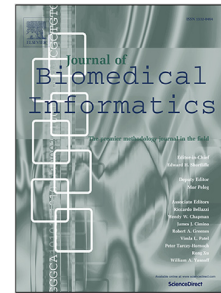


Journal Pre-proof

GERNERMED++: Semantic annotation in German medical NLP through transfer-learning, translation and word alignment

Johann Frei, Ludwig Frei-Stubber, Frank Kramer



PII: S1532-0464(23)00234-4
DOI: <https://doi.org/10.1016/j.jbi.2023.104513>
Reference: YJBIN 104513

To appear in: *Journal of Biomedical Informatics*

Received date: 14 November 2022
Revised date: 27 September 2023
Accepted date: 4 October 2023

Please cite this article as: J. Frei, L. Frei-Stubber and F. Kramer, GERNERMED++: Semantic annotation in German medical NLP through transfer-learning, translation and word alignment, *Journal of Biomedical Informatics* (2023), doi: <https://doi.org/10.1016/j.jbi.2023.104513>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

GERNERMED++: Semantic Annotation in German Medical NLP through Transfer-Learning, Translation and Word Alignment

Johann Frei^{a,*}, Ludwig Frei-Stuber^b, Frank Kramer^a

^a *IT-Infrastructure for Translational Medical Research, University of Augsburg
Alter Postweg 101, 86159 Augsburg, Germany
firstname.lastname@informatik.uni-augsburg.de*

^b *Institute and Outpatient Clinic for Occupational, Social and Environmental Medicine
80336 Munich, Germany
firstname.lastname@med.uni-muenchen.de*

Abstract

We present a statistical model, GERNERMED++, for German medical natural language processing trained for named entity recognition (NER) as an open, publicly available model. We demonstrate the effectiveness of combining multiple techniques in order to achieve strong results in entity recognition performance by the means of transfer-learning on pre-trained deep language models (LM), word-alignment and neural machine translation, outperforming a pre-existing baseline model on several datasets. Due to the sparse situation of open, public medical entity recognition models for German texts, this work offers benefits to the German research community on medical NLP as a baseline model. The work serves as a refined successor to our first GERNERMED model. Similar to our previous work, our trained model is publicly available to other researchers. The sample code and the statistical model is available at:

<https://github.com/frankkramer-lab/GERNERMED-pp>

Keywords: natural language processing, medical NLP, medical named entity recognition, transfer learning, German NLP, artificial intelligence

*Corresponding author

1. Introduction

Extraction and processing of key information from medical notes and doctors' letters pose a common challenge in the advanced digitization of healthcare systems. In particular, research-oriented data mining of non-research-centric data sources (often referred to as *second use*) often requires expensive data harmonization processes in order to transform unstructured or semi-structured data into strictly structured, uniform data representations such as HL7 or FHIR. While manually solving these processes can be carried out for document analysis on certain studies, it is rendered impractical for large-scale text analysis on legacy data or processing day-to-day clinical data.[42, 47]

Handling heterogeneous data from text-based documents is a central subject of natural language processing. In recent years deep learning-inspired approaches have been applied successfully to tackle various NLP tasks effectively. However, training deep language models requires proper datasets in regard to aspects like corpus size, annotation work, data diversity and overall dataset quality, in order to retrieve well-performing models. In medical NLP, obtaining such annotated datasets remains rather difficult for various reasons.[7] For instance, the use and publication of medical data is highly restricted for the reasons of privacy and country-dependent data protection legislation.[7] Even though medical datasets have been published in English, such datasets for German texts in contrast are still frequently unavailable to external researchers.[42]

In this paper, we propose an approach of combining multiple ideas to obtain a German medical NLP model, which we refer to as *GERNERMED++* and which serves as a successor to our previous *GERNERMED*[13] model:

- **Translation:** The state of German medical corpora is limited and the use of internal datasets for training and publication of such models is legally unclear. In contrast, medical datasets in English have already been published and therefore, neural machine translation (NMT) can be applied to obtain German data from English datasets.
- **Annotation Projection:** Annotation of large corpora is crucial for su-

pervised learning and determines the quality of the final performance of the model. However the cost of obtaining gold-standard annotations from scratch is prohibitively expensive. Given our set of NMT-based German data, word alignment estimation can be used to project token-level annotations from English data to German data without manual intervention.

- **Transfer-Learning through Model Fine-Tuning:** To further improve the downstream performance of the NLP model under the constraints of our small, task-specific dataset, a larger, pre-trained German LM is used for advanced semantic, context-aware feature extraction and further fine-tuning.

Our method and our results highlight the effectiveness of non-German data sources for training a German NER model for medical semantic annotation such as medication detection. Our model can surpass the performance of the prior German NLP model GGPONC 2[4] which is traditionally trained on German text data. In principle, the method is not inherently limited to German because NMT and word alignment techniques also exist for several other languages and therefore, it could be applied to other languages as well.

1.1. Related Work

In the recent decade, in particular in the last five years, the field of natural language processing has been radically transformed by the use of data-driven, neural methods that are able to surpass previous state-of-the-art performances.[43, 47] This development is likewise reflected by several empirical facts such as quantity of published research or project funding[47]. The introduction of the attention-based transformer model[45] in the field of NLP led to various follow-up works such as BERT[9] and similar deep language models that are trained and applied on domain-specific contexts[33, 36, 27, 1, 2, 28]. All these domain-specific works share in common that their research focus lies primarily on English application and use.

The training of novel transformer-based German NLP models requires large, well-suited datasets with respect to size and quality. In purely supervised scenarios, this also includes the need for gold-standard annotation labels. While several works with internal datasets exist, their datasets are not shared among the research community and remain undisclosed[46, 12, 5, 44, 24, 39, 8, 23, 17, 30, 25, 41], and thus this presents major hurdles for open research and independent reproducibility. The situation on public, English datasets is more convenient and several large datasets like MIMIC-III[34] or the i2b2 challenges with datasets such as the n2c2 2018 dataset[18] have been published, as well as the multilingual Mantra GSC[22] dataset from the biomedical domain. Only in recent years has the German medical NLP research community addressed this issue and developed novel German medical datasets that are publicly accessible as foundation for future NLP work.[3, 21]. Regarding the GGPONC[3], an updated iteration has been presented[4].

