

MaiBaam: A Multi-Dialectal Bavarian Universal Dependency Treebank

Verena Blaschke,[▲] Barbara Kovačić,[▲] Siyao Peng,[▲]
Hinrich Schütze,[▲] Barbara Plank[▲]

- ▲ Center for Information and Language Processing, LMU Munich, Germany
 - Munich Center for Machine Learning (MCML), Munich, Germany
 - Department of Computer Science, IT University of Copenhagen, Denmark
- {verena.blaschke, b.plank}@lmu.de

Abstract

Despite the success of the Universal Dependencies (UD) project exemplified by its impressive language breadth, there is still a lack in ‘within-language breadth’: most treebanks focus on standard languages. Even for German, the language with the most annotations in UD, so far no treebank exists for one of its language varieties spoken by over 10M people: Bavarian. To contribute to closing this gap, we present the first multi-dialect Bavarian treebank (MaiBaam) manually annotated with part-of-speech and syntactic dependency information in UD, covering multiple text genres (wiki, fiction, grammar examples, social, non-fiction). We highlight the morphosyntactic differences between the closely-related Bavarian and German and showcase the rich variability of speakers’ orthographies. Our corpus includes 15k tokens, covering dialects from all Bavarian-speaking areas spanning three countries. We provide baseline parsing and POS tagging results, which are lower than results obtained on German and vary substantially between different graph-based parsers. To support further research on Bavarian syntax, we make our dataset, language-specific guidelines and code publicly available.

Keywords: Less-Resourced Languages, Treebank, Part-of-Speech Tagging, Corpus Creation & Annotation

1. Introduction

In the recent decade, the Universal Dependencies (UD) project (Zeman et al., 2023; de Marneffe et al., 2021) has significantly pushed the frontier in multilingual Natural Language Processing (NLP). UD aims to use consistent syntactic representations for the world’s languages. As of today, UD provides over 240 treebanks in 140+ languages. Despite this coverage, there is still a gap in ‘within-language breadth’ – namely, a lack of diversity within high-resource languages and their closely related non-standard languages and dialects. For example, while Standard German currently has the largest treebank support in UD (with close to 3.8M annotated words as of UD version 2.13), so far UD lacks a treebank for one of the German language varieties spoken by over 10M people¹ in three different countries: Bavarian. In this paper, we present MaiBaam,² the first UD treebank for Bavarian.

Overall, manually annotated corpora are scarce for regional dialects. In large parts, this is due to the fact that collecting and annotating data for non-standard languages and dialects is especially difficult: it is hard to obtain and collect texts, it is challenging to recruit native-speaking annotators with sufficient linguistic background (Miletic

¹The exact speaker population is not known, but Rowley (2011) estimates around 11M.

²EN: ‘maypole’ (lit. ‘May tree’), ‘Maibaum’ in German.

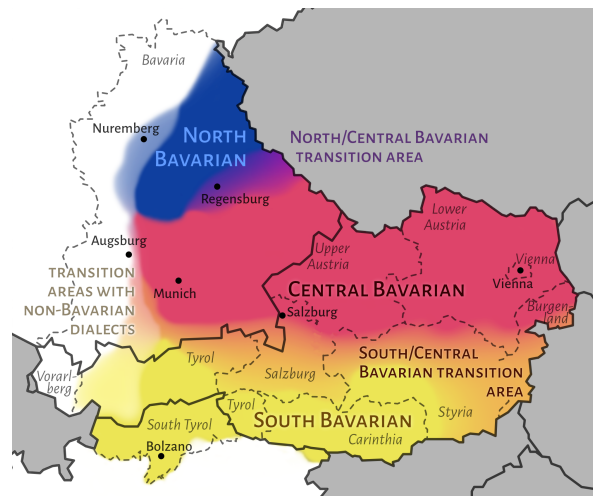


Figure 1: **Bavarian dialect groups in Germany, Austria and Italy**, based on the classification by Wiesinger (1983, map 47.4). Names of dialect groups are in SMALL CAPS, names of provinces and states in *italics*.

et al., 2020), and it requires more expert knowledge and time to adopt guidelines developed for standard languages. Despite all of these challenges, collecting and annotating non-standard and dialectal data is rewarding for at least two reasons. From a linguistic perspective, it allows contrastive analyses of language varieties and studying of morphosyntactic differences to standard lan-

- Trotzdean das'e's moch, hairon tou'e's niat.
- Obwoi i's mog, heirodn dua e's ned.
- Trotz dass i's mog, hairatn tua i's net.

DEU Obwohl ich sie mag, heiraten tu ich sie nicht.

Table 1: **Dialectal and orthographic variation in Bavarian.** A ● North Bavarian grammar example from our corpus, ‘Although I like her I won’t marry her,’ translated into ● Central and ● South Bavarian as well as German (DEU).

guage. From an NLP perspective, dialects provide a unique test bed with a dire need for technological innovations in light of data scarcity (Blaschke et al., 2023b).

Bavarian is a Germanic language variety closely related to German and spoken in parts of Germany, Austria, and Italy. The language status of Bavarian is complicated, as it could be called a language distinct from German on linguistic grounds, but is perceived as a dialect by its speakers (Rowley, 2011). It nevertheless has an ISO 639-3 code: BAR. Bavarian comprises a number of local varieties belonging to three major dialect groups: ● North, ● Central, and ● South Bavarian, connected by transition areas to the ● North and ● South of the Central Bavarian area (Wiesinger, 1983, p. 839). None of these varieties are standardized. Figure 1 shows where they are spoken and Table 1 provides an example of the linguistic and orthographic variation.

