

Evaluating Knowledge Fusion Models on Detecting Adverse Drug Events in Text

Philipp Wegner^{1,2}, Holger Fröhlich^{1,3}, Sumit Madan^{1,*}

¹ Department of Bioinformatics, Fraunhofer Institute for Algorithms and Scientific Computing (SCAI), Schloss Birlinghoven, 53757 Sankt Augustin, Germany

² German Center for Neurodegenerative Diseases (DZNE), Venusberg Campus 1, 53127 Bonn, Germany

³ Bonn-Aachen International Center for Information Technology (B-IT), University of Bonn, 53115 Bonn, Germany

*Corresponding authors: sumit.madan@scai.fraunhofer.de

Abstract

Background: Detecting adverse drug events (ADE) of drugs that are already available on the market is an essential part of the pharmacovigilance work conducted by both medical regulatory bodies and the pharmaceutical industry. Concerns regarding drug safety and economic interests serve as motivating factors for the efforts to identify ADEs. Hereby, social media platforms play an important role as a valuable source of reports on ADEs, particularly through collecting posts discussing adverse events associated with specific drugs.

Methodology: We aim with our study to assess the effectiveness of knowledge fusion approaches in combination with transformer-based NLP models to extract ADE mentions from diverse datasets, for instance, texts from Twitter, websites like askapatient.com, and drug labels. The extraction task is formulated as a named entity recognition (NER) problem. The proposed methodology involves applying fusion learning methods to enhance the performance of transformer-based language models with additional contextual knowledge from ontologies or knowledge graphs. Additionally, the study introduces a multi-modal architecture that combines transformer-based language models with graph attention networks (GAT) to identify ADE spans in textual data.

29 **Results:** A multi-modality model consisting of the ERNIE model with knowledge on drugs
30 reached an F_1 -score of 71.84% on CADEC corpus. Additionally, a combination of a graph
31 attention network with BERT resulted in an F_1 -score of 65.16% on SMM4H corpus. Impressively,
32 the same model achieved an F_1 -score of 72.50% on the PSYTAR corpus, 79.54% on the ADE
33 corpus, and 94.15% on the TAC corpus. Except for the CADEC corpus, the knowledge fusion
34 models consistently outperformed the baseline model, BERT.

35 **Conclusion:** Our study demonstrates the significance of context knowledge in improving the
36 performance of knowledge fusion models for detecting ADEs from various types of textual data.

37

38 **Keywords:** Adverse Drug Reaction Detection, Transformers, Graph Neural Network,
39 Knowledge Fusion.

40 **Author Summary:** Adverse Drug Events (ADEs) are one of the main aspects of drug safety
41 and play an important role during all phases of drug development, including post-marketing
42 pharmacovigilance. Negative experiences with medications are frequently reported in textual
43 form by individuals themselves through official reporting systems or social media posts, as well
44 as by doctors in their medical notes. Automated extraction of ADEs allows us to identify these
45 in large amounts of text as they are produced every day on various platforms. The text sources
46 vary highly in structure and the type of language included which imposes certain challenges on
47 extraction systems. This work investigates to which extent knowledge fusion models may
48 overcome these challenges by fusing structured knowledge coming from ontologies with
49 language models such as BERT. This is of great interest since the scientific community provides
50 highly curated resources in the form of ontologies that can be utilized for tasks such as extracting
51 ADEs from texts.

52 Introduction

53 An adverse drug event (ADE) can be defined as “an injury resulting from a medical intervention
54 related to a drug” [1]. ADEs as a major aspect of drug safety are objective of interest in the

55 pharmacovigilance efforts done by pharmacological companies as well as medical regulatory
56 bodies. Negative experiences with certain medications are frequently reported in textual form
57 by individuals themselves through official reporting systems or social media posts, as well as by
58 doctors in their medical notes. The mentioned ADEs are often hidden in unstructured text, and
59 the process of identifying and extraction of ADE entities from such text requires a significant
60 amount of a medical professional's time. Performing large-scale automatic extraction from a
61 variety of text sources could help domain experts in quickly identifying new ADEs. However, this
62 extraction process requires robust and highly accurate text mining methods.

63 In recent years, the natural language processing (NLP) field has made significant advancements
64 with transformer-based language models such as BERT [2] or GPT [3]. These models have set
65 new benchmarks in several NLP tasks. Furthermore, these models have been successfully
66 applied to detect ADEs from textual documents [1, 4–6]. There are mainly two different types of
67 texts mentioning ADEs such as reports or scientific publications written by medical professionals
68 and reports provided by the patient or their relatives themselves. Social media texts differ from
69 medical reports as they often contain informal language, slang, abbreviations, and
70 colloquialisms. Additionally, these texts predominantly consist of opinions of people and contain
71 fewer factual statements. Due to the continuously growing quantity and significance of social
72 media texts, we place particular attention on analyzing patient-reported texts. In this work, we
73 considered the CADEC corpus [5] that contains annotated texts from <https://askapatient.com>,
74 which is a forum dedicated to collecting drug experiences and a corpus, here referred to as
75 SMM4H [6], that comprises annotated Twitter postings. Moreover, we evaluate our models on
76 three more corpora, namely PSYTAR [7], TAC [8], and ADE [9]. The CADEC, SMM4H, and
77 PSYTAR were derived from sources where patients authored the texts themselves, whereas
78 the ADE and TAC were composed by medical experts written in formal and scientific language.
79 Further details on the corpora are given in Section Datasets.

80 It is important to highlight previous scientific initiatives that have aimed to extract ADEs from
81 texts. Sboev et al. [10] elaborated on the performance of various transformer models evaluated

82 on CADEC, where they reported an F₁-score of 69.68% for strict matches (exact matching
83 between true and predicted instances) using the XLM-Roberta-large model that ranked best
84 among all considered models. Additionally, Portelli et al. [11] provided a performance overview
85 of different transformer models on CADEC and SMM4H, in which they reported F₁-scores of
86 67.95% and 62.15%, respectively. Portelli et al. [11] reported that a SpanBERT-based approach
87 yielded the best results. Furthermore, Ge et al. [4] offered a federated learning methodology for
88 the ADE detection problem and evaluated it on both datasets. This approach was able to
89 achieve for relaxed matches (partial overlap of true and predicted instances) an F₁ of 84.55%
90 on CADEC and 67.8% on SMM4H corpus. For strict matches, 65.16% and 32.69% were
91 reported for the same corpora by the authors. Ramesh et al. [12] presented their solution to the
92 2021 SMM4H shared task 1 that adopts the roBERTa base model to extract ADE mentions,
93 which reached a relaxed F₁-score of 50% on the final test set. Furthermore, Raval et al. [13]
94 presented an interesting strategy by tackling text classification concerning ADEs as well as the
95 actual ADE span extraction with a multi-task learning approach that used the T5 as a pre-trained
96 encoder-decoder transformer model. They could reach the strict F₁-score of 69.8% on CADEC
97 and 71.3% on SMM4H corpus as well as the relaxed F₁-scores of 79.1% and 75.1%,
98 respectively. Another notable work that deserves mention is of Haq et al. [14] as they evaluated
99 their NLP pipeline on the ADE corpus [9] as well as on CADEC and SMM4H. The end-to-end
100 system proposed by Haq et al. was able to report strict macro-averaged F₁-scores of 91.7%,
101 78.7%, and 76.7% on the ADE, CADEC, and SMM4H corpora respectively. Furthermore,
102 Miftahutdinov & Tutubalina [15] evaluated BERT on the PSYTAR corpus and were able to reach
103 an accuracy of 83.07% during the task of normalizing the ADE entities to a controlled
104 vocabulary. Analogously the authors reported accuracy scores of 88.84% on CADEC as well as
105 89.64% on SMM4H during the entity normalization task. Finally, in the 2017 Text Analysis
106 Conference (TAC) a team from the University of Texas Health Science Center at Houston was
107 able to achieve a micro-averaged F₁-score of 82.48% over all entities of the TAC corpus
108 including ADE mentions. The participants from Houston were able to reach that score by utilizing
109 a bi-directional LSTM model.