With regards to novel German medical NLP systems, commercial software like *Averbis Health Discovery*[15]¹ and *German Spark NLP for Healthcare*[19]² are proprietary and require licenses. As an exception, *mEx*[37] is freely available, but the model weights can only be requested and used under data use agreement. An updated iteration has been presented as well[38]. For German medical NER tasks, only few public, open neural models are available to the best of our knowledge, such as *GGPONC*[4] and *GERNERMED*[13].

2. Methods

2.1. Dataset Acquisition

The dataset retrieval pipeline for German texts follows the approach proposed in GERNERMED[13]: As a starting point, the *2018 n2c2 shared task on ADE and medication extraction in EHR* dataset serves as an English source dataset of medical entities from anonymized electronic health records. The

¹<https://averbis.com/de/health-discovery/>

²https://nlp.johnsnowlabs.com/2021/03/31/ner_healthcare_de.html

Sample 1	
Raw	History of Present Illness: Ms. [**Known lastname 99778**] is a 41 year old female a history of warm autoantibody hemolytic anemia diagnosed...
Mask Replacement	History of Present Illness: Ms. Zahn is a 41 year old female a history of warm autoantibody hemolytic anemia...
Translated Sentence	Geschichte der gegenwärtigen Krankheit: Frau Zahn ist eine 41-jährige Frau, bei der eine hämolytische Anämie mit warmen Autoantikörpern diagnostiziert wurde,....

Sample 2	
Raw	Mr. [**Known lastname 1794**] was admitted from [**2185-4-23**] - [**2185-5-1**] for left sided chest pain...
Mask Replacement	Mr. Hartmann was admitted from 1985-04-13 - 2007-01-03 for left sided chest pain...
Translated Sentence	Herr Hartmann wurde von 1985-04-13 - 2007-01-03 wegen linksseitiger Brustschmerzen aufgenommen...

Figure 1: Effect of mask replacements on the English and German sentences for two exemplary samples.

English source dataset is decomposed into sentences as the initial preprocessing step. During that process, text spans that have been replaced with an anonymized identifier text bracket by the editors of the source dataset are detected and replaced with randomized synthetic data from the *Faker* Python module in order to reduce the number of irregular text occurrences while updating the initial annotation span indices accordingly. For instance, this includes text entities like first and family name, dates and postal addresses. For illustration purposes, two samples from the corpus are shown in Figure 1.

We apply the publicly available FAIRseq *transformer.wmt19.en-de* [32] NMT model for sentence-wise automatic translation, which features a transformer-based neural model for translating sentences from English to German. Since the annotation information from the English source dataset cannot be directly preserved for German sentences, the reconstruction of the annotation spans for the translated German sentences can be estimated by the means of a bitext word alignment as a postprocessing step. Artifacts in translation and alignment

have been discussed for GERNERMED[14]. In contrast to the approach in GERNERMED, we refine the word alignment estimation step in regard to the following aspects:

- **Improved Tokenization:** The tokenization of sentences for the word alignment differs from modern tokenizers that generate sub-word-level tokens optimized through techniques such as byte pair encoding schemes. Most word alignment methods operate on word-level tokenization with whitespace-based token splitting. In order to reduce the number of misaligned words, we further refined the word-level tokenization by separating punctuation from words instead of only relying on tokenization splits on whitespace characters. In our previous work[13], the projected German label spans often included trailing punctuation because a whitespace-based tokenization does not separate trailing punctuation from words and therefore, the label span reconstruction algorithm is unable to differentiate between words and punctuation within a token. This effect impedes subsequent model training but is countered by the improved, punctuation-aware tokenization.
- **Word Alignment Technique:** In NLP bitext word alignment is the task of determining the semantic correspondence between words from a bilingual sentence pair consisting of the source and translated sentence. In previous work, the *Fast-Align*[11] implementation has been used for establishing such correspondences. It uses the IBM 2 alignment model for alignment estimation in a purely unsupervised fashion. While there are also other models inspired by statistical machine translation[31, 48], recent work has been done towards neural approaches[20, 10]. For this work, we use the pre-trained model from *Awesome-Align*[10]. In short, the model tackles the task by encoding both sentences through a pre-trained cross-lingual language model in order to obtain contextualized word vector embeddings. Although the words of the sentence pairs largely differ with respect to their syntactic and linguistic features, the implementa-

tion makes use of the assumption that corresponding words are similar in terms of their word vectors in embedding space in order to find the word correlations in each sentence.

After the translation of the sentences, applying the word alignment estimation on the set of sentence pairs given the refinements for tokenizer and word alignment yields essential information on the relationship between the annotation spans of the English entity labels and their German counterparts. This step is crucial because potentially misaligned labels are further propagated and impede the quality of the dataset and NER scores of the final model. The process is illustrated by Figure 2.

As a minor disadvantage of the common *Pharaoh* alignment format, the difference in annotation granularity cannot be preserved completely on character level. Even though the annotation spans of the source dataset are provided as character-level indices, the word-level tokenization restricts the ability to reconstruct sub-word-level annotation spans in the German target data when the backprojection of the word-level indices from the word alignment estimation onto the character-level indices of the target sentence text string is evaluated.

2.2. Entity Recognition Training

The training of our entity recognition model employs the entity recognition parser from the *SpaCy* library which follows a transducer-based parsing approach[26] with a BILOU[16] scheme (*Begin, Inside, Last, Outside, Unit*; an extension to the IOB[35] scheme) instead of a state-agnostic token tagging approach.

Slim model: Without the use of a transfer-learning-based approach, in *SpaCy* the transformation from discrete tokens into a dense vector representation is implemented by a model that is usually trained from scratch. Such model includes the embedding of the tokens into vectors via Bloom[29] embeddings and further uses convolutional and dense layers to establish context-awareness and feature abstraction.

