha, et al. 2020. Sex and gender differences and biases in artificial intelligence for biomedicine and healthcare. *NPJ digital medicine*, 3(1):81.
- Donald C Comeau, Chih-Hsuan Wei, Rezarta Islamaj Doğan, and Zhiyong Lu. 2019. Pmc text mining subset in bioc: about three million full-text articles and growing. *Bioinformatics*, 35(18):3533–3535.
- Paula Czarnowska, Yogarshi Vyas, and Kashif Shah. 2021. Quantifying social biases in nlp: A generalization and empirical comparison of extrinsic fairness metrics. *Transactions of the Association for Computational Linguistics*, 9:1249–1267.
- Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. 2019. Racial bias in hate speech and abusive language detection datasets. In *Proceedings of the Third Workshop on Abusive Language Online*, pages 25–35, Florence, Italy. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019a. 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.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019b. 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. Association for Computational Linguistics.

- Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 67–73.
- Lizzy Farrugia and Charlie Abela. 2020. Mining drugdrug interactions for healthcare professionals. *Proceedings of the 3rd International Conference on Applications of Intelligent Systems*.
- Casey Fiesler, Michael Zimmer, Nicholas Proferes, Sarah Gilbert, and Naiyan Jones. 2024. Remember the human: A systematic review of ethical considerations in reddit research. *Proceedings of the ACM on Human-Computer Interaction*, 8(GROUP):1–33.
- Joel Escudé Font and Marta R Costa-jussà. 2019. Equalizing gender bias in neural machine translation with word embeddings techniques. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 147–154.
- Andrew Gaut, Tony Sun, Shirlyn Tang, Yuxin Huang, Jing Qian, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, et al. 2020. Towards understanding gender bias in relation extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2943–2953.
- Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ricardo Muñoz Sánchez, Mugdha Pandya, and Adam Lopez. 2021. Intrinsic bias metrics do not correlate with application bias. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1926–1940.
- Evan Hepler-Smith. 2015. "just as the structural formula does": Names, diagrams, and the structure of organic chemistry at the 1892 geneva nomenclature congress. *Ambix*, 62(1):1–28.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Svetlana Kiritchenko and Saif Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 43–53.
- Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in large language models. In *Proceedings of The ACM Collective Intelligence Conference*, pages 12–24.
- Martin Krallinger, Obdulia Rabal, Florian Leitner, Miguel Vazquez, David Salgado, Zhiyong Lu, Robert Leaman, Yanan Lu, Donghong Ji, Daniel M Lowe, et al. 2015. The chemdner corpus of chemicals and drugs and its annotation principles. *Journal of cheminformatics*, 7(1):1–17.

- Brian Larson. 2017. Gender as a variable in naturallanguage processing: Ethical considerations. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 1–11.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.
- Jiao Li, Yueping Sun, Robin J Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J Mattingly, Thomas C Wiegers, and Zhiyong Lu. 2016. Biocreative v cdr task corpus: a resource for chemical disease relation extraction. *Database*, 2016.
- Jing Li, Aixin Sun, Jianglei Han, and Chenliang Li. 2020. A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering*, 34(1):50–70.
- Theo Lieven, Bianca Grohmann, Andreas Herrmann, Jan R Landwehr, and Miriam Van Tilburg. 2015. The effect of brand design on brand gender perceptions and brand preference. *European Journal of Marketing*.
- Pilar López-Úbeda, Manuel Carlos Díaz-Galiano, L Alfonso Ureña-López, and M Teresa Martín-Valdivia. 2021. Combining word embeddings to extract chemical and drug entities in biomedical literature. *BMC bioinformatics*, 22(1):1–18.
- Brandon Lwowski and Anthony Rios. 2021. The risk of racial bias while tracking influenza-related content on social media using machine learning. *Journal of the American Medical Informatics Association*, 28(4):839–849.
- Maria Mammì, Rita Citraro, Giovanni Torcasio, Gennaro Cusato, Caterina Palleria, and Eugenio Donato di Paola. 2013. Pharmacovigilance in pharmaceutical companies: An overview. *Journal of Pharmacology & Pharmacotherapeutics*, 4:S33 S37.
- Lisa A Martin, Harold W Neighbors, and Derek M Griffith. 2013. The experience of symptoms of depression in men vs women: analysis of the national comorbidity survey replication. *JAMA psychiatry*, 70(10):1100–1106.
- Ninareh Mehrabi, Thamme Gowda, Fred Morstatter, Nanyun Peng, and A. G. Galstyan. 2020. Man is to person as woman is to location: Measuring gender bias in named entity recognition. *Proceedings of the 31st ACM Conference on Hypertext and Social Media*.
- Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. In 1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings.

- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Shubhanshu Mishra, Sijun He, and Luca Belli. 2020. Assessing demographic bias in named entity recognition. *arXiv preprint arXiv:2008.03415*.
- Marzieh Mozafari, Reza Farahbakhsh, and Noël Crespi. 2020. Hate speech detection and racial bias mitigation in social media based on bert model. *PloS one*, 15(8):e0237861.
- David Mueller, Nicholas Andrews, and Mark Dredze. 2020. Sources of transfer in multilingual named entity recognition. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8093–8104.
- World Health Organization and Key Centre for Women's Health in Society. 2009. Mental health aspects of women's reproductive health: a global review of the literature.
- Karen O'Connor, Pranoti Pimpalkhute, Azadeh Nikfarjam, Rachel Ginn, Karen L Smith, and Graciela Gonzalez. 2014. Pharmacovigilance on twitter? mining tweets for adverse drug reactions. In *AMIA annual symposium proceedings*, volume 2014, page 924. American Medical Informatics Association.
- Ji Ho Park, Jamin Shin, and Pascale Fung. 2018. Reducing gender bias in abusive language detection. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 2799–2804.
- James Pustejovsky and Amber Stubbs. 2012. *Natural Language Annotation for Machine Learning: A guide to corpus-building for applications*. "O'Reilly Media, Inc.".
- S Pyysalo, F Ginter, H Moen, T Salakoski, and S Ananiadou. 2013. Distributional semantics resources for biomedical text processing. In *Proceedings of LBM 2013*, pages 39–44.
- Joseph L Riley III, Michael E Robinson, Emily A Wise, Cynthia D Myers, and Roger B Fillingim. 1998. Sex differences in the perception of noxious experimental stimuli: a meta-analysis. *Pain*, 74(2-3):181–187.
- Anthony Rios. 2020. Fuzze: Fuzzy fairness evaluation of offensive language classifiers on african-american english. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 881–889.
- Tim Rocktäschel, Michael Weidlich, and Ulf Leser. 2012. Chemspot: a hybrid system for chemical named entity recognition. *Bioinformatics*, 28(12):1633–1640.

- Rachel H Salk, Janet S Hyde, and Lyn Y Abramson. 2017. Gender differences in depression in representative national samples: Meta-analyses of diagnoses and symptoms. *Psychological bulletin*, 143(8):783.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, pages 1668–1678.
- Mirsada Serdarevic, Catherine W Striley, and Linda B Cottler. 2017. Gender differences in prescription opioid use. *Current opinion in psychiatry*, 30(4):238.
- Matthew Shardlow, Nhung Nguyen, Gareth Owen, Claire O'Donovan, Andrew Leach, John McNaught, Steve Turner, and Sophia Ananiadou. 2018. A new corpus to support text mining for the curation of metabolites in the chebi database. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. 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)*, pages 3407–3412.
- David RM Smith, F Christiaan K Dolk, Timo Smieszek, Julie V Robotham, and Koen B Pouwels. 2018. Understanding the gender gap in antibiotic prescribing: a cross-sectional analysis of english primary care. *BMJ open*, 8(2):e020203.
- Cong Sun, Zhihao Yang, Lei Wang, Yin Zhang, Hongfei Lin, and Jian Wang. 2021. Deep learning with language models improves named entity recognition for pharmaconer. *BMC bioinformatics*, 22(1):1–16.
- Lena Thunander Sundbom, Kerstin Bingefors, Kerstin Hedborg, and Dag Isacson. 2017. Are men undertreated and women over-treated with antidepressants? findings from a cross-sectional survey in sweden. *BJPsych bulletin*, 41(3):145–150.
- Amir Pouran Ben Veyseh, Franck Dernoncourt, Bonan Min, and Thien Huu Nguyen. 2022. Generating complement data for aspect term extraction with gpt-2. In *Proceedings of the Third Workshop on Deep Learning for Low-Resource Natural Language Processing*, pages 203–213.
- Zhen Wan, Fei Cheng, Zhuoyuan Mao, Qianying Liu, Haiyue Song, Jiwei Li, and Sadao Kurohashi. 2023. Gpt-re: In-context learning for relation extraction using large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3534–3547.
- Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. Gpt-ner: Named entity recognition via large language models. *arXiv preprint arXiv:2304.10428*.