110 Moreover, Stanovsky et al. [16] adopted a fusion learning approach by combining contextual
111 knowledge from DBpedia with a Bi-LSTM. By doing so the authors reported an F1-score of
112 93.4% on the CADEC corpus. Fusion model approaches are often able to increase performance
113 in comparison with standalone transformer models. Zhang et al. [17] reported a performance
114 increase from 73.5% F₁-score using a BERT model to 75.5% adopting ERNIE as a fusion
115 learning model evaluating however on the Open Entity dataset [18]. Liu et al. [19] published an
116 alternative approach that demonstrates the advantages of transformer-based language
117 encoding with contextual knowledge, Their K-BERT model achieved a notable increase of 0.04
118 in the F1-score on a question-answering task.

119 In this study, we conducted a series of experiments to assess the effectiveness of knowledge
120 fusion methods in combination with transformer-based NLP models for extracting ADEs from
121 unstructured texts. We performed these experiments on a total of five diverse text corpora. To
122 incorporate contextualized knowledge, we constructed a knowledge graph (KG) that included
123 drug brand names and integrated a symptom ontology. This combination proved to be well-
124 suited for analyzing ADE-related texts. Additionally, we utilized graph neural network (GNN)
125 techniques, specifically a graph attention network, to learn representations of drug and symptom
126 entities within the KG. These representations were subsequently integrated into transformer
127 models through a fusion learning approach. We compared our proposed model architecture
128 against ERNIE, a well-established knowledge fusion language model, as well as two non-
129 knowledge fusion models, namely BERT and BioBERT.

130 Materials and Methodology

131 First, we introduce different datasets and knowledge resources used in our work and
132 subsequently we present the knowledge fusion models that have been developed for the
133 purpose of detecting ADEs from textual corpora.

134 Datasets

135 CADEC

136 The CSIRO Adverse Drug Event Corpus (CADEC) [5] is an annotated text corpus published in
137 2015 that consists of forum posts from askapatient.com and comes with 5 different types of
138 annotations: ADE, Drug, Disease, Symptom, and Finding (any other clinical finding).

139

140 The whole CADEC corpus includes reports on 12 drugs such as Diclofenac or Lipitor.
141 Diclofenac (<https://go.drugbank.com/drugs/DB00586>) is a non-steroidal anti-inflammatory drug
142 that is used to treat pain and inflammation from different sources while Lipitor
143 (<https://go.drugbank.com/drugs/DB01076>) lowers lipid levels and reduces the risk of
144 cardiovascular diseases. The CADEC corpus is composed of 1,253 posts with 7,398 sentences
145 in total, where 1,107 posts contain at least one ADE mention (see Table 1). This adds up to
146 7,409 ADE spans with an average post length of six sentences. Finally, all posts were written
147 between January 2001 and September 2013 by patients between 17 and 84.

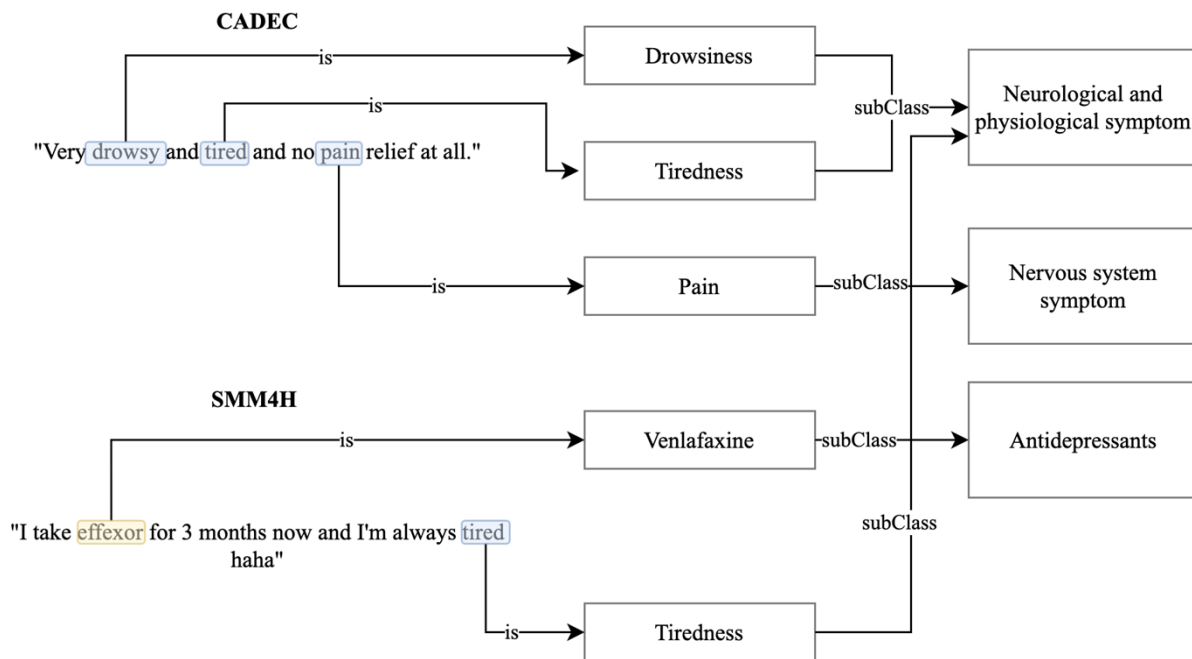
148 SMM4H

149 The second dataset used in this work is the SMM4H corpus [6], which is one of the datasets
150 provided to the participants of the Social Media Mining for Health Applications (#SMM4H)
151 Shared Task 2021 (<https://healthlanguageprocessing.org/smm4h-shared-task-2021/>). In this
152 work, we focus on the corpus for Subtask 1b, which is about extracting ADE mentions from
153 Twitter posts. We ignore Subtasks 1a which dealt with classifying Tweets containing an ADE
154 and 1c which tackled the normalization of ADEs to MedDRA.

155

156 There are differences between the SMM4H Subtask 1b corpus and the CADEC, while the
157 biggest difference might be that CADEC has annotations of 5 different types whereas the corpus
158 of Subtask 1b of SMM4H has only adverse drug reaction mentions tagged. The corpus is
159 composed of 1,300 tweets with 1,800 annotated ADE spans (see Table 1). On average each
160 tweet has 21 words and two sentences.

161



162

163

164

165

166

167

168

Fig. 1: CADEC and SMM4H example phrases that are enriched with contextual knowledge about drugs and symptoms. The sentence from CADEC “Very drowsy and tired and no pain relief at all.” can be equipped with symptom classes such as Drowsiness and Tiredness, which are subclasses of Neurological and physiological symptom class, as well as Pain, which is a subclass of Nervous system symptom.