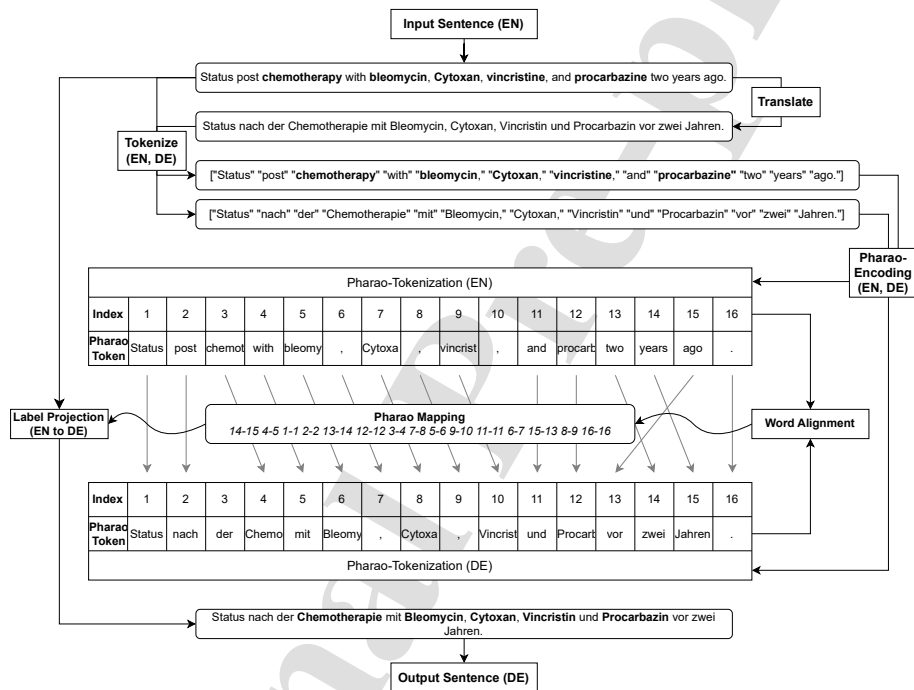


Figure 2: Whitespace-based tokenization and additional Pharaoh-based tokenization for word alignment with subsequent annotation projection. Annotations in the text samples are highlighted by bold font. Only *Drug* annotations are shown in this example.

160 **Transfer-learning:** Inspired by the success of transformer-based neural networks and their effectiveness on language modeling through pre-training on large-scale text corpora, transfer-learning-based methods using deep transformer models can also contribute to stronger entity recognition performance by providing contextualized token embeddings through earlier pre-training without the
165 need to train such large models from scratch. As one instance, the masked language model BERT and several descendants have been released with pre-trained weights for various different languages including German, making it well-suited for transfer-learning.

170 **Entity Parsing:** The entity parser from the SpaCy implementation is strongly influenced by the state-based text chunking algorithm from Lample et al. [26]. The parser uses the feature vectors from previous stages (such as from the slim model or the transfer-learning approach) and aggregates a feature vector from the current parsing state to predict the next valid action which likewise annotates the current token during NER parsing. The whole process is
175 shown in Figure 3.

3. Results

3.1. Dataset Acquisition

The English source dataset from the *2018 n2c2 shared task on ADE and medication extraction in EHR* consists of 404 annotated text documents. The
180 annotation includes the labels *Strength, Form, Dosage, Route, Frequency, Drug, Duration, Reason, ADE*. The documents are split into sentences using the SpaCy sentencizer for English texts. After the sentence-wise translation we apply the word alignment step. During this process we discard sentences whenever an annotation label cannot be reconstructed due to incomplete word alignment
185 mappings. We obtain our raw German dataset with 17938 sentences. The annotation distribution of the raw German dataset is shown in Table 1.

For further clean-up of the raw dataset, sentences that do not contain any entity label at all are discarded from the set of sentences, resulting in a total of

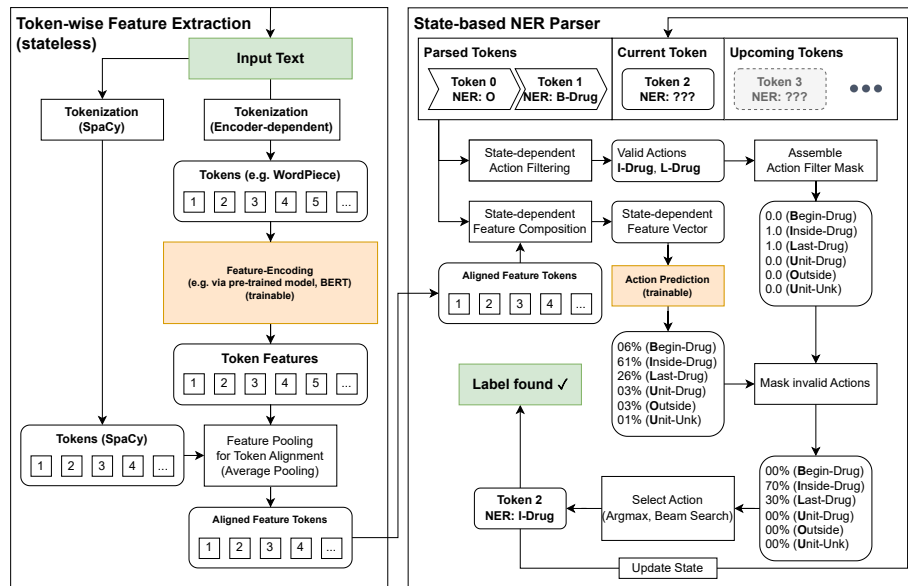


Figure 3: Logical text processing steps for text encoding and entity parsing in SpaCy. The feature encoding can utilize pre-trained deep embeddings via transfer-learning or SpaCy's native Bloom embeddings[29]. Abbreviations: named entity recognition (NER)

NER Tag	Count
Drug	26003
Route	8560
Reason	6244
Strength	10546
Frequency	9794
Duration	956
Form	10546
Dosage	6700
ADE	1557

Table 1: The distribution of annotations in the (raw) synthesized German dataset in absolute numbers. Note that a single tag sample count may include multiple tokens. The dataset consists of 16632 sentences. Abbreviations: named entity recognition (NER)

16632 sentence samples.

190 3.2. Entity Recognition Training

For the training of the NER model, we ignore the following annotation labels for the following reasons:³

- **ADE**: The scope of the English source dataset covers the analysis of medical texts with respect to adverse drug effects. We consider the task of detecting adverse drug effects in texts as of lesser general interest and observed low scores in preliminary experiments when we trained a NER model on all labels including *ADE*. In general, the decision on text phrases in the *ADE* class is complex and context-dependent across datasets.
- **Reason**: Similiar to *ADE*, its usefulness depends on the nature of the dataset and the context, and in preliminary experiments the label class

³Experimental results on NER model training for all label classes as well the visualization of class-specific label text distributions are provided as supplementary data.

yielded low scores.