Our contributions are as follows:

- We introduce MaiBaam, the first Bavarian UD treebank (§3.1).
- We present annotation guidelines for Bavarian morphosyntactic structures that differ from German ones (§3.3–3.4).
- We analyse transfer performance of multiple parsers trained on German data (§4).

We make our data available at github.com/UniversalDependencies/UD_Bavarian-MaiBaam, to be released in UD v2.14 (May 2024). Additionally, we share our guidelines (Blaschke et al., 2024) and the code we use for preprocessing and annotating the data and for the transfer experiments.³

2. Related Work

The UD project hosts four German treebanks: HDT (Borges Völker et al., 2019), GSD (McDonald et al., 2013), PUD (Zeman et al., 2017), and

³<https://github.com/mainlp/maibaam-code>

Genre	Tokens	Sents	Tok/Sent	Dialect
W	7 988	417	19.2	●●●●●●?
✎	2 485	285	8.7	●●●●●●?
📖	2 019	238	8.5	●●●●●●?
📺	1 599	87	18.4	●●●●●●?
📖	932	43	21.7	●●●●●●?
All	15 023	1 070	14.0	●●●●●●?

Table 2: **Genre distribution in MaiBaam.** Genres: W Wiki, ✎ grammar examples, 📖 non-fiction, 📺 social, 📖 fiction. Dialects: ● North Bavarian, ● North/Central B. transition dialect, ● Central B., ● South/Central B. transition dialect, ● South B., ? un(der)specified dialect.

LIT (Salomoni, 2017). There also exist two dependency treebanks for related non-standard varieties: the Swiss German UZH (Aepli and Clematide, 2018) and the Low Saxon LSDC (Siewert et al., 2021), as well as a non-UD Swiss German corpus with phrase structure annotations (Schönenberger and Haerberli, 2019).

A few NLP datasets include Bavarian data. Both Kontatto (Dal Negro and Ciccolone, 2020) and DiDi (Frey et al., 2015) contain part-of-speech (POS) tags for Bavarian data from South Tyrol; the former was tagged manually or semi-automatically, the latter automatically based on German glosses. The BarNER dataset provides named entity annotations for Bavarian wiki and social media data (Peng et al., 2024). Van der Goot et al. (2021a) and Winkler et al. (2024) have collected South and Central Bavarian slot and intent detection data. The Kontatti corpus (Ghilardi, 2019) and the multi-dialectal Zwirner corpus (IDS, n. d.) contain Bavarian speech data.

3. The MaiBaam Treebank

3.1. Data and Corpus Statistics


Our data come from several different sources that all allow public re-sharing. They span several UD genres (Table 2):


W *Wiki* sentences are taken from Bavarian Wikipedia articles.⁴ We select articles on various different topics to avoid over-representing the template structures that many location-related articles tend to follow. Additionally, we nearly exclusively choose articles tagged as being written in a specific dialect.⁵ As such, while there still might be some overlapping au-


⁴<https://bar.wikipedia.org>; CC BY-SA 4.0


⁵https://bar.wikipedia.org/wiki/Kategorie:Artikel_nach_Dialekt

thors/editors, we expect that the overall set of wiki writers in MaiBaam is fairly diverse.


 *Grammar examples* come from three sources: i) Tatoeba⁶ sentences contributed by users who self-report as Bavarian native speakers, ii) Wikipedia articles that contain collections of linguistic samples, and iii) UD’s Cairo CICLing Corpus,⁷ translated by a Bavarian native speaker.

 *Fiction*: We include parts of non-encyclopedic Wikipedia pages recounting fairy tales.

 The *non-fiction* genre includes questions and commands for a hypothetical digital assistant from the South Tyrolean validation split of xSID (van der Goot et al., 2021a), as well as from the natural (untranslated) queries and the Central Bavarian test split of xSID from Winkler et al. (2024).

 *Social*: We annotate sections of Wikipedia discussion pages and replace usernames mentioned in the text with USERNAME.

Inspired by the analysis of Müller-Eberstein et al. (2021), we annotate each sentence with genre metadata so that relevant patterns can also be analyzed on a genre level rather than only a treebank level. Table 2 shows the distribution of genres in our dataset. Currently, wiki articles represent the largest group, making up 53% of the tokens and 39% of the sentences in MaiBaam. The average sentence length differs across genres, with fiction, wiki articles and discussions having much longer sentences than grammar examples or queries for digital assistants. This is consistent with statistics for other treebanks (Peng, 2023, p. 63).

All of our data sources have metadata indicating that the text is in Bavarian. In many cases, the metadata also mention a more specific dialect or location. Table 3 presents the token-level geographical distribution of MaiBaam across the dialect groups displayed in Figure 1. Just under half of the tokens are in sentences that we can clearly assign to one of the dialect areas. The  Central Bavarian group is the dialect area with the most tokens in MaiBaam (22%). This group also contains the two best-represented sub-regions: the cities of Vienna and Munich.

A significant part of our data does not contain any location or dialect information, or refers to larger regions in which multiple dialects are spoken. In our treebank, we tag each sentence with the most specific dialect and location information available.