169

PsyTAR

170

The third corpus considered in this work is the corpus presented by Zolnoori et al. [7]. The

171

“Psychiatric Treatment Adverse Reactions” (PsyTAR) corpus contains 891 drug reviews from

172

askapatient.com which is the same source as the previously mentioned CADEC corpus. The

173

corpus contains reviews for four drugs (Zoloft, Lexapro, Cymbalta, and Effexor XR) and holds a

174

total of 6009 sentences with 4813 ADE mentions (see Table 1). On average each post contains

175

7 (6.7) sentences. Further note that the PsyTAR text corpus contains, besides ADE mentions,

176

6 other annotation types, which are Withdrawal Symptoms (WDs), Signs/Symptoms/Illness

177

(SSIs), Drug Indications (DIs), Drug Effectiveness (EF), and Drug Ineffectiveness (INF) and other,

178

not applicable, mentions.

179

TAC

180

The TAC corpus [8] was assembled from drug labels and was used in the 2017 text annotation

181

conference (TAC). The corpus consists of a set of drug labels in which ADE mentions among

182 other entities are annotated. In that conference participants were provided with the corpus and
183 challenged to extract adverse drug reactions from these drug labels. This task was referred to
184 as Task 1 within TAC. Each drug label contains on average 79 sentences and hence was split
185 into sentences to fit it into the transformer models used in this work. Each sentence contains on
186 average 33 (32.69) words. Besides ADE entities the corpus comes with annotations for Severity,
187 Factor (additional aspects of the ADE entity), Drug Class, Negation, and Animal.

188 ADE

189 The 5th and final corpus used was published by Gurulingappa et al. [9] and was constructed
190 from 3000 MEDLINE case reports. After an exhaustive annotation and harmonization process
191 that involved three annotators, the corpus holds 2972 reports. The final corpus comprises a total
192 of 5063 drugs and 5776 ADE annotations distributed over 4272 sentences (see Table 1). On
193 average each sentence contains 20 (20.09) words. Besides drug and ADE entities the corpus
194 further contains annotations for Dosage. Other than some of the corpora previously introduced,
195 the authors of the ADE corpus did not restrict the retrieved documents to a certain set of drugs
196 but rather retrieved 30.000 documents and randomly selected the 3000 case reports that were
197 further used for the annotation process.

Dataset	Document class	# Documents	# Sentences	# Drugs	# ADEs
CADEC [5]	Drug reviews	1,253	7,398	1,800	7,409
SMM4H Subtask 1b [6]	Tweets	1,300	2,107	-	1,496
PsyTAR [7]	Drug reviews	891	6,009	792	4,813
TAC Task 1 [8]	Drug labels	101	3,154	249	13,795
ADE [9]	Medline case reports	2,972	4,272	5,063	5,776

198 Table 1. Overview of the ADE datasets used in this study. Note that the SMM4H corpus does
199 not contain any drug annotations.

200 Knowledge Bases

201 In our work, we explored the enhancement of transformer models by incorporating contextual
202 knowledge through fusion models to improve the detection of adverse drug events. We utilized

203 two knowledge resources: one for encoding knowledge about symptoms and the other for
204 modeling the domain of drug space.

205 Symptom Ontology

206 The symptom ontology (SYMP) is a publicly available ontology developed in the context of the
207 Gemina system [20]. The creators designed the ontology while understanding a symptom as a
208 “perceived change in function, sensation or appearance reported by a patient indicative of a
209 disease” [20]. The ontology consists of 860 classes as well as a total of 1,586 cross-references
210 to other databases like UMLS (<https://www.nlm.nih.gov/research/umls/index.html>) or ICD (<https://www.who.int/standards/classifications/classification-of-diseases>). Furthermore, the
211 ontology comprises 5,445 axioms and class annotations such as definitions, synonyms, and
212 labels of symptoms. We use the symptoms ontology to provide context knowledge about
213 symptoms. An example of how a model can enrich sentences with symptom classes is shown
214 in Fig. 1.
215

216 Drug Resources

217 Contextual knowledge about drugs and how they function in the human body can be valuable
218 for tackling the task of ADE detection. We decided to assemble such knowledge in a structured
219 way and store it in the form of an ontology. The resulting ontology inherits information from the
220 ATC ontology and is further enriched with selected information about drugs. Fig. 1 illustrates an
221 example of how a model can enhance sentences by incorporating drug resource information.
222 Fig. 1 depicts the utilization of contextual knowledge exemplarily for CADEC and SMM4H but
223 works equally for the other three corpora.

224
225 We used three different resources to collect various information on approved drugs. Firstly, the
226 DrugBank database (version 5.1.9) [21] was used to extract drug descriptions, synonyms, and
227 product names, as well as information about drug targets. Fortunately, DrugBank provides
228 cross-references to the anatomical therapeutic chemical classification system (ATC), which
229 divides active ingredients into classes based on anatomical properties like the organ they act

230 on, chemical properties, as well as therapeutic properties [22]. DrugMechDB [23] is another drug
231 resource, which contains information about the mechanism of action of a drug in the body. This
232 mechanism is represented as a graph where each node can be of several types (such as
233 disease, drug, protein, or cell). A sub-graph was taken from this graph to obtain information
234 about the proteins that are involved in the drug mechanism, which we added to our ontology.
235 Furthermore, since this ontology is used to extract drug entities from text based on the drug
236 product names it is important to add as many brand names to the ontology as known. To
237 accomplish that, the website drugs.com was a highly useful resource for adding brand names
238 for each drug in ATC.

239
240 Finally, all of the collected knowledge on drugs was added to the ATC ontology at its respective
241 position and stored as an OWL (web ontology language) file. The resulting ontology, in this work
242 referred to as DRUGO, provides knowledge about drug names, definitions, synonyms, drug
243 targets, and information about proteins involved in the drug's action mechanism. The final
244 DRUGO ontology comprises a total of 6,441 classes.

245 Detection of Adverse Drug Events

246 Our experimental strategy to create models that can detect ADEs in texts builds upon knowledge
247 fusion models that integrate transformer-based models with knowledge graph embeddings. As
248 transformer-based models, we focus on using BERT [24] and BioBERT [25]. These models are
249 also used to create baseline results. Furthermore, we experiment with multiple fusion
250 approaches such as ERNIE and the graph concat model, which are introduced in the next
251 sections.

252 Knowledge Fusion

253 To incorporate the information from the aforementioned knowledge bases (DRUGO and SYMP)
254 into the language models, a numerical representation is necessary that effectively captures the
255 encoded knowledge. We experimented with two approaches, the first one uses the well-

256 established TransE method [26] to embed the underlying graphs of the two ontologies into a
257 vector space. Whereas, in the second approach, a GNN was incorporated for this task. More
258 specifically a graph attention network (GAT) was trained with a node classification task, which
259 provided the final node-level embeddings for the integration in the language model.

260
261 A total of three GATs were trained on the DRUGO and SYMP ontologies, as well as on an
262 ontology generated by combining SYMP and DRUGO. In this approach, ontologies are treated
263 as graphs, without taking into account any logical axioms, similar to other ontology embedding
264 approaches. All GATs have been trained identically by initially considering the ontologies as
265 graphs and assembling a set of nodes (V) from the classes of the ontology and a set of edges
266 (E) from the relations between the classes. Specifically, we derived E by treating every
267 'subClass' property as an edge. As a result, we obtained a circle-free, fully connected, directed
268 graph with 6,441 nodes and 6,440 edges for DRUGO, 860 nodes and 859 edges for SYMP,
269 and, 7301 nodes and 7300 edges for the combined KG of DRUGO and SYMP.