- **Route:** While we consider *Route* to be of potential general interest, we found that the label diversity in the English source dataset is quite low. For instance, 5356 times (out of 8560 total *Route* annotations) the phrases' value is "PO". The second most frequent value is "IV" (874 times). We decided to refrain from including the *Route* label class because its lack of diversity yields to high scores on the test set and could lead readers to draw misleading conclusions about the actual annotation capabilities of a model for this label class.

Before the entity recognition model is trained, we split the previously described, filtered German dataset into training, validation and test set (80%,10%,10%). The split statistics are provided in Table 2. Since the IOB-based entity recognition parser requires the annotated dataset to contain only non-overlapping annotation spans, annotation overlaps are resolved by removing the annotation span of shorter length while only preserving the longest span.

Dataset	Split	# Tokens	# Entities	# Sentences
Train Set	0.8	293693	50955	13306
Validation Set	0.1	37218	6420	1663
Test Set	0.1	36168	6064	1663
Total	1.0	367079	63439	16632

Table 2: Information on the filtered German dataset. Overlapping annotation spans were removed. The following named entity recognition (NER) tags were omitted: Route, Reason, ADE

We investigate the ability of improving the entity recognition performance by the means of transfer-learning on deep language models on the basis of two German models:

- **German BERT**[6]⁴ (*bert-base-german-cased*): The model from Deepset

⁴<https://www.deepset.ai/german-bert>

220 AI follows the default architecture of BERT and has been specifically pre-trained on German data. The pre-training dataset stems from German Wikipedia, OpenLegalData, and German news articles.

- **GottBERT**[40]: The model is based on the RoBERTa architecture and has been trained on the OSCAR dataset using the fairseq implementation.

225 OSCAR is a German subset of CommonCrawl.

Both language models are publicly available. We retrieve both models from the Huggingface platform. For fine-tuning the entity recognizer on top of the language model, we utilize SpaCy for training. In this context, the model-specific tokenizer is inherited from the language model.

230 The training was performed on a single Nvidia Titan RTX. The training took 8-47 minutes (*German BERT*: 47m, *GottBERT*: 26m, *Slim*: 8m). Due to our observations from the preliminary hyperparameter search, we chose to stick to the default hyperparameters from SpaCy (Adam with weight decay, $\alpha = 0.00005$ (GottBERT, GermanBERT) / 0.001 (Slim), $\beta_1 = 0.9$, $\beta_2 = 0.999$, batch size = 128 (GottBERT, GermanBERT) / 1000 (Slim)) as we did not find major score-wise improvements. In order to measure the differences in performance scores, we also compare the SpaCy Slim model using the same training and test set as baseline model, as well as the publicly available GERNERMED model as static model evaluated on the test set. It should be noted that the GERNERMED model scores must be considered as tainted because its weights are trained on a dataset that might partially contain samples from our test set. For evaluation, the NER procedure is considered as a token-wise multi-class classification problem. We computed the precision (*Pr*), recall (*Re*) and F1 score (*F1*) for each individual label class as well as its respective (class-frequency-weighted) average score (*Total*). The final results on the test set are depicted in Table 3.

240 Both transfer-learning-based approaches exhibit strong performance in absolute numbers. Though to our surprise, German BERT achieves notably inferior performance scores in direct comparison to GottBERT by 0.7% total F1 score

<i>Scores on Test Set</i>		NER Tags						
Model		Str	Dur	Form	Dos	Drug	Freq	Total
GERNERMED++ (GottBERT)	Pr	0.971	0.806	0.947	0.967	0.969	0.880	0.942
	Re	0.964	0.825	0.969	0.971	0.923	0.953	0.950
	F1	0.967	0.815	0.958	0.969	0.945	0.915	0.946
GERNERMED++ (GermanBERT)	Pr	0.944	0.791	0.956	0.963	0.969	0.859	0.932
	Re	0.973	0.825	0.962	0.971	0.933	0.924	0.947
	F1	0.958	0.807	0.959	0.967	0.951	0.890	0.939
GERNERMED++ (SpaCy Slim)	Pr	0.965	0.823	0.965	0.958	0.929	0.855	0.926
	Re	0.967	0.749	0.950	0.971	0.884	0.966	0.941
	F1	0.966	0.784	0.957	0.964	0.906	0.907	0.932
GERNERMED[13] [†]	Pr	0.916	0.613	0.842	0.915	0.644	0.739	0.790
	Re	0.917	0.697	0.882	0.959	0.634	0.901	0.841
	F1	0.917	0.652	0.861	0.937	0.639	0.812	0.814

Note: [†] specific training set might be tainted by samples from the test set.

Table 3: Evaluation of models’ performance scores on test set for the labels **Strength**, **Duration**, **Form**, **Dosage**, **Drug** and **Frequency**. **Precision**, **Recall** and **F1**-scores are evaluated. Abbreviations: named entity recognition (NER)

250 difference. We attribute this performance gap to the differences in pre-training dataset sizes for German BERT (12GB) and GottBERT (145GB) and the use of the RoBERTa architecture as for NER such observation and conclusion have been reported and drawn by the authors of GottBERT as well for monolingual models[40].

255 To verify the robustness of our observations and estimate the degree of a test set selection bias, we re-trained the GottBERT model using 10-fold cross-validation on the dataset. GottBERT was chosen due to its strongest total F1 score in Table 3. The mean and standard deviation of the 10-fold models are provided as well as the distance to the GottBERT results from Table 3 to the mean scores. The results are shown in Table 4.