- Leon Weber, Mario Sänger, Jannes Münchmeyer, Maryam Habibi, Ulf Leser, and Alan Akbik. 2021. Hunflair: an easy-to-use tool for state-of-the-art biomedical named entity recognition. *Bioinformatics*, 37(17):2792–2794.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Tongxuan Zhang, Hongfei Lin, Yuqi Ren, Zhihao Yang, Jian Wang, Xiaodong Duan, and Bo Xu. 2021. Identifying adverse drug reaction entities from social media with adversarial transfer learning model. *Neurocomputing*, 453:254–262.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 2019. Gender bias in contextualized word embeddings. 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 629–634, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2979–2989.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20.

A Appendix

A.1 Datasets

CDR (Li et al., 2016) We use the BioCreative V CDR shared task corpus. The CDR corpus comprises 1,500 PubMed articles with 4,409 annotated chemicals, 5,818 diseases, and 3,116 chemical disease interactions. This corpus is designed to address two distinct tasks: Relation classification and NER. For this study, we focus on the NER for chemical entities. The annotator agreement for this corpus was .87. Finally, we used the same train,

Data	Embedding	Fine-tuning	hidden_size	dropout	lr
	general word	TRUE	128	0.4	0.1
	general flair	TRUE	256	0.3	0.1
CDR	BERT	TRUE	256	0.2	0.05
CDK	Pubmed word	FALSE	128	0.2	0.1
	pubmed flair	FALSE	128	0.4	0.1
	BioBERT	TRUE	1024	0.5	0.05
	general word	TRUE	1024	0.2	0.1
	general flair	TRUE	512	0.5	0.1
CHEMD	BERT	TRUE	1024	0.3	0.025
СПЕМІ	Pubmed word	TRUE	256	0.3	0.1
	pubmed flair	FALSE	128	0.2	0.05
	BioBERT	TRUE	1024	0.2	0.025
	general word	TRUE	1024	0.2	0.1
	general flair	TRUE	512	0.5	0.1
Askdoc	bert	TRUE	1024	0.3	0.025
ASKUOC	Pubmed word	TRUE	256	0.3	0.1
	pubmed flair	FALSE	128	0.2	0.05
	biobert	TRUE	128	0.2	0.01
	general word	TRUE	128	0.4	0.1
CHEBI	General Flair	TRUE	128	0.3	0.1
	BERT	TRUE	1024	0.5	0.05
СПЕВІ	Pubmed word	TRUE	128	0.4	0.1
	Pubmed flair	FALSE	512	0.3	0.1
	BioBERT	TRUE	256	0.4	0.05

Table 6: Comprehensive List of Hyperparameters Investigated in the Search for the Optimal Model

validation, and test splits from the shared task for our experiments.

CHEMDNER (Krallinger et al., 2015) The CHEMDNER corpus includes abstracts from 10000 chemistry-related journals published in 2013 on PubMed. Each abstract was manually annotated for chemical mentions. These mentions were categorized into seven subtypes: abbreviation, family, formula, identifier, multiple, systematic, and trial. The BioCreative organizers divided the corpus into training (3500 abstracts), development (3500 abstracts), and test (3000 abstracts) sets. The BioCreative IV CHEMDNER corpus comprises 84,355 chemical mention annotations across 10,000 abstracts, with an inter-annotator agreement of .91 (Krallinger et al., 2015). For this study, we only use the major Chemical annotations and ignore the subtypes for consistency across corpora. Finally, we use the same train, validation, and test splits used in the shared task for our experiments.

CHEBI (Shardlow et al., 2018). We also use the ChEBI corpus, an extensive dataset consisting of 199 annotated abstracts and 100 full papers. This corpus contains over 15,000 named entity annotations and more than 6,000 inter-entity relations, specifically aligned with the needs of the ChEBI database curators. The dataset has annotated chemicals, proteins, species, biological activities, and spectral data. Moreover, it has a high inter-annotator agreement of 0.80-0.89 (F1 score, strict-matching). It also categorizes relationships into several types such as Isolated From, Associated With, Binds With, and Metabolite Of, offering a detailed view of the interactions between metabolites and other entities. This corpus is not only a rich source for exploring lexical characteristics of metabolites and associated entities but also serves as a critical resource for training machine learning algorithms in the recognition of these entities and their relations in the biochemical context.