270
271 In the following step, initial representations for all nodes were generated. This was performed
272 by using the annotation properties of each ontology class/node and embedding these using a
273 pre-trained language model. For all graphs, this was done by using either BERT or BioBERT,
274 depending on the exact experimental setup. This lead to the representation of each node as a
275 768-dimensional real vector. Graphs derived from DRUGO and SYMP provided a top-level
276 classification with 14 classes, enabling the assignment of each node to one of these classes
277 based on its position in the graph. The third graph obtained from combining the two ontologies
278 yielded 28 classes.

279
280 Finally, a GNN was trained to predict the assigned class of each node in the graph. Note that in
281 our work, we specifically favored GAT over other GNN architectures because of its capability for
282 self-attention. The self-attention mechanism in GAT allows nodes to attend to the features of
283 their neighboring nodes. With the usage of GAT, we would like to address the issue that certain

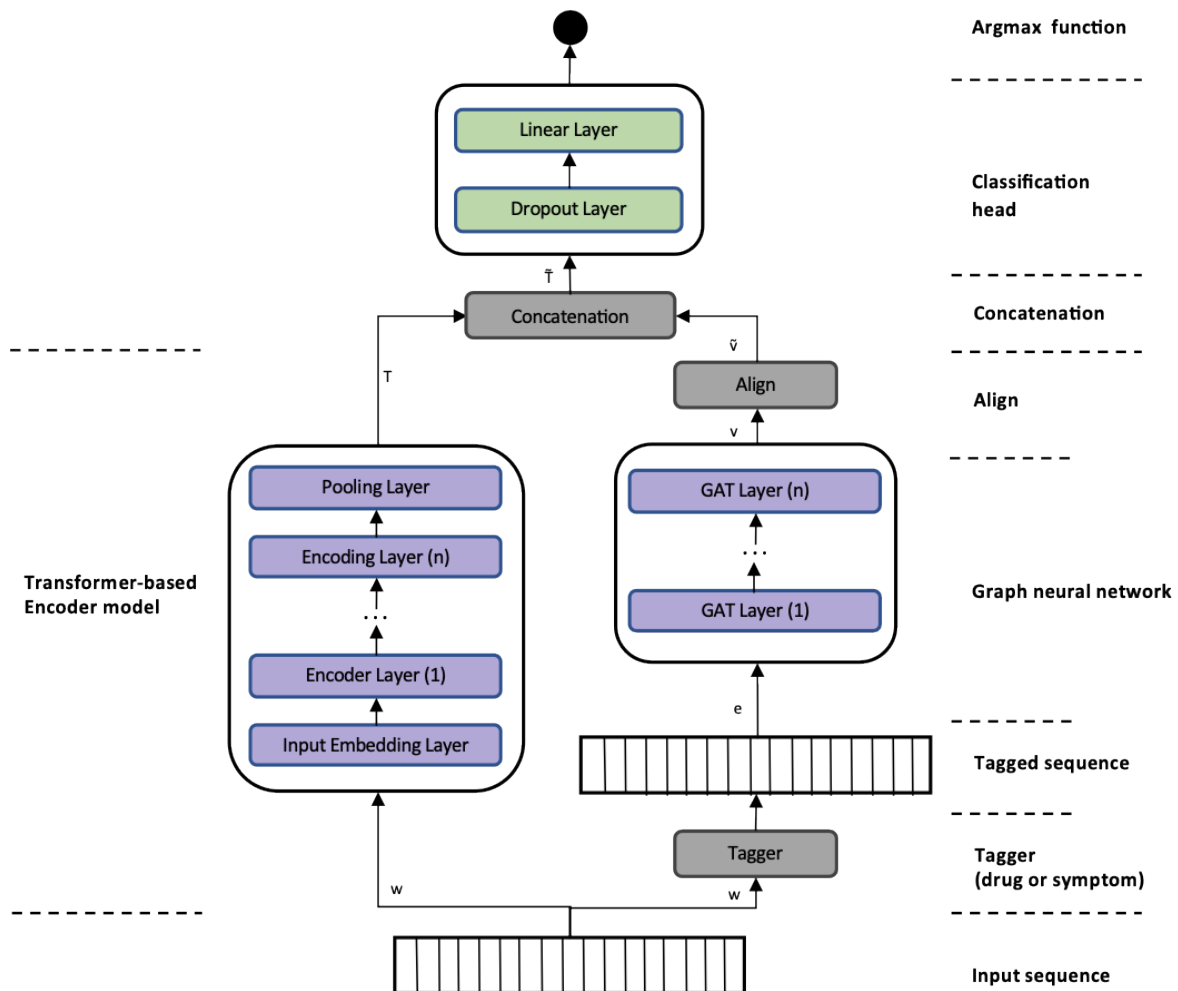
284 classes of the ontology may lack valuable information due to a lack of class annotations. As a
285 result, nodes can assign lower weights to neighbors without valuable information due to the
286 attention mechanism [27]. The aforementioned methodology of generating knowledge graph
287 embeddings corresponds to what Yang et al. [28] refer to as cascaded model architecture. In
288 this architecture, initial node features are generated using language models and then further
289 processed by GNNs [28].

290 Integrating Transformer-based Models with GNNs

291 We propose a knowledge fusion model to combine node embeddings learned via a graph neural
292 network with a transformer-based model. We begin by taking an input sentence and using a
293 rule-based tagger to identify symptoms and/or drug entities depending on the given knowledge
294 graph. The KG can be either SYMP, DRUGO, or a combination of both. The tagged input
295 sequence has the same length as the original input sequence but holds additional information
296 for those input tokens that were tagged by the rule-based annotator. Furtheron, the tagged input
297 sequence is passed through a GNN and returns a vector that holds zeroes for tokens that do
298 not belong to any tagged entity and the corresponding node embedding for tokens that were
299 tagged by the previous tagger. Subsequently, the resulting vector v is aligned with the
300 representation of the transformer (by adding zeros wherever a padding token was added or
301 where words were split into word pieces). This aligned vector \tilde{v} is then concatenated with T to
302 create a final knowledge-enriched representation \tilde{T} of the input sequence. This final
303 representation is further passed into a linear layer, which serves as the classification head (Fig.
304 2).

305
306 Additionally, we set the GNN weights as fixed by default, resulting in the usage of GNN as a
307 lookup table within the underlying embedding space. We refer to this architecture as a graph
308 concat model. Nevertheless, we have implemented an additional model variant called the graph
309 concat adaptive weights model, in which we treat the GNN weights as trainable parameters that

310 are adjusted during the training of the entire model. Fig. 2 illustrates the architecture of the graph
 311 concat model without (orange box) and with (purple box) adaptive GNN weights.
 312
 313 Furthermore, instead of using the entire graph as in the setting presented above, we explored
 314 an additional GNN configuration where only a subgraph of the knowledge graph is used and
 315 passed through the GNN. This subgraph is constructed from the k-hop neighborhood of the
 316 tagged entity. Finally, instead of concatenating the node representation to the transformer
 317 representation, a graph pooling layer (concatenation of global max and average pooling) is
 318 added and its output is concatenated to the transformer representations. The just presented
 319 architecture will be noted as graph concat (graph concat AW for adaptive GNN weights and
 320 graph concat AWS for graph concat with adaptive weights and subgraph modification) from now
 321 on.