260

<i>10-fold Cross-validation</i> (GottBERT model)		NER Tags						Total
		Str	Dur	Form	Dos	Drug	Freq	
Precision	μ (mean)	0.967	0.798	0.964	0.962	0.938	0.961	0.950
	σ (std dev)	0.008	0.043	0.012	0.015	0.012	0.009	0.004
	Δ (diff to ref)	-0.004	-0.008	0.017	-0.005	-0.031	0.081	0.008
Recall	μ (mean)	0.967	0.841	0.953	0.958	0.958	0.863	0.939
	σ (std dev)	0.010	0.066	0.010	0.010	0.010	0.014	0.008
	Δ (diff to ref)	0.003	0.016	-0.016	-0.013	0.035	-0.09	-0.011
F1	μ (mean)	0.967	0.817	0.958	0.960	0.948	0.909	0.944
	σ (std dev)	0.006	0.033	0.006	0.010	0.008	0.008	0.004
	Δ (diff to ref)	0.000	0.002	0.000	-0.009	-0.003	0.003	-0.002

Table 4: Averaged scores of test folds from 10-fold cross-validation for labels **Strength**, **Duration**, **Form**, **Dosage**, **Drug** and **Frequency**. All fold-wisely trained models are based on GottBERT. For reference, the score differences to the presented GottBERT model from Table 6 are given. Abbreviations: named entity recognition (NER), standard deviation (std dev), difference to reference (diff to ref)

3.3. Out-of-Distribution Evaluation

The evaluation on the test set does not provide valuable information on how a model can maintain its scores beyond the scope of the train and test set. A known property of neural networks as statistical models is their ability to overfit to the training dataset. While strong performance on the test set indicates the ability to abstract from individual samples without blunt sample memorization, it cannot measure the model’s reliance on the inherent bias of the dataset and its ability to generalize to *out-of-distribution* (OoD) samples. To investigate the OoD generalization ability, we retrieved 30 text samples provided by independent physicians annotated with equivalent labels to our dataset and evaluated the models’ performance on this separated dataset. Since the physicians were instructed to use the class labels from our initial dataset, the OoD samples are annotated with matching label classes and can be directly used for full evaluation of our models. The results are shown in Table 5.

The results display the impact of the transfer-learning-based NER models in order to preserve strong performance on OoD data samples. However similar to

<i>Scores on OoD Dataset</i>		NER Tags						
Model		Str	Dur	Form	Dos	Drug	Freq	Total
GERNERMED++ (GottBERT)	Pr	0.866	1.000	1.000	0.125	0.891	0.923	0.883
	Re	0.960	0.400	0.632	0.250	0.932	0.615	0.835
	F1	0.911	0.571	0.774	0.167	0.911	0.738	0.845
GERNERMED++ (GermanBERT)	Pr	0.955	1.000	0.909	0.077	0.830	0.456	0.817
	Re	0.832	0.800	0.526	0.250	1.000	0.667	0.797
	F1	0.889	0.889	0.667	0.118	0.907	0.542	0.794
GERNERMED++ (SpaCy Slim)	Pr	0.951	0.000	1.000	0.111	0.690	0.486	0.778
	Re	0.772	0.000	0.316	0.250	0.659	0.462	0.623
	F1	0.852	0.000	0.480	0.154	0.674	0.474	0.679
GERNERMED	Pr	0.851	0.000	0.500	0.045	0.460	0.390	0.619
	Re	0.624	0.000	0.158	0.250	0.523	0.410	0.500
	F1	0.720	0.000	0.240	0.077	0.489	0.400	0.541
#Labels		37	3	19	4	36	20	119

Table 5: Evaluation of models’ performance scores on separated out-of-distribution (OoD) dataset for the labels **Strength**, **Duration**, **Form**, **Dosage**, **Drug** and **Frequency**. **Precision**, **Recall** and **F1**-scores are evaluated. Abbreviations: named entity recognition (NER)

the results on the test set, German BERT performs inferior to the GottBERT-based model by an increased margin according to the weighted F1 score. In contrast, the baseline models suffer from substantially degraded scores in comparison to their scores on the test set.

Due to the sparseness and independent origin of the OoD dataset, the number of labels is imbalanced across individual class labels and explains that the evaluation scores can yield 1.0 or 0.0 in several situations. While the reliability of the scores in these cases remains a major limitation, the scores still indicate the degree of abstraction beyond the in-distribution bias in other cases, because the evaluation on the test set is unable to quantify such in-distribution biases.

3.4. Related Datasets

We select three relevant datasets in order to further evaluate our models. To put our results in perspective, we also evaluate the reference model from

²⁹⁰ GGPONC[4] on these datasets. The entity labels from the datasets differ from the labels of our training dataset and our OoD dataset. This limits our ability to perform a complete comparison of our model with respect to all label classes. All related datasets provide annotation information on entities that we consider to be semantically strongly related to the class label *Drug*, although the datasets ²⁹⁵ commonly lack clear and homogeneous definitions on their label classes. We evaluate the scores as a classification task on token- and character-level. The results are shown in Table 6.

Scores on Related Datasets		F1 Scores	
Model / Dataset		Drug (char-wise)	Drug (token-wise)
Medline Dataset[22]		Drug=CHEM	
GERNERMED++ (GottBERT)	Pr	0.858	0.837
	Re	0.701	0.706
	F1	0.772	0.766
GERNERMED++ (GermanBERT)	Pr	0.885	0.875
	Re	0.638	0.686
	F1	0.742	0.769
GERNERMED++ (SpaCy Slim)	Pr	0.437	0.500
	Re	0.182	0.216
	F1	0.257	0.301
GERNERMED	Pr	0.477	0.414
	Re	0.207	0.235
	F1	0.288	0.300
GGPONC[4]	Pr	0.822	0.771
	Re	0.488	0.529
	F1	0.612	0.628
GGPONC Dataset[3]		Drug=Chemicals_Drugs	
GERNERMED++ (GottBERT)	Pr	0.535	n/a
	Re	0.664	n/a
	F1	0.592	n/a
GERNERMED++ (GermanBERT)	Pr	0.522	n/a
	Re	0.645	n/a
	F1	0.577	n/a
GERNERMED++ (SpaCy Slim)	Pr	0.185	n/a
	Re	0.433	n/a
	F1	0.260	n/a
GERNERMED	Pr	0.089	n/a
	Re	0.303	n/a
	F1	0.138	n/a
GGPONC[4]	Pr	0.636	n/a
	Re	0.737	n/a
	F1	0.683	n/a
BRONCO Dataset[21]		Drug=MEDICATION	
GERNERMED++ (GottBERT)	Pr	0.673	0.726
	Re	0.789	0.752
	F1	0.726	0.739
GERNERMED++ (GermanBERT)	Pr	0.684	0.730
	Re	0.677	0.637
	F1	0.680	0.680
GERNERMED++ (SpaCy Slim)	Pr	0.320	0.378
	Re	0.512	0.486
	F1	0.394	0.425
GERNERMED	Pr	0.155	0.148
	Re	0.478	0.482
	F1	0.234	0.227
GGPONC[4]	Pr	0.573	0.346
	Re	0.449	0.430
	F1	0.504	0.384

Table 6: Evaluation of models' F1 scores on related dataset. The GGPONC reference model[4] is evaluated for comparison. To allow fair comparison, only Drug-related label classes are selected. Annotations from the GGPONC[3] dataset do not align onto the tokens from the SpaCy tokenizer and are therefore omitted. **P**recision, **R**ecall and **F1**-scores are evaluated.