⁶https://tatoeba.org/en/sentences/show_all_in/bar/none; CC BY 2.0 FR

⁷github.com/UniversalDependencies/cairo



























Dialect group with location	Tokens	Sents
 North Bavarian	833	65
Western North Bavarian area	308	34
Unspecified North Bavarian	525	31
 North/Central Bavarian	793	47
Bavarian Forest	793	47
 Central Bavarian	3 303	221
Munich	1 166	60
Cent. Bav. in Upper Bavaria	613	76
Salzburg (city)	102	5
Upper Austria	43	5
Vienna	1 356	73
Unspecified Central Bavarian	23	2
 South/Central Bavarian	1 130	50
Bad Reichenhall	206	11
Berchtesgaden	110	5
Pongau	515	21
Pinzgau	299	13
 South Bavarian	995	70
Carinthia	99	4
South Tyrol	896	66
 Underspecified	7 969	617
   ? Upper Bavaria	438	28
   ? Austria	1 491	173
   ? Other C., S./C. or S.	901	122
  ? South East Upper B.	182	10
  ? East Austria	1 687	86
  ? Styria	292	15
     ? Unspecified	2 980	183

Table 3: **Dialect groups in MaiBaam.** ‘Underspecified’ refers to cases where we do not have any dialect or location information (‘unspecified’) or where the specified geographic area encompasses multiple dialect groups.

We include a full data statement (Bender and Friedman, 2018) in Appendix A. Appendix B provides an overview of the POS tag and dependency label distributions in our data.

3.2. Annotation Procedure

The annotation procedure includes training and adjudication sessions, first on a sample of German texts from existing UD treebanks (§2) and then on the target Bavarian data. We train the annotator initially on universal part-of-speech (UPOS) tags and later also on dependencies. We use a modified version of ConlluEditor (Heinecke, 2019) to annotate the data.

The annotator is a computational linguistics student who is a native speaker of German and a (non-Bavarian) Upper German dialect. We also consult three Bavarian native speakers from the

South and Central Bavarian dialect areas for grammaticality judgments, lexical disambiguation and translations. The annotator and consultants involved in this project are hired and compensated according to local standards.

Manually correcting automatically predicted labels is a common strategy in UD annotation to save time and labour (Salomoni, 2017; Borges Völker et al., 2019). To minimize bias from model outputs, we use a simple rule-based pre-tokenizer. Inspired by Berzak et al. (2016), we pre-annotate POS tags where the UDPipe model trained on GSD and the one trained on HDT agree and jointly achieve a precision >95% on an initial test set of 4k tokens.⁸ We do not pre-annotate dependency arcs or labels. In terms of dependency labels, the ones that can easily be predicted are also the ones that can trivially be expressed in a rule-based way as suggestions within ConlluEditor (e.g., if a DET is marked as a dependent of a NOUN, the relation will be *det*).

We largely follow UD’s annotation guidelines for German.⁹ To the extent it is possible, we also follow the annotation decisions made in the closely related treebanks (§2), which however often disagree in more particular grammatical contexts.¹⁰ When the grammatical structures are similar to English ones, we also consult the English EWT (Silveira et al., 2014) and GUM (Zeldes, 2017) treebanks. We use Grew-match (Guillaume, 2021) for querying these treebanks.

We discuss and resolve difficult annotation cases in weekly meetings. As additional approaches to finding annotation errors, we use Udapi (Popel et al., 2017) and UD’s validation scripts.¹¹ Furthermore, we manually double-check the word forms that we annotated with closed-class POS tags.

The annotation time – excluding training time, research into the annotation decisions made for other treebanks, discussions of grammatical phenomena, adjudication meetings and subsequent corrections – totals 165 hours. The average annotation time per sentence varies greatly, depending on the text genre and dialect, as well as on the

⁸The POS tags with high precision scores are: AUX, CCONJ, DET, NOUN, NUM, PART, PRON, and PUNCT. About 41% of the tokens in the documents we pre-annotate receive a POS tag.

⁹<https://universaldependencies.org/de>

¹⁰This observation is not new; see Hovy et al. (2014) and Wisniewski and Yvon (2019) for investigations of POS inconsistencies (the latter specifically within UD) and Zeldes and Schneider (2023) for a comparison of decisions made in two large English UD treebanks. Our annotation guidelines contain more details on the German case.

¹¹<https://github.com/UniversalDependencies/tools>

level of familiarity with the guidelines and annotation tools.

We decide against normalizing the data to an artificial Bavarian standard since no actual written, or even spoken, standard exists. Such a decision would ultimately have been biased towards certain Bavarian dialects, thus conflicting with our goal of curating a diverse set of Bavarian varieties.

3.3. Tokenization

In Bavarian, prepositions and determiners are often contracted, e.g., *beim* ‘at the.DAT’. We follow the UD guidelines for German (see also Grünewald and Friedrich, 2020) and treat such cases as multi-word tokens: *beim* becomes *bei*_{ADP} plus *m*_{DET}. This decision is also consistent with how the Low Saxon guidelines (Siewert et al., 2022)¹² handle tokenization but differs from the decision made for Swiss German to leave merged word sequences as they are written (Aeppli and Clematide, 2018).¹³ Since there is variation in the way determiners are pronounced and written, we simply split the words into substrings (rather than normalizing them to an arbitrary standard).¹⁴

Other commonly fused sequences that we split are determiners followed by common or proper nouns (*d’*_{DET} *Rundn*_{NOUN} ‘the round’) and verbs or complementizers followed by pronouns or neuter determiners (*houd*_{AUX} *s*_{DET} ‘has the’; *habn*_{AUX} *se*_{PRON} *s*_{PRON} ‘they have [...] it’).