	Total Male	FNR Male	Total Female	FNR Female	
Contraceptives	33	1.0000	408	.9730	
Hormones	170	.0882	230	.1565	
Analgesics	571	.1489	952	.2048	
Antibiotics	326	.2454	347	.4438	
Antihistamines	270	.5593	295	.6780	
Stimulants	522	.3065	390	.5051	
Antidepressants	781	.4110	1043	.3365	
Minerals	605	.3983	785	.3312	
Opioids	43	.5814	95	.2316	
Organic Chemical	441	.3764	346	.3902	
Illicit drug	353	.5666	311	.5048	
Vaccine	108	1.0000	78	1.0000	
Stomach Drug	55	.5455	44	.4545	
Antipsychotics	47	.6170	95	.1368	
Anxiolytics	126	.5603	100	.2300	
Sexual Function Drug	78	.9615	8	1.0000	
PCC between Total and FNR	5	58	26		

Table 7: False negative rate (FNR) for female and male-related AskDoc datasets. The pearson correlation coefficient (PCC) between the frequency of each chemical type and the FNR for teach group is marked in the last row.

A.2 Hyper-Parameter Settings

In this section, we report the best hyperparameter for each model, shown in Table 6. Similar to random hyperparameter search (Bergstra and Bengio, 2012), we generate 100 samples using different parameters for each dataset-model combination (e.g., we generate 100 versions of BERT for the CDR dataset). For the specific hyper-parameters, we used sample dropout from .1 to .9, hidden layer sizes from {128, 256, 512, 1024}, learning rates selected from 1e-4 to 1e-1 at random, and the option of whether to fine-tune the embedding layers (i.e., True vs. False). In addition, we trained all models for 25 epochs with a mini-batch size set to 32, where only the best model on the validation dataset is saved after each epoch. Finally, all experiments were run on four NVidia GeForce GTX 1080 Ti GPUs.

A.3 Error Analysis and Discussion

Interestingly, we find that the prevalence of chemicals across gender-related posts matches the prevalence found in traditional biomedical studies. Previous research report that women have been prescribed analgesics (e.g., pain killers such as opioids) twice as often as men (Chilet-Rosell, 2014; Serdarevic et al., 2017). While there is still limited understanding about whether men are under-

	FNR	wFNR
Male	.3948	.6875
Female	.4064	.8088
Gap	.0116	.1213
Ratio	1.0294	1.1764

Table 8: FNR and weighted FNR (wFNR) results.

prescribed or women are over-prescribed, the disparities in prescriptions are evident. Thus, the finding in Figure 1 that we receive twice as many analgesics FNs for female data is important. Depending on the downstream application of the Chemical NER system, these performance disparities may potentially increase harm to women. For example, if more varieties of drugs are prescribed to women, but our system does not detect them, then an ADR detection system will not be able to detect important harms.

We also find differences in Antibiotic FNs in Figure 1. There have also been medical studies showing gender differences in Antibiotic prescriptions. For example, a recent meta-analysis of primary care found that women received more antibiotics than men, especially women aged 16-54,

receiving 36%–40% more than males of the same age (Smith et al., 2018). Again, if we do not detect many of the antibiotics prescribed to women, this can cause potential health disparities in downstream ADR (and other) systems.

Next, in Table 7, we report the false negative rate (FNR) for each category along with the general frequency of each category. Using the Pearson correlation coefficient, we relate the frequency of each category with the false negative rate for the male and female groups, respectively. Intuitively, we would expect the false negative rate to go down as the frequency increases, which matches our findings. However, we find that the correlation is much stronger for the male group than the female group.

In Table 8, we report the FNR for the female and male groups, respectively. We also introduce a new metric, weighted FNR, which assigns importance scores for each of the FNRs shown to create a macro-averaged metric. Intuitively, the distribution of categories is different for both the male and female groups. So, we want to test whether the FNR scores are distributed uniformly across all categories, irrespective of, whether the errors are more concentrated for gender-specific categories. More errors in gender-specific categories can adversely impact a group that is not captured with the global FNR metric. Formally, we define wFNR for the

female group as

$$wFNR^f = \sum_{i}^{N} w_i^f FNR_i^f$$

where FNR_i^f represents the female false negative rate for category i. Likewise, w_i^f is defined as

$$w_i^f = \frac{1}{\sum_i w_i^f} \cdot \frac{N_i^f / N^f}{N_i^m / N^m}$$

where N_i^f and N_f^m represent the total number of times a category i appears for the female and male groups, respectively. Intuitively, we are dividing the ratio of each category for female and male groups. So, if a category appears more often for females than males, proportionally, then the score will be higher. We normalize these scores for each group so they sum to one. Overall, we find an absolute gap of more than 1% (3% relative difference) between the FNR for male and female groups. But, even worse, there is a much larger gap (.1213 vs .0116) when using wFNR. This result suggests that many of the false negatives are concentrated for gender-specific categories (e.g., contraceptives) for the female group more than the male group.