To no surprise, the GGPONC reference model archives better performance on its native GGPONC dataset[3], yet all our models with transfer-learning-
300 based, pre-trained BERT encoder outperform the reference model, our slim model and the baseline GERNERMED model. Considering that the baseline GGPONC model was developed in traditional fashion using a manually crafted German dataset, the archived performance margins from both GottBERT- and GermanBERT-based models are unexpected. Throughout the tasks, the GottBERT-
305 based model beats the GermanBERT-based model which is consistent with previous observations.

4. Discussion

Our results indicate strong performance of all models on the test set, however our evaluation on the OoD dataset as well as on external, related datasets shows
310 the impact of using the transfer-learning abilities of pre-trained BERT-based feature encoders to solidify the robust performance on such external datasets. Considering the fact that our models were developed without additional manual work of annotating datasets and only a public non-German dataset was used, the obtained models compete surprisingly well with the pre-existing reference
315 model and are able to outperform it on independent datasets. The lack of more independent annotated datasets, lacking matching annotation labels and unclear label class definitions still limit the possibility to deeper evaluate and compare novel models and methods. In this context, the small sample size of our OoD dataset remains a major limitation of our work and emphasizes the
320 continuous need for German medical corpora with diverse label annotations.

In general, considering the current poor availability of open medical NLP systems for non-English natural languages as well as for German in particular, our refined approach demonstrates a powerful opportunity to build a strong medical NER model solely by the use of a public English dataset.

325 5. Conclusion

In this work, we presented a fine-tuned German NER model for semantic medical entity annotation using deep pre-trained language models by the means of transfer-learning. We demonstrated its ability to outperform the basic baseline model on the test set and on an out-of-distribution dataset. In comparison to the existing GGPONC reference model, we showed competitive results on external datasets and outperformed the reference model on all independent datasets. Furthermore, we described the process and its relevant improvements to obtain a medical-specific German dataset without the use of internal data. Our open NER model is publicly available for third-party use on GitHub.

335 Acknowledgment

This work is a part of the DIFUTURE project funded by the German Ministry of Education and Research (Bundesministerium für Bildung und Forschung, BMBF) grant FKZ01ZZ1804E.

References

- 340 [1] Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78. Association for Computational Linguistics, 2019.
- 345 [2] Iz Beltagy, Kyle Lo, and Arman Cohan. SciBERT: Pretrained language model for scientific text. In *EMNLP*, 2019. eprint: arXiv:1903.10676.
- [3] Florian Borchert, Christina Lohr, Luise Modersohn, Thomas Langer, Markus Follmann, Jan Philipp Sachs, Udo Hahn, and Matthieu-P Schapranow. GGPONC: A corpus of german medical text with rich metadata based on clinical practice guidelines. In *Proceedings of the 11th International*

Workshop on Health Text Mining and Information Analysis, pages 38–48, 2020.

- [4] Florian Borchert, Christina Lohr, Luise Modersohn, Jonas Witt, Thomas Langer, Markus Follmann, Matthias Gietzelt, Bert Arnrich, Udo Hahn, and Matthieu-P. Schapranow. GGPONC 2.0 - the german clinical guideline corpus for oncology: Curation workflow, annotation policy, baseline NER taggers. In *Proceedings of the Language Resources and Evaluation Conference*, pages 3650–3660. European Language Resources Association, 2022.
- [5] Claudia Bretschneider, Sonja Zillner, and Matthias Hammon. Identifying pathological findings in german radiology reports using a syntacto-semantic parsing approach. In *Proceedings of the 2013 Workshop on Biomedical Natural Language Processing*, pages 27–35, 2013.
- [6] Branden Chan, Timo Möller, Malte Pietsch, and Tanay Soni. German BERT - state of the art language model for german NLP, 2019.
- [7] Wendy W. Chapman, Prakash M. Nadkarni, Lynette Hirschman, Leonard W. D’Avolio, Guergana K. Savova, and Ozlem Uzuner. Overcoming barriers to NLP for clinical text: the role of shared tasks and the need for additional creative solutions. *Journal of the American Medical Informatics Association: JAMIA*, 18(5):540–543, 2011.
- [8] Viviana Cotik, Roland Roller, Feiyu Xu, Hans Uszkoreit, Klemens Budde, and Danilo Schmidt. Negation detection in clinical reports written in german. In *Proceedings of the Fifth Workshop on Building and Evaluating Resources for Biomedical Text Mining (BioTextM2016)*, pages 115–124, 2016.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. eprint: 1810.04805.

- [10] Zi-Yi Dou and Graham Neubig. Word alignment by fine-tuning embeddings on parallel corpora. *arXiv:2101.08231 [cs]*, 2021.
- 380 [11] Chris Dyer, Victor Chahuneau, and Noah A. Smith. A simple, fast, and effective reparameterization of IBM model 2. In Lucy Vanderwende, Hal Daumé III, and Katrin Kirchhoff, editors, *Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, June 9-14, 2013, Westin Peachtree Plaza Hotel, Atlanta, Georgia, USA*, pages 644–648. The Association for
385 Computational Linguistics, 2013.
- [12] Georg Fette, Maximilian Ertl, Anja Wörner, Peter Kluegl, Stefan Störk, and Frank Puppe. Information extraction from unstructured electronic health records and integration into a data warehouse. *INFORMATIK 2012*,
390 2012. Publisher: Gesellschaft für Informatik eV.
- [13] Johann Frei and Frank Kramer. GERNERMED: An open german medical NER model. *Software Impacts*, 11:100212, 2022.
- [14] Johann Frei and Frank Kramer. German medical named entity recognition model and data set creation using machine translation and word alignment: Algorithm development and validation. *JMIR Formative Research*,
395 7(1):e39077, 2023. Company: JMIR Formative Research Distributor: JMIR Formative Research Institution: JMIR Formative Research Label: JMIR Formative Research Publisher: JMIR Publications Inc., Toronto, Canada.
- 400 [15] Averbis Gmbh. Averbis health discovery - analyse von patienten daten.
- [16] R. Grishman and Andrew Borthwick. A maximum entropy approach to named entity recognition. 1999.
- [17] Udo Hahn, Franz Matthies, Christina Lohr, and Markus Löffler. 3000pa-
405 towards a national reference corpus of german clinical language. In *MIE*, pages 26–30, 2018.