When a vowel-initial word is appended to a vowel-final word, a linking consonant can be inserted in between (Merkle, 1993, pp. 30–33). In this case, we include the consonant with the first word (e.g., we analyze *wiera* ‘how he’ with its linking *-r-* as *wier*_{SCONJ} and *a*_{PRON}).

In order to enable comparisons with datasets tokenized like the Swiss German UD treebank, we will include a script upon data release that reverts the token splits, assigns tags to the unsplit tokens (e.g., DET+NOUN becomes NOUN and VERB+PRON becomes VERB) and adjusts the dependencies accordingly.

When it comes to hyphenated compound words, we follow the German HDT treebank and do *not* split them apart: e.g., *Fabel-Viech* ‘mythical creature’ is a single word.

¹²<https://universaldependencies.org/nds>

¹³<https://universaldependencies.org/gsw>

¹⁴Full forms of the dative definite determiner in our corpus include *dem*, *am*, *im* and, due to partial case syncretism of dative and accusative forms (Merkle, 1993, pp. 85, 98), also *den*, *an*, *in*.

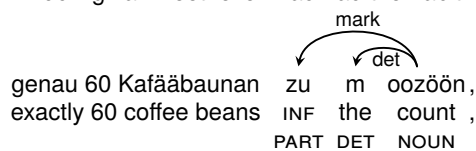
3.4. POS Tags and Dependencies

We use the same set of dependency relations and subrelations¹⁵ as defined for German and refer to the German guidelines and treebanks where possible.¹⁶ However, Bavarian permits syntactic structures that are not licensed in Standard German.¹⁷ We here discuss a range of such structures in our data, along with our annotation decisions.

3.4.1. Verbs

Infinitives In German, many infinitive constructions require the marker *zu*_{PART}. In Bavarian, two similar constructions appear: one where a cliticized form of the marker (*z*) is followed by a verbal infinitive, and one where the infinitive is nominalized and the a cliticized dative determiner (*m/n*) is added to the marker: *zum* or *zun* (Bayer, 1993; Bayer and Brandner, 2004). In both cases, we annotate *z(u)*_{PART} with *mark* (as in the German treebanks), and in the latter, we separately annotate *m/n*_{DET} with *det*:

- (1) Ludwig van Beethoven hod de Gwohnheit ghobt,
Ludwig van Beethoven had had the habit



- um si draus a Schalal Mokka z mochn.
so as to REFL out of it a cup of coffee INF make .

‘Ludwig van Beethoven had a habit of counting exactly 60 coffee beans in order to brew a cup of coffee from them.’ (Wiki *Kafää* ‘Coffee’)

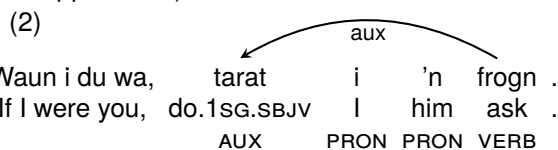
Auxiliary *tua* ‘do’ In addition to the auxiliary verbs named in the German guidelines, we include *tua* ‘do’, which is used in several periphrastic constructions in conjunction with a lexical verb, both

¹⁵These subrelations include dative objects (*obl:arg*), possessive pronouns (*det:poss*), lexicalized reflexive pronouns (*expl:pv*), particle verbs (*compound:prt*), passive constructions (*nsubj:pass*, *csubj:pass*, *aux:pass*, *obl:agent*) and relative clauses (*acl:relcl*, *advcl:relcl*).

¹⁶In cases not mentioned by the guidelines and where the German treebanks disagree, we make decisions based on other sources. For instance, in the case of *selbst/selber* ‘self’ being added to a sentence for emphasis, we follow the analysis by Hole (2002) who distinguishes between *selbst* being used as an adnominal or adverbial intensifier, and attach the word to the corresponding noun or clausal head respectively. Other such cases are detailed in Blaschke et al. (2024).

¹⁷Analytic possession and articles before person names do appear in colloquial German, but are uncommon in written Standard German.

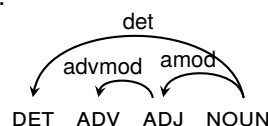
in indicative and subjunctive constructions (Merkle, 1993, pp. 65–67).