322 Fig. 2. The architecture of the graph concat model with fixed and trainable GNN weights.
 323

324 Compared Method: ERNIE

325 Enhanced Language Representation with Informative Entities (ERNIE) is a fusion model
326 introduced by Zhang et al. [17]. However, ERNIE handles the knowledge injection differently
327 than other models. Instead of calculating the representation of context knowledge within ERNIE
328 itself, it is computed separately. In ERNIE, TransE is utilized to generate and retrieve
329 embeddings for the knowledge. To have a fair comparison, we also adopted this approach in
330 our work. For a more detailed explanation of the working principle of ERNIE, we refer to the
331 original study published by Zhang et al. [17].

332 The implementation used in this work is obtained from the GitHub repository
333 <https://github.com/thunlp/ERNIE> , which provides a pre-trained ERNIE model.

334 Experimental Setup and Training Strategy

335 To perform an unbiased final evaluation on a completely independent test set, we randomly
336 chose and reserved 20% from each dataset. The remaining 80% of each dataset was divided
337 into a train and validation set, with a ratio of 4-to-1. This means that 64% of the entire dataset
338 served as a training set used to train the model, while the remaining 16% was used as the
339 validation set hyperparameter tuning. After hyperparameter tuning, we trained the final model
340 by combining both training and validation sets, which were used to evaluate the performance of
341 the aforementioned independent test set. Furthermore, to have maximum comparability along
342 all the different model architectures, those splits were consistently applied throughout all
343 experiments.

344

345 Each experiment conducted in our study was constructed from the four categories listed in Table
346 2. The categories encompass the model architecture, the pre-trained transformer-based
347 language model, the ADE text corpus, and the contextual knowledge resource. The selected
348 model architectures further categorize the results into ERNIE, graph concat model with fixed
349 GNN weights, graph concat model with adaptive GNN weights, and graph concat model with
350 adaptive GNN weights and k-hop subgraph. Additionally, baseline experiments are considered

351 as a separate category that only uses the pre-trained language models BERT and BioBERT. It
352 is important to note that we utilized BERT as a general language model to assess the
353 performance achievable by a transformer-based encoder that was not pre-trained on domain-
354 specific documents. On the other hand, BioBERT is a domain-specific model that was pre-
355 trained on biomedical documents [25]. All models were evaluated on all five ADE text corpora.
356 Finally, each model was equipped with either contextual knowledge about drugs, symptoms, or
357 both. In addition to the 10 baseline experiments, the various options for experiment
358 configurations resulted in a total of 115 experiments.

359

Experiment categories	Values
Knowledge fusion model architecture	ERNIE (not used in combination with pre-trained transformer), Graph concat model with fixed GNN weights, Graph concat model with adaptive GNN weights, and Graph concat model with adaptive GNN weights and k-hop subgraph
Pre-trained language model	BERT and BioBERT
ADE corpora	SMM4H, CADEC, PSYTAR, TAC, and ADE
Knowledge resource	SYMP, DRUGO, and DRUGO + SYMP

360 Table 2: Overview of experiment categories. Their combination results in a total of 115
361 experiments in addition to 10 baseline experiments.

362

363 To ensure unbiased and comparable results, the same overall strategy for training, validation,
364 hyperparameter tuning, and testing was employed in each experiment. The optimal
365 hyperparameters were deduced by performing Bayesian hyperparameter optimization [29]. To
366 determine the optimal hyperparameters for each model, multiple models with different
367 hyperparameter configurations were trained on the training set. These models were then
368 evaluated on the validation set, to maximize the F_1 -Score. The cross-entropy loss function was
369 employed for all models in the context of NER. The AdamW [30] optimization algorithm was
370 chosen to adjust the model's weights during training. Finally, the optimal hyperparameters were

371 used to train models on the combination of training and validation sets. These new models were
372 then tested on held-out, independent test sets of each corpora.

373 Evaluation Scheme

374 We used the precision, recall, and F_1 measures to assess the performance of models. Each
375 dataset was labeled in the IOB scheme with which each token of a sequence is labeled either
376 as outside (O) of a named entity, as the beginning (B), or as an inside (I) token of a named
377 entity. Hence, the classification head of each of the models had three output neurons and the
378 NER problem was formulated as a classification task with three classes. However, we are
379 interested in ADE spans that can consist of multiple tokens, therefore, for the final evaluation
380 the IOB labeling was discarded, and the sequences were aggregated into real ADE mentions.
381 The final scores were then calculated by taking into account the exact overlap of the full spans
382 of ADE mentions.

383 Implementation

384 The experiments conducted in this study were implemented using PyTorch and PyTorch
385 Lightning. An essential component are transformer-based models for which we used the
386 Huggingface transformers library. To perform hyperparameter tuning Optuna was chosen as the
387 library. Finally, for processing and handling the considered datasets we used Pandas and
388 Spacy. The baseline models as well as the graph concat model experiments are using BERT
389 and BioBERT, which come in different sizes and configurations. We used uncased BERT,
390 commonly known as 'bert-base-uncased', which contains a total of 110M parameters. The
391 BioBERT model is specified as 'dmis-lab/biobert-v1.1', which has the equivalent amount of
392 parameters as 'bert-base-uncased'. The model training and testing was performed using
393 Nvidia's V100 and A100 GPUs.

Results

394
395
396
397
398
399
400
401
402
403
404
405
406
407

We evaluated the aforementioned five different model architectures¹ on each of the ADE datasets. Table 3 provides an overview of the final evaluation results providing the F1-score obtained by applying the models within a certain configuration on the independent test sets. Here the configuration refers to the choice of context knowledge resource and underlying transformer-based model, where applicable. Please take note that the graph concat k-hop subgraph experiments were omitted from Table 1 since this architecture did not achieve the top ranking on any of the corpora. For a comprehensive overview of results including this architecture as well as precision and recall measures for all models, we refer to Supplementary Table 1.

Model	Knowledge resource	F ₁ (in %) on ADE Corpora				
		CADEC	SMM4H	PSYTAR	ADE	TAC
BERT	-	71.84	62.30	70.02	75.37	92.06
BioBERT	-	70.81	61.95	68.80	79.42	93.87
ERNIE + TransE	DRUG	71.84	63.23	70.63	75.44	92.57
ERNIE + TransE	DRUGO_SYMP	69.32	61.76	70.95	76.04	92.55
ERNIE + TransE	SYMP	68.70	61.32	71.40	76.58	92.06
Graph concat + BERT	DRUG	70.45	62.65	71.38	79.79	93.80
Graph concat + BERT	DRUGO_SYMP	70.70	65.16	72.32	78.84	93.49
Graph concat + BERT	SYMP	71.05	62.83	72.03	78.13	93.87
Graph concat + BioBERT	DRUG	70.28	61.51	68.24	76.73	94.15
Graph concat + BioBERT	DRUGO_SYMP	69.57	62.75	70.05	78.90	93.88
Graph concat + BioBERT	SYMP	69.40	62.48	69.32	78.59	93.31
Graph concat AW + BERT	DRUG	70.55	63.96	72.50	79.03	93.02
Graph concat AW + BERT	DRUGO_SYMP	71.82	63.99	71.38	79.54	93.87
Graph concat AW + BERT	SYMP	70.59	64.22	72.02	78.4	93.22
Graph concat AW + BioBERT	DRUG	71.23	61.05	70.08	78.11	93.75
Graph concat AW + BioBERT	DRUGO_SYMP	68.87	62.12	69.62	78.62	93.78
Graph concat AW + BioBERT	SYMP	70.16	57.78	69.00	76.01	93.45

¹ Baseline, ERNIE, Graph concat, Graph concat AW, Graph concat AWSUB

408 Table 3: Final evaluation results on test set from all experiments. F_1 stands for F_1 -score. All
409 scores are strict scores and given in %. The best score on each corpus is given in bold.
410 AW=adaptive weights

411
412 When examining the results on the CADEC corpus (Table 3), one may observe that the best
413 baseline experiment utilizing BERT already demonstrates a strong performance in terms of the
414 F_1 -score (71.84%). None of the other models evaluated on CADEC were able to improve upon
415 this score. However, ERNIE equipped with contextual knowledge about drugs achieved the
416 same score of 71.84%. Additionally, the graph concat AW model incorporated with drugs and
417 symptom knowledge came quite close with a score of 71.82%.