- [18] Sam Henry, Kevin Buchan, Michele Filannino, Amber Stubbs, and Ozlem Uzuner. 2018 n2c2 shared task on adverse drug events and medication extraction in electronic health records. *Journal of the American Medical Informatics Association : JAMIA*, 27(1):3–12, 2020.
- 410 [19] John Snow Labs Inc. Detect symptoms, treatments and other entities in german- spark NLP model, 2021.
- [20] Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. SimAlign: High quality word alignments without parallel training data using static and contextualized embeddings. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1627–1643. Association for Computational Linguistics, 2020.
- 415 [21] Madeleine Kittner, Mario Lamping, Damian T Rieke, Julian Götze, Bariya Bajwa, Ivan Jelas, Gina Rüter, Hanjo Hautow, Mario Sängler, Maryam Habibi, Marit Zettwitz, Till de Bortoli, Leonie Ostermann, Jurica Ševa, Johannes Starlinger, Oliver Kohlbacher, Nisar P Malek, Ulrich Keilholz, and Ulf Leser. Annotation and initial evaluation of a large annotated german oncological corpus. *JAMIA Open*, 4(2):ooab025, 2021.
- 420 [22] Jan A Kors, Simon Clematide, Saber A Akhondi, Erik M van Mulligen, and Dietrich Rebholz-Schuhmann. A multilingual gold-standard corpus for biomedical concept recognition: the mantra GSC. *Journal of the American Medical Informatics Association*, 22(5):948–956, 2015.
- 425 [23] Jonathan Krebs, Hamo Corovic, Georg Dietrich, Max Ertl, Georg Fette, Mathias Kaspar, Markus Krug, Stefan Störk, and Frank Puppe. Semi-automatic terminology generation for information extraction from german chest x-ray reports. *GMDS*, 243:80–84, 2017.
- 430 [24] Markus Kreuzthaler and Stefan Schulz. Detection of sentence boundaries and abbreviations in clinical narratives. In *BMC medical informatics and decision making*, volume 15, pages 1–13. BioMed Central, 2015.

- [25] Maximilian König, André Sander, Ilja Demuth, Daniel Diekmann, and Elisabeth Steinhagen-Thiessen. Knowledge-based best of breed approach for automated detection of clinical events based on german free text digital hospital discharge letters. *PloS one*, 14(11):e0224916, 2019. Publisher: Public Library of Science San Francisco, CA USA.
- [26] Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. Neural architectures for named entity recognition, 2016.
- [27] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 2019.
- [28] Fei Li, Yonghao Jin, Weisong Liu, Bhanu Pratap Singh Rawat, Pengshan Cai, and Hong Yu. Fine-tuning bidirectional encoder representations from transformers (BERT)-based models on large-scale electronic health record notes: An empirical study. *JMIR Med Inform*, 7(3):e14830, 2019.
- [29] Lester James Miranda, Ákos Kádár, Adriane Boyd, Sofie Van Landeghem, Anders Søggaard, and Matthew Honnibal. Multi hash embeddings in spaCy, 2022.
- [30] Jose A Miñarro-Giménez, Ronald Cornet, Marie-Christine Jaulent, Heike Dewenter, Sylvia Thun, Kirstine Rosenbeck Gøeg, Daniel Karlsson, and Stefan Schulz. Quantitative analysis of manual annotation of clinical text samples. *International journal of medical informatics*, 123:37–48, 2019. Publisher: Elsevier.
- [31] Franz Josef Och and Hermann Ney. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51, 2003.
- [32] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for

sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.

- [33] Yifan Peng, Shankai Yan, and Zhiyong Lu. Transfer learning in biomedical
465 natural language processing: An evaluation of BERT and ELMo on ten
benchmarking datasets. In *Proceedings of the 18th BioNLP Workshop and
Shared Task*, pages 58–65, 2019.
- [34] Tom J Pollard and Alistair EW Johnson. The MIMIC-III clinical database,
2016.
- 470 [35] Lance Ramshaw and Mitch Marcus. Text chunking using transformation-
based learning. In *Third Workshop on Very Large Corpora*, 1995.
- [36] Laila Rasmy, Yang Xiang, Ziqian Xie, Cui Tao, and Degui Zhi. Med-BERT:
pretrained contextualized embeddings on large-scale structured electronic
health records for disease prediction. *npj Digital Medicine*, 4(1):1–13, 2021.
475 Publisher: Nature Publishing Group.
- [37] Roland Roller, Christoph Alt, Laura Seiffe, and He Wang. mEx - an infor-
mation extraction platform for german medical text. In *Proceedings of the
11th International Conference on Semantic Web Applications and Tools for
Healthcare and Life Sciences (SWAT4HCLS'2018). Semantic Web Appli-
cations and Tools for Healthcare and Life Sciences (SWAT4HCLS-2018),
480 December 3-5, Antwerp, Belgium*, 2018.
- [38] Roland Roller, Laura Seiffe, Ammer Ayach, Sebastian Möller, Oliver
Marten, Michael Mikhailov, Christoph Alt, Danilo Schmidt, Fabian Hal-
leck, Marcel Naik, Wiebke Duettmann, and Klemens Budde. A medical
485 information extraction workbench to process german clinical text, 2022.
- [39] Roland Roller, Hans Uszkoreit, Feiyu Xu, Laura Seiffe, Michael Mikhailov,
Oliver Staeck, Klemens Budde, Fabian Halleck, and Danilo Schmidt. A
fine-grained corpus annotation schema of german nephrology records. In