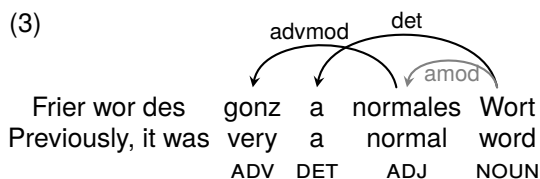
‘If I were you, I would ask him.’ (Tatoeba 5166978)

3.4.2. Noun Phrases

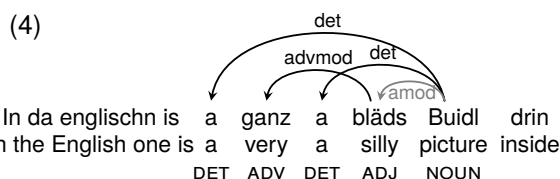
Order of determiner and adverb In German, if an adverb modifies an adjective in a noun phrase, the adverb appears between the determiner and the adjective:



For a small set of Bavarian intensifiers, alternative orders are possible (typically when the determiner is indefinite): the order of adverb and determiner can be reversed (ADV DET ADJ NOUN) and the determiner can be doubled (DET ADV DET ADJ NOUN) (Lenz et al., 2014; Merkle, 1993, pp. 89–90, 158). In such cases, we allow non-projective dependencies:



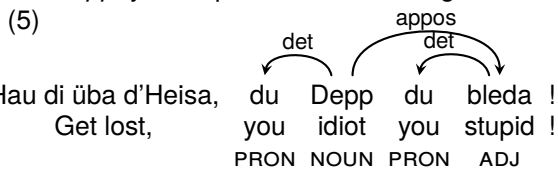
‘It used to be a completely normal word.’ (Wiki *Walsch* ‘Italian/Romance’)



‘The English [wiki] contains a very silly picture [...]’ (Wiki discussion *Ottoman* ‘sofa’)

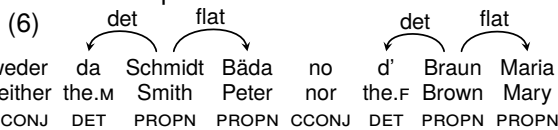
Postponed adjectives For emphasis (and especially when voicing annoyance), phrases of the pattern (ADP) DET ADJ NOUN can be rearranged into (ADP) DET NOUN (ADP) DET ADJ (Merkle, 1993, p. 168). We consider the postponed adjective to be an apposition of the noun. This structure is often combined with constructions where a first or second-person pronoun is used in lieu of a determiner. In such cases, we tag the pronoun as PRON and, following the recommendation by Höhn (2021), label the relation *det*. In the following

sentence in our corpus (pardon our Bavarian), *du bleida Depp* ‘you stupid idiot’ is re-arranged:



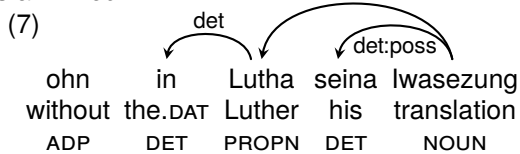
‘Get lost [lit.: scam over the houses], you stupid idiot!’
(Tatoeba 5657152)

Personal names In Bavarian, personal names are preceded by a determiner matching in case and gender (Weiß, 1998, pp. 69–70), and the family name is often put before the given name (Weiß, 1998, p. 71). Following the general UD guidelines, we connect the parts of the name via a *flat* relation:



‘neither Peter Smith nor Mary Brown [...]’
(Cairo CILing 12)

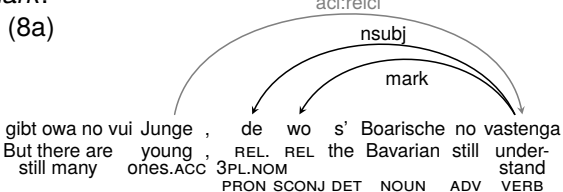
Possession Bavarian, like many German dialects and colloquial variants, eschews the genitive in favour of analytic possessive constructions (Fleischer, 2019; Bülow et al., 2021). These can be prepositional phrases or the prenominal dative construction, in which we analyze the possessor as an *nmod*:



[...] without Luther’s translation [...]
(Wiki discussion *Ödenburg* ‘Sopron’)

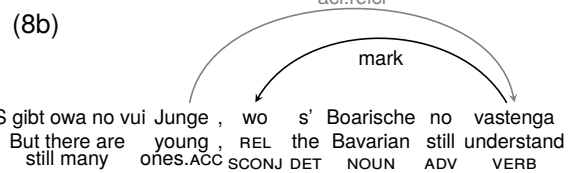
3.4.3. Subordinate Clauses

Relative markers Where German uses the relative pronouns *der/die/das* ‘that, which’, Bavarian can append the invariant relative marker *wo* (in some dialects *was*) (Moser, 2023). We tag the relative pronoun as PRON (as in the German treebanks) and the relative marker as CONJ with the relation *mark*:

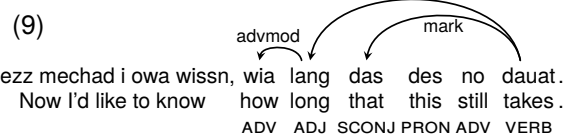


‘However, there are still many young people who still understand Bavarian [...]’
(Wiki *Minga* ‘Munich’)

In certain situations, the relative pronoun can be dropped in Bavarian if the relative marker is present (Pittner, 1996). This can for instance happen when the relative pronoun would be in the nominative case:

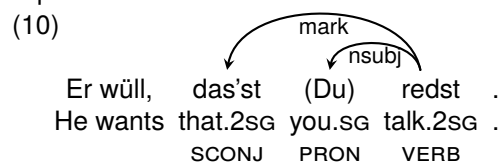


Additional complementizer The adverb, relative pronoun, or question word introducing a subordinate clause can be followed by an additional conjunction *dass* ‘that’ (Weiß, 1998, pp. 29–30; Merkle, 1993, pp. 190–191), which we consider a *marker*:

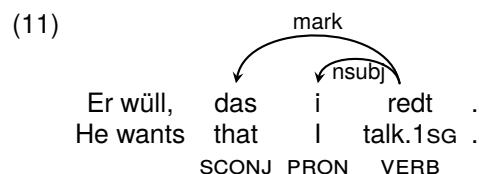


‘Now I’d like to know how long this will still take.’
(Wiki *Pronomen* ‘Pronouns’)

Complementizer agreement In Bavarian, reduced forms of 2nd person (and, optionally, 1PL) pronouns are used when they appear in the Wackernagel position immediately after complementizers. These reduced forms are immediately attached to the previous word and can still be followed by a full pronoun for additional stress (Weiß, 1998, p. 119). Whether these constructions should be analyzed as a word followed by an enclitic pronoun or as inflected complementizers is debatable (for an overview of the different arguments, see Weiß, 1998, pp. 123–133). For our annotations, we follow Bayer (2013) and adopt the interpretation of inflection:



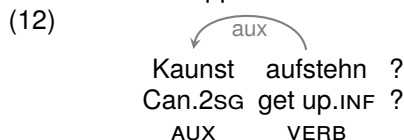
‘He wants you to talk.’
(Wiki *Konjunktiona* ‘Conjunctions’)



‘He wants me to talk.’
(Wiki *Konjunktiona* ‘Conjunctions’)

3.4.4. Other

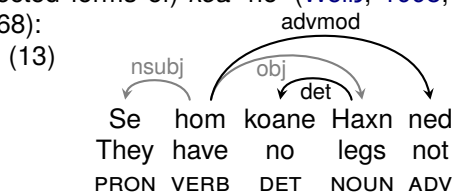
Dropped 2nd person pronouns Similarly, second person pronouns can be omitted when they occur after a correspondingly inflected verb. Consider for instance the sentence *Kaunst du aufstehn?* ‘Can you get up?’ where the pronoun *du* ‘you.sg’ can be dropped:



‘Can you get up?’ (Tatoeba 10673747c)

1PL inflection endings The 1PL.PRES inflection of verbs is typically straightforward, e.g., *mia schbui+n* ‘we play+1PL’. However, it is also possible to add *-ma* to the stem: *mia schbui+ma* (Merkle, 1993, p. 127). Although this ending historically comes from a cliticized form of the pronoun (and some analyze it as such; Weiß, 1998, p. 123, fn. 48), we simply analyze it as inflection: *mia*_{PRON} *schbuima*_{VERB}. This decision also lends itself well to UD annotation: if we were to treat *-ma* as its own word, it is unclear what dependency label it should get, since the independent pronoun *mia* is already the subject.

Negative concord Unlike German, Bavarian allows for negative concord in constructions with (inflected forms of) *koa* ‘no’ (Weiß, 1998, pp. 167–168):



‘They have no legs [...]’ (Wiki *Fiisch* ‘Fish’)

4. Transfer Experiments

This section establishes baselines to gauge how well dependency parsers trained on German perform on our Bavarian treebank.

4.1. Data

Only two German UD treebanks (v2.12) include a training partition: HDT, which comprises almost 3.5M words from web news articles, and GSD, which contains 292k words from news articles, reviews, Wikipedia articles and other webpages. Both are released under a CC BY-SA license. We train and validate on HDT and GSD and test on the entire gold-tokenized MaiBaam.¹⁸

¹⁸Since MaiBaam contains various multi-word tokens (§3.3), this makes for an easier test condition than if we

System or language model	Pretraining language(s)	License
UDPipe 2.12-230717 (Straka, 2018)	see <i>mBERT</i>	CC BY-NC-SA
Stanza 1.6.1 (Qi et al., 2020)	DEU	Apache 2.0
MaChAmp (van der Goot et al., 2021b)	with the language models below:	
<i>mBERT</i> (Devlin et al., 2019)	multi, incl. DEU & BAR	Apache 2.0
<i>bert-base-multilingual-cased</i>		
XLM-R (Conneau et al., 2020)	multi, incl. DEU	MIT
<i>xlm-roberta-base</i>		
GBERT (Chan et al., 2020)	DEU	MIT
<i>deepset/gbert-base</i>		

Table 4: **Systems and pretrained language models used for parsing experiments.** Key: DEU is German, BAR is Bavarian. All systems are finetuned on German data.

4.2. Models

We compare the POS tagging and dependency parsing results of several models. Version and license details can be found in Table 4.

This includes two already trained systems, each trained once on GSD and once on HDT. Firstly, we use UDPipe (Straka, 2018), which combines *mBERT* embeddings (Devlin et al., 2019)¹⁹ with custom word and character embeddings. The architecture of UDPipe is specifically built for dependency parsing and contains steps for ensuring a proper graph structure of the predicted dependencies. Secondly, we investigate the predictions made by Stanza (Qi et al., 2020), which uses German word embeddings (via Zeman et al., 2017) and includes a graph-based dependency parser.