418
419 The performance of the models on the SMM4H corpus, in general, was lower than on all other
420 corpora. The difference of performance can already be observed in the results of the baseline
421 experiments that show a noticeable gap of almost 8-20% points. Furthermore, ERNIE, equipped
422 with prior knowledge about drugs, was able to perform better on SMM4H with an F_1 -score of
423 63.23% than the best baseline experiment using BERT, which reached an F_1 -score of 62.3%.
424 Moreover, the graph concat AW model with contextual knowledge about symptoms adopting
425 BioBERT as the underlying transformer was also able to report better F_1 -scores (64.22%) than
426 the baseline experiments and better than the best-performing ERNIE model (Table 3). Finally,
427 the graph concat model with fixed GNN weights using BERT as its underlying pre-trained
428 transformer while equipped with joint prior knowledge about symptoms and drugs reports the
429 overall best score on SMM4H with an F_1 -score of 65.16%.

430
431 On PSYTAR, the ERNIE model equipped with prior knowledge about symptoms, reaching an
432 F_1 -score of 71.40%, was able to slightly improve the performance of the BERT baseline
433 experiment that was able to achieve an F_1 -score of 70.02%. The graph concat model using
434 BioBERT and drugs and symptoms knowledge was able to improve this score to 72.32% F_1 -
435 score. The graph concat AW model with BERT and the drug knowledge graph further improves
436 this score to 72.50% F_1 -score.

437

438 On the ADE corpus, the ERNIE model was not able to reach the score reported by the best
439 baseline model BioBERT (79.42% F_1 -score). However, the graph concat AW model using BERT
440 and adopting prior knowledge about drugs and symptoms was able to slightly increase this score
441 to 79.54% F_1 -score. The graph concat model with fixed GNN weights while also using BERT
442 as its transformer and equipped with prior knowledge about drugs further improved this score
443 to 79.79% F_1 -score.

444

445 Finally, on the TAC corpus, all models considered in the results were able to score F_1 -scores
446 above 90%. The best baseline model, BioBERT, was able to reach an F_1 -score of 93.87%. The
447 ERNIE and the graph concat AW model were not able to outperform the best baseline model.
448 However, the graph concat model with fixed GNN weights using BioBERT as its transformer
449 and equipped with contextual knowledge about drugs was able to increase upon the baseline
450 performance achieving the highest F_1 -score of 94.15% on TAC corpus.

451

452 We performed an additional analysis to determine the different attributes of each of the 5 corpora
453 that could shed some light on explaining the modeling performance. Table 4 depicts the results
454 of this corpus analysis comprising three measures. Firstly, the wordpiece diversity, which was
455 assembled by counting how many unique wordpieces could be found in each sentence of a
456 corpus normalized by the total amount of wordpieces in a sentence. The second measure
457 calculates the sentence length on wordpiece level and the number of hits in the DRUGO_SYMP
458 knowledge graph. A hit is defined as an entity in the sentence corresponding to a node in the
459 knowledge graph. All values presented in Table 4 are averaged over all sentences in the
460 corresponding corpus. The CADEC corpus is a clear outlier in terms of the mean number of KG
461 hits, the mean sentence length, and wordpiece diversity. CADEC is the only corpus where we
462 did not observe any advantage of using a knowledge fusion model in terms of F_1 -score.

463

Corpus	Vocabulary/Model	Mean wordpiece diversity	Mean sentence length (in wordpieces)	Mean number of KG hits	Difference of best model to baseline (in F_1 % points)
--------	------------------	--------------------------	--------------------------------------	------------------------	--

CADEC	BERT	0.74	113.87	4.81	-0.02 ↓
CADEC	BioBERT	0.75	121.42	4.81	-0.02 ↓
ADE	BERT	0.94	33.21	1.45	0.12 ↑
ADE	BioBERT	0.94	35.56	1.45	0.12 ↑
PSYTAR	BERT	0.93	22.74	1.13	2.48 ↑
PSYTAR	BioBERT	0.93	23.75	1.13	2.48 ↑
TAC	BERT	0.81	47.36	1.85	0.28 ↑
TAC	BioBERT	0.82	52.93	1.85	0.28 ↑
SMM4H	BERT	0.91	30.46	1.30	2.86 ↑
SMM4H	BioBERT	0.91	31.81	1.30	2.86 ↑

464 Table 4: Corpora characterization in terms of average wordpiece diversity, average sentence
465 length, and average number of knowledge graph hits.

466 Discussion

467 Extracting meaningful insights about ADEs from unstructured text offers the chance to enhance
468 our knowledge of ADEs and in the long run contributes to drug safety. Specifically, the extraction
469 of ADEs from patient-reported texts allows for gathering great amounts of negative drug
470 experiences since vast amounts of data are published every day on social media. In our work,
471 we evaluate various knowledge fusion modeling approaches on the ADE extraction task using
472 five relevant text corpora, namely CADEC, SMM4H, PSYTAR, TAC, and ADE. Additionally, we
473 utilized a rich knowledge base in terms of drugs and symptoms, which provided valuable
474 contextual knowledge to these models. Knowledge graph embeddings derived from GNNs have
475 ensured a knowledge representation well suited for the fusion with linguistic representations
476 obtained using transformer-based large language models. The final results on independent test
477 sets showed that using models with contextual knowledge can help to gain performance on ADE
478 corpora.

479
480 We observed a significant variation in performance scores and model behavior across different
481 datasets. There was no clear advantage of adopting a knowledge fusion methodology over the

482 baseline model BERT on the CADEC dataset. Using graph concat adaptive weights model
483 resulted in an F_1 -score quite similar to the BERT and ERNIE model. However, on the SMM4H
484 corpus, we observed a performance increase from top-scoring baseline (BERT) to ERNIE to the
485 graph concat model. BERT reached an F_1 -score of 62.30%, and equipping it with contextual
486 knowledge about drugs and symptoms raised this score to 65.16%. When examining the results
487 for PSYTAR, the top-performing baseline model (BERT) achieved an F_1 -score of 70.02% for
488 extracting ADE entities. ERNIE was able to improve this score by approximately 1.5%. By
489 enabling BERT to utilize contextual knowledge about drugs through the graph concat
490 architecture, the score further increased to 72.5%. When considering the ADE corpus, there
491 was a notable difference in scores between baseline models (75.37% for BERT and 79.42% for
492 BioBERT). None of the ERNIE models were able to match the baseline score achieved by
493 BioBERT. However, the graph concat model with fixed GNN weights that utilizes BERT and
494 contextual knowledge about drugs was able to slightly increase the baseline performance to a
495 79.79% F_1 -score. Similarly, in the case of the TAC dataset, BioBERT was able to reach a high
496 F_1 -score of 93.87% that was not surpassed by any ERNIE model. The graph concat model was
497 able to slightly increase the baseline performance on TAC to an F_1 -score of 94.15%.