- 490 *Proceedings of the Clinical Natural Language Processing Workshop (ClinicalNLP)*, pages 69–77, 2016.
- [40] Raphael Schreible, Fabian Thomeczyk, P. Tippmann, V. Jaravine, and M. Boeker. GottBERT: a pure german language model. *ArXiv*, 2020.
- [41] Anton Schäfer, Nils Blach, Oliver Rausch, Maximilian Warm, and Nils Krüger. Towards automated anamnesis summarization: BERT-based models for symptom extraction. *arXiv:2011.01696 [cs]*, 2020.
- 495 [42] Johannes Starlinger, Madeleine Kittner, Oliver Blankenstein, and Ulf Leser. How to improve information extraction from german medical records. *it - Information Technology*, 59(4):171–179, 2017. Publisher: De Gruyter Oldenbourg.
- 500 [43] Tian-Xiang Sun, Xiang-Yang Liu, Xi-Peng Qiu, and Xuan-Jing Huang. Paradigm shift in natural language processing. *Machine Intelligence Research*, 19(3):169–183, 2022.
- [44] Martin Toepfer, Hamo Corovic, Georg Fette, Peter Klügl, Stefan Störk, and Frank Puppe. Fine-grained information extraction from german transthoracic echocardiography reports. *BMC medical informatics and decision making*, 15(1):1–16, 2015. Publisher: Springer.
- 505 [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- 510 [46] Joachim Wermter and Udo Hahn. An annotated german-language medical text corpus as language resource. In *LREC*. Citeseer, 2004.
- [47] Honghan Wu, Minhong Wang, Jinge Wu, Farah Francis, Yun-Hsuan Chang, Alex Shavick, Hang Dong, Michael T. C. Poon, Natalie Fitzpatrick, Adam P. Levine, Luke T. Slater, Alex Handy, Andreas Karwath, Georgios V. Gkoutos, Claude Chelala, Anoop Dinesh Shah, Robert Stewart,
- 515

Nigel Collier, Beatrice Alex, William Whiteley, Cathie Sudlow, Angus Roberts, and Richard J. B. Dobson. A survey on clinical natural language processing in the united kingdom from 2007 to 2022. *npj Digital Medicine*, 5(1):1–15, 2022. Number: 1 Publisher: Nature Publishing Group.

- ⁵²⁰ [48] Robert Östling and J. Tiedemann. Efficient word alignment with markov chain monte carlo. *Prague Bull. Math. Linguistics*, 2016.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

GERNERMED++: Semantic Annotation in German Medical NLP through Transfer-Learning, Translation and Word Alignment

Johann Frei ^{a,*}, Ludwig Frei-Stuber^b, Frank Kramer ^a

^a *IT-Infrastructure for Translational Medical Research, University of Augsburg
Alter Postweg 101, 86159 Augsburg, Germany
firstname.lastname@informatik.uni-augsburg.de*

^b *Institute and Outpatient Clinic for Occupational, Social and Environmental Medicine
80336 Munich, Germany
firstname.lastname@med.uni-muenchen.de*

Abstract

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53

We present a statistical model, GERNERMED++, for German medical natural language processing trained for named entity recognition (NER) as an open, publicly available model. We demonstrate the effectiveness of combining multiple techniques in order to achieve strong results in entity recognition performance by the means of transfer-learning on pre-trained deep language models (LM), word-alignment and neural machine translation, outperforming a pre-existing baseline model on several datasets. Due to the sparse situation of open, public medical entity recognition models for German texts, this work offers benefits to the German research community on medical NLP as a baseline model. The work serves as a refined successor to our first GERNERMED model. Similar to our previous work, our trained model is publicly available to other researchers. The sample code and the statistical model is available at:

<https://github.com/frankkramer-lab/GERNERMED-pp>

Keywords: natural language processing, medical NLP, medical named entity recognition, transfer learning, German NLP, artificial intelligence

*Corresponding author

Medical Semantic Annotation in German NLP

Using Transfer-Learning, Translation and Word Alignment

Obstacles and Challenges

Relevant information only encoded in large, unstructured data

Limited availability of **annotated** datasets and pre-trained models for (bio-)medical semantic annotation in non-English/German NLP

Models trained on internal datasets are often excluded from sharing due to potential **data extraction leaks**.

Public, annotated datasets available **only in English**

Methods

Joint combination of **Machine Translation** and **Word Alignment** enable German dataset synthesis from English source data

Training of **robust** NLP model by using **Transfer-Learning** through deep pre-trained language models

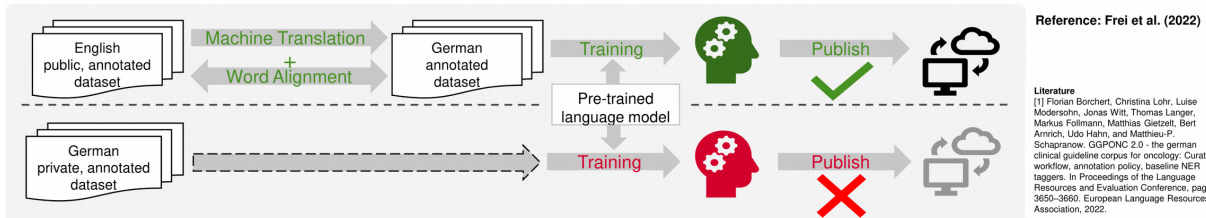
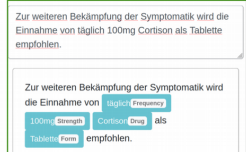
Method relies **only** on publicly available English dataset

Validation on multiple, external datasets

Results & Conclusion

Obtained model **outperforms** prior baseline model (GGPONC 2.0 [1])

Model can be shared **openly, freely and publicly** on GitHub.



Public Significance Statement: Training data for NLP annotation models is a major limiting factor for successful model training. For several reasons, matching datasets are often not available in a certain target language. We combine multiple techniques to utilize data from outside of the target language to obtain an annotation model for our selected target language. Our results show the model's ability to surpass the performance of the baseline model trained traditionally with internal data. Consequently, our work highlights a way to utilize datasets of non-target languages for a certain target language. We apply our method in the context of medical semantic text annotation in German which is a novel contribution to the field.

Johann Frei: Conceptualization, Methodology, Software, Investigation, Validation,
Formal analysis, Writing - Original Draft
Ludwig Frei-Stuber: Resources and Data Curation (OoD dataset author and annotator),
Clinical partner
Frank Kramer: Supervision, Project administration, Funding acquisition, Writing -
Review & Editing

Journal Pre-proof

The authors have declared that no competing interests exist.

Journal Pre-proof