For comparison, we use MaChAmp (van der Goot et al., 2021b) to train models from scratch. We finetune the multilingual models *mBERT* and XLM-R (Conneau et al., 2020) as well as the German model GBERT (Chan et al., 2020), and otherwise use MaChAmp’s default settings. Both multilingual models contain German data in their pre-training data, and *mBERT* was also pretrained on Bavarian Wikipedia data.²⁰

Since the UDPipe and Stanza pipelines trained on HDT perform markedly worse than their GSD counterparts, we only use GSD for finetuning our models. We use the regular GSD treebank, as

tested systems that include a tokenization step. Using the UDPipe models on plain text input, the tokenization F₁ scores are 96.76% for the version trained on GSD, and 95.14% for the HDT version.

¹⁹https://ufal.mff.cuni.cz/udpipe/2/models#universal_dependencies_212_models

²⁰<https://github.com/google-research/bert/blob/master/multilingual.md>

Model	Train	POS		Dependency	
		Acc	F ₁	LAS	UAS
UDPipe	GSD	80.29	62.45	65.79	79.60
UDPipe	HDT	76.36	59.30	61.55	73.59
Stanza	GSD	42.30	36.73	24.89	40.39
Stanza	HDT	39.80	36.10	20.67	29.04
mBERT	GSD _{n40}	78.74	58.74	54.96	66.38
mBERT	GSD	77.47	57.19	52.48	64.30
GBERT	GSD _{n50}	74.68	57.15	50.57	62.67
GBERT	GSD	58.86	47.21	36.40	50.51
XLM-R	GSD _{n40}	72.45	55.33	48.81	61.25
XLM-R	GSD	55.00	44.03	31.42	43.42

Table 5: **Prediction scores on MaiBaam, in %.** The F₁ scores are macro-averaged. Except for the UDPipe and Stanza results, all values are averaged over three runs. The subscript additions _{n40} and _{n50} refer to noise levels of 40 and 50%, respectively (see §4.2 for details).

well as a version with character-level noise inspired by Aeppli and Sennrich (2022) to see if it improves over vanilla fine-tuning. For character-level noise, we select a certain ratio of the words ('noise level') in each sentence and randomly inject noise into each word by replacing, deleting or inserting a character. We use the *split word ratio difference* heuristic (Blaschke et al., 2023a) for selecting appropriate noise levels: For each noise level in {0, 10, 20, ..., 100} and each pretrained language model's tokenizer, we compare the proportion of words that the tokenizer splits into multiple subword tokens in the (noised) German and (untouched) Bavarian data, and select the noise level that minimizes the difference between the two ratios.²¹ For mBERT and XLM-R, this means we inject noise into 40% of the words in a sentence, and 50% for GBERT. We train each model on three random seeds and report the mean results in the next section.

4.3. Results

Table 5 shows the different models' dependency parsing and POS tagging scores. For dependency parsing, we use the unlabelled and labeled at-

²¹We determine the noise levels based on comparisons with just 10% of the sentences in MaiBaam, in order not to overfit to our test data. However, we observe that the split word ratios are very similar if we use the full dataset.

tachment scores (UAS and LAS, respectively),²² and for POS tagging we consider both accuracy and macro-averaged F₁ scores. Figure 2 shows a parse produced by the best system. While most of its predictions are correct, they contain several wrong POS tags and dependencies (arcs and labels). Half of the prediction errors in this sentence affect the phrase with the doubled determiner, which is not licensed by Standard German grammar (§3.4.2; *a rechl a sauwas Wossa* 'lit.: a fairly a clean water').

For both UDPipe and Stanza, we observe that the versions trained on the GSD treebank outperform those trained on the larger HDT. The UDPipe models (regardless of training data) attain the highest scores for all four metrics (except for the HDT version's POS tagging accuracy). The single best model is the UDPipe version trained on GSD, reaching LAS and UAS scores of 65.79% and 79.60%, and, for POS tagging, an accuracy 80.29% and F₁ score of 62.45%. Conversely, the Stanza models do not generalize well to the Bavarian data and achieve the worst scores of all models across all metrics.²³

Focusing on the models we trained on the vanilla GSD data, we observe that mBERT outperforms the other two language models. Presumably, it benefits from the overlap between some of our corpus data and its (unlabelled) pretraining data. Including many languages *other* than Bavarian in the pretraining data does not appear to be nearly as advantageous, as the German model GBERT produces better results than the multilingual XLM-R.

Injecting noise into the training data consistently improves our models' performance. The improvements are especially large for the worse-performing XLM-R and GBERT. For XLM-R, the scores are between 26 and 55% (11.31–17.83 percentage points; pp) higher with noise than without, and GBERT sees improvements of 21 to 39% (9.94–15.82 pp). Conversely, mBERT's scores improve only by between 2 and 5% (1.27–2.48 pp).

4.4. Discussion

The POS tagging scores of our best re-trained model (mBERT with noised data) are competitive with those of the UDPipe models. We hypothesize that UDPipe's character embeddings are of advantage when processing the orthographically very variable Bavarian input. Stanza, on the other hand, uses static German embeddings for entire words –

²²For LAS, we ignore relation subtypes, as in UD's official evaluation script: github.com/UniversalDependencies/tools/blob/master/eval.py.

²³When evaluating the GSD and HDT Stanza models on the test splits of their respective training sets, all scores are very high and similar to those of UDPipe.