498
499 There was no clear indication of whether the graph concat models work better with BERT or
500 BioBERT as the underlying transformer model. However, we observed that on CADEC, utilizing
501 BioBERT in knowledge fusion could improve the baseline BioBERT performance (BioBERT:
502 70.81% F_1 and 71.23% F_1 graph concat with adaptive GNN weights and contextual knowledge
503 about drugs), whereas this could not be observed for BERT (71.84% F_1 is best score on
504 CADEC). When considering the usefulness of knowledge resources, it is noteworthy to mention
505 that all models that outperformed the baseline experiments relied either on DRUGO or
506 DRUGO_SYMP contextual knowledge. Based on this observation, it suggests that contextual
507 knowledge about drugs may hold greater importance for the knowledge fusion models
508 compared to knowledge about symptoms. The trend was apparent in both the graph concat
509 model and ERNIE.

510

511 As mentioned, our observations indicate that the effectiveness of knowledge fusion models
512 varies across different corpora. We did not observe any performance improvement using
513 knowledge fusion models on the CADEC corpus. This aligns with the findings in Table 4, which
514 highlights CADEC being an outlier in the textual analysis in terms of wordpiece diversity,
515 sentence length, and KG hits. Further investigation is necessary to determine the causal
516 relationship between these metrics and the potential improvement of pure linguistic models with
517 knowledge fusion. However, based on our interpretation of the results, it can be reasoned that
518 knowledge fusion models are most beneficial for relatively short text, such as postings found in
519 SMM4H and PSYTAR (<24 wordpieces on average in PSYTAR and <32 in SMM4H). Notably,
520 the CADEC corpus stands out in terms of the number of hits in the knowledge graph. This
521 suggests that an excessive amount of contextual knowledge may not contribute positively to the
522 model's accuracy. Liu et al. [19] introduced the concept of knowledge noise (KN), which refers
523 to the phenomenon that an excess of context can disrupt the original meaning of the sentence.
524 However, further investigation is needed to find whether during knowledge fusion KN played a
525 role in the lack of performance improvement on CADEC. Additionally, since PSYTAR and
526 SMM4H are derived from Twitter, it is reasonable to assume that these corpora deviate from
527 formal, scientific English. In this context, knowledge fusion can potentially compensate for the
528 informality in language and for the lack of linguistic context by providing valuable information on
529 specific ADEs.

530

531 The current workflow infuses context knowledge into models for the words that are identified as
532 drugs or symptoms by a rule-based NER tagger. For this purpose, we preferred a rule-based
533 system to avoid false positives in terms of context knowledge. However, a more advanced
534 machine learning-based tagger with a better performance may produce even higher results,
535 which we will explore in our future work. One possible machine learning-based model for such
536 an approach would be Med7 [31], which reports good results in terms of F_1 -score on the task of
537 extracting drug entities from text. Although the used knowledge resources have shown

538 performance gains while using the knowledge fusion approach, they are far from being complete
539 and perfect. Encoding even more knowledge about drugs and symptoms could improve the
540 current models of ADE detection.

541

542 Although this study performed a comprehensive analysis, it is important to note existing
543 limitations. Further knowledge fusion approaches such as K-BERT, K-Adapters, or SKILL [32–
544 34] are worth exploring in future experiments for evaluating knowledge fusion models on the
545 ADE extraction task. Some of the training datasets used in this work comprise only a relatively
546 small number of postings, around 1,000 for both the SMM4H and CADEC corpora. It is well-
547 known that deep learning-based NLP models generally tend to perform better when trained on
548 larger datasets. Therefore, to further enhance the performance of the knowledge fusion models
549 employed in this study, having access to large and diverse corpora of patient-reported texts that
550 include annotated ADE entities, particularly in the style of CADEC, would be beneficial.
551 Consequently, future efforts should be directed toward creating, collecting, and annotating a
552 comprehensive ADE corpus of diverse texts, which could contribute to the advancement of this
553 research.

554 Conclusion

555 The presented work elaborates on the approach to enriching transformer models such as BERT
556 and its relative, BioBERT, with contextual knowledge about the texts fed into them. Two types
557 of prior knowledge on drugs and symptoms were considered in this work. The drug knowledge
558 resource provides rich, structured knowledge about drugs and their working principles and was
559 especially created for this work. We conducted a great number of experiments and reported the
560 combinations of transformer models, knowledge fusion architectures, and context knowledge
561 that yielded the highest F_1 -scores. The presented results allow the conclusion that contextual
562 knowledge encoded suitably and provided to a transformer model is a valid approach to improve
563 performance in an NER task scenario. Also observable is that this prior knowledge is especially
564 of great use when the data at hand is rather unstructured and composed of short texts as is the

565 case in the SMM4H and PSYTAR corpus. Finally, one can conclude that knowledge resources
566 that provide well-structured domain knowledge, encoded as knowledge graphs respectively
567 ontologies can provide valuable context for transformer models. Graph neural networks have
568 shown to be a well-suited method to derive a numerical representation of the ontologies used in
569 this work capable of being concatenated with the linguistic representation created by a
570 transformer model. The architecture of the graph concat model with and without adaptive GNN
571 weights implemented in this work has shown to be advantageous compared to pure
572 transformers (BERT and BioBERT) as well as to another, well-established, knowledge fusion
573 model, ERNIE. Hence, that architecture deserves additional development to further improve its
574 performance on tasks such as ADE extraction in structured and unstructured texts. Huge
575 potential lies in the idea of fusing large language models with appropriate domain knowledge
576 and definitively deserves further research that includes whether the presented approach
577 generalizes on tasks further than detecting adverse drug events in texts.

578 Availability

579 The code is available in the repository: [https://github.com/SCAI-BIO/adr-detection-with-](https://github.com/SCAI-BIO/adr-detection-with-knowledge-fusion)
580 [knowledge-fusion](https://github.com/SCAI-BIO/adr-detection-with-knowledge-fusion).

581 Funding

582 The author(s) received no specific funding for this work.

583 Competing interests

584 The authors have declared that no competing interests exist.

585 Figure Captions

586 **Fig. 1.** CADEC and SMM4H example phrases that are enriched with contextual knowledge
587 about drugs and symptoms. The sentence from CADEC “Very drowsy and tired and no pain
588 relief at all.” can be equipped with symptom classes such as Drowsiness and Tiredness, which

589 are subclasses of Neurological and physiological symptom class, as well as Pain, which is a
590 subclass of Nervous system symptom.

591

592 **Fig. 2.** The architecture of the graph concat model with fixed and trainable GNN weights.

593

594 References

595 1. Jain H, Raj N, Mishra S (2021) A Sui Generis QA Approach using RoBERTa for Adverse
596 Drug Event Identification. *BMC Bioinformatics* 22:330

597 2. Devlin J, Chang M-W, Lee K, Toutanova K (2018) BERT: Pre-training of Deep
598 Bidirectional Transformers for Language Understanding.
599 <https://doi.org/10.48550/ARXIV.1810.04805>

600 3. Radford A, Narasimhan K, Salimans T, Sutskever I (2018) Improving language
601 understanding by generative pre-training. *OpenAI*

602 4. Ge S, Wu F, Wu C, Qi T, Huang Y, Xie X (2020) FedNER: Privacy-preserving Medical
603 Named Entity Recognition with Federated Learning.
604 <https://doi.org/10.48550/ARXIV.2003.09288>

605 5. Karimi S, Metke-Jimenez A, Kemp M, Wang C (2015) Cadec: A corpus of adverse drug
606 event annotations. *J Biomed Inform* 55:73–81

607 6. Magge A, Klein A, Miranda-Escalada A, et al (2021) Overview of the Sixth Social Media
608 Mining for Health Applications (#SMM4H) Shared Tasks at NAACL 2021. In: *Proc. Sixth*
609 *Soc. Media Min. Health SMM4H Workshop Shar. Task. Association for Computational*
610 *Linguistics, Mexico City, Mexico*, pp 21–32

611 7. Zolnoori M, Fung KW, Patrick TB, et al (2019) The PsyTAR dataset: From patients
612 generated narratives to a corpus of adverse drug events and effectiveness of psychiatric
613 medications. *Data Brief* 24:103838

614 8. Tonning et al. KR (2017) Overview of the TAC 2017 Adverse Reaction Extraction from
615 Drug Labels Track.

616 9. Gurulingappa H, Rajput AM, Roberts A, Fluck J, Hofmann-Apitius M, Toldo L (2012)
617 Development of a benchmark corpus to support the automatic extraction of drug-related
618 adverse effects from medical case reports. *J Biomed Inform* 45:885–892

619 10. Sboev A, Selivanov A, Rylkov G, Rybka R (2021) On the accuracy of different neural
620 language model approaches to ADE extraction in natural language corpora. *Procedia*
621 *Comput Sci* 190:706–711

622 11. Portelli B, Lenzi E, Chersoni E, Serra G, Santus E (2021) BERT Prescriptions to Avoid
623 Unwanted Headaches: A Comparison of Transformer Architectures for Adverse Drug
624 Event Detection. In: *Proc. 16th Conf. Eur. Chapter Assoc. Comput. Linguist. Main Vol.*
625 *Association for Computational Linguistics, Online*, pp 1740–1747

626 12. Ramesh S, Tiwari A, Choubey P, Kashyap S, Khose S, Lakara K, Singh N, Verma U
627 (2021) BERT based Transformers lead the way in Extraction of Health Information from

- 628 Social Media. In: Proc. Sixth Soc. Media Min. Health SMM4H Workshop Shar. Task.
629 Association for Computational Linguistics, Mexico City, Mexico, pp 33–38
- 630 13. Raval S, Sedghamiz H, Santus E, Alhanai T, Ghassemi M, Chersoni E (2021) Exploring
631 a Unified Sequence-To-Sequence Transformer for Medical Product Safety Monitoring in
632 Social Media. In: Find. Assoc. Comput. Linguist. EMNLP 2021. Association for
633 Computational Linguistics, Punta Cana, Dominican Republic, pp 3534–3546
- 634 14. Haq HU, Kocaman V, Talby D (2022) Mining Adverse Drug Reactions from Unstructured
635 Mediums at Scale. <https://doi.org/10.48550/arXiv.2201.01405>
- 636 15. Miftahutdinov Z, Tutubalina E (2019) Deep Neural Models for Medical Concept
637 Normalization in User-Generated Texts. In: Proc. 57th Annu. Meet. Assoc. Comput.
638 Linguist. Stud. Res. Workshop. pp 393–399
- 639 16. Stanovsky G, Gruhl D, Mendes PN (2017) Recognizing Mentions of Adverse Drug
640 Reaction in Social Media Using Knowledge-Infused Recurrent Models. Proc. 2017 Conf.
641 Eur. Chapter Assoc. Comput. Linguist.
- 642 17. Zhang Z, Han X, Liu Z, Jiang X, Sun M, Liu Q (2019) ERNIE: Enhanced Language
643 Representation with Informative Entities. <https://doi.org/10.48550/ARXIV.1905.07129>
- 644 18. Choi E, Levy O, Choi Y, Zettlemoyer (2018) Ultra-Fine Entity Typing. Proc. ACL
- 645 19. Liu W, Zhou P, Zhao Z, Wang Z, Ju Q, Deng H, Wang P (2019) K-BERT: Enabling
646 Language Representation with Knowledge Graph.
647 <https://doi.org/10.48550/ARXIV.1909.07606>
- 648 20. Schriml LM, Arze C, Nadendla S, et al (2010) GeMInA, Genomic Metadata for Infectious
649 Agents, a geospatial surveillance pathogen database. *Nucleic Acids Res* 38:D754-764
- 650 21. Wishart DS, Feunang YD, Guo AC, et al (2018) DrugBank 5.0: a major update to the
651 DrugBank database for 2018. *Nucleic Acids Res* 46:D1074–D1082
- 652 22. Anatomical therapeutic chemical (Atc) classification - [https://www.who.int/tools/atc-ddd-](https://www.who.int/tools/atc-ddd-toolkit/atc-classification)
653 [toolkit/atc-classification](https://www.who.int/tools/atc-ddd-toolkit/atc-classification).
- 654 23. Mayers, Michael, Steinecke, Dylan, Su, Andrew I. (2020) Database of mechanism of
655 action paths for selected drug-disease indications.
656 <https://doi.org/10.5281/ZENODO.3708278>
- 657 24. Devlin J, Chang M-W, Lee K, Toutanova K (2019) BERT: Pre-training of Deep
658 Bidirectional Transformers for Language Understanding. In: Proc. 2019 Conf. North Am.
659 Chapter Assoc. Comput. Linguist. Hum. Lang. Technol. Vol. 1 Long Short Pap. pp 4171–
660 4186
- 661 25. Lee J, Yoon W, Kim S, Kim D, Kim S, So CH, Kang J (2019) BioBERT: a pre-trained
662 biomedical language representation model for biomedical text mining. *CoRR*
663 [abs/1901.08746](https://arxiv.org/abs/1901.08746):
- 664 26. Bordes A, Usunier N, Garcia-Duran A, Weston J, Yakhnenko O (2013) Translating
665 Embeddings for Modeling Multi-relational Data. *Adv. Neural Inf. Process. Syst.* 26:
- 666 27. Veličković P, Cucurull G, Casanova A, Romero A, Liò P, Bengio Y (2017) Graph
667 Attention Networks. <https://doi.org/10.48550/ARXIV.1710.10903>

- 668
669
28. Yang J, Liu Z, Xiao S, Li C, Sun G, Xie X (2021) GraphFormers: GNN-nested Language Models for Linked Text Representation. CoRR abs/2105.02605:
- 670
671
672
29. Bergstra J, Yamins D, Cox D (2013) Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In: Int. Conf. Mach. Learn. PMLR, pp 115–123
- 673
674
30. Loshchilov I, Hutter F (2017) Decoupled Weight Decay Regularization. <https://doi.org/10.48550/ARXIV.1711.05101>
- 675
676
677
31. Kormilitzin A, Vaci N, Liu Q, Nevado-Holgado A (2020) Med7: a transferable clinical natural language processing model for electronic health records. ArXiv Prepr. ArXiv200301271
- 678
679
32. Liu W, Zhou P, Zhao Z, Wang Z, Ju Q, Deng H, Wang P (2019) K-BERT: Enabling Language Representation with Knowledge Graph. ArXiv190907606 Cs
- 680
681
33. Moiseev F, Dong Z, Alfonseca E, Jaggi M (2022) SKILL: Structured Knowledge Infusion for Large Language Models. <https://doi.org/10.48550/arXiv.2205.08184>
- 682
683
684
34. Wang R, Tang D, Duan N, Wei Z, Huang X, Ji J, Cao G, Jiang D, Zhou M (2020) K-Adapter: Infusing Knowledge into Pre-Trained Models with Adapters. <https://doi.org/10.48550/arXiv.2002.01808>
- 685