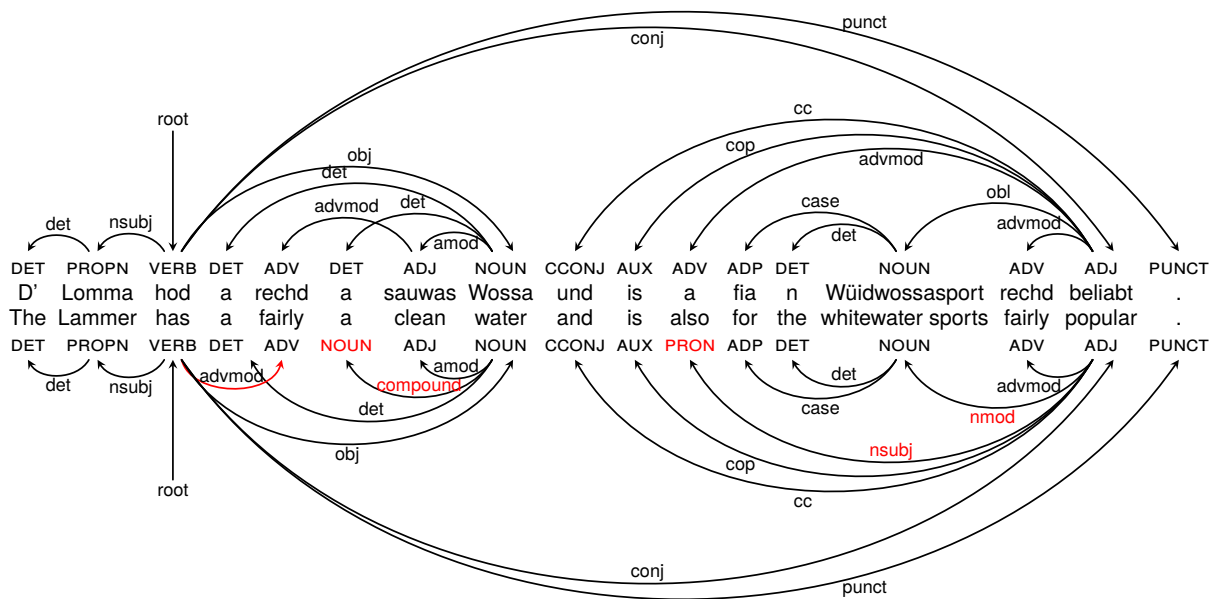


Figure 2: **Gold-standard (top) and predicted (bottom) annotations.** Predictions are produced by the UDPipe model trained on GSD, the best system in our evaluation. Wrong predictions are in red. ‘The Lammer (river) has fairly clean water and is also pretty popular for whitewater sports.’ (Wiki *Lamma* ‘Lammer’)

a mismatch for the Bavarian input. The subword tokens used by mBERT, GBERT and XLM-R provide an intermediary input representation, and the versions finetuned on the noised data are better geared towards processing short subwords. However, the dependency parsing results of our models clearly lag behind UDPipe. We assume that this is owed to UDPipe’s processing steps for properly constructing directed spanning dependency graphs.

5. Conclusion

We present MaiBaam, the first Bavarian treebank, which we manually annotated with POS tags and syntactic dependencies in UD. It comprises a range of dialects and genres. We share MaiBaam as a resource for analyzing and modelling Bavarian data and, more broadly, non-standard language data.

We also conduct transfer learning experiments with models trained on German data to provide parsing and POS tagging baselines for MaiBaam. Even the best model has ample room for improvement. This shows that processing Bavarian data is not as simple as merely using zero-shot transfer from German.

Acknowledgements

We thank Miriam Winkler and Marie Kolm for lending us their native speaker expertise. We also

thank the anonymous reviewers for their feedback.

This research is supported by the ERC Consolidator Grant DIALECT 101043235. We also gratefully acknowledge partial funding by the European Research Council (ERC #740516).

6. Bibliographical References

- Noëmi Aepli and Simon Clematide. 2018. [Parsing approaches for Swiss German](#). In *Proceedings of the 3rd Swiss Text Analytics Conference (SwissText)*, Winterthur, Switzerland.
- Noëmi Aepli and Rico Sennrich. 2022. [Improving zero-shot cross-lingual transfer between closely related languages by injecting character-level noise](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 4074–4083, Dublin, Ireland. Association for Computational Linguistics.
- Josef Bayer. 1993. [Zum in Bavarian and scrambling](#). In Werner Abraham and Josef Bayer, editors, *Dialektsyntax*. Westdeutscher Verlag.
- Josef Bayer. 2013. [Klitisierung, Reanalyse und die Lizenzierung von Nullformen: zwei Beispiele aus dem Bairischen](#). In Werner Abraham and Elisabeth Leiss, editors, *Dialektologie in neuem Gewand. Zu Mikro-/Varietätenlinguistik, Sprachenvergleich und Universalgrammatik*,

- volume 19 of *Linguistische Berichte, Sonderhefte*. Buske.
- Josef Bayer and Ellen Brandner. 2004. Klitisiertes zu im Bairischen und Alemannischen. In *Morphologie und Syntax deutscher Dialekte und Historische Dialektologie des Deutschen: Beiträge zum 1. Kongress der Internationalen Gesellschaft für Dialektologie des Deutschen*.
- Emily M. Bender and Batya Friedman. 2018. [Data statements for natural language processing: Toward mitigating system bias and enabling better science](#). *Transactions of the Association for Computational Linguistics*, 6:587–604.
- Yevgeni Berzak, Yan Huang, Andrei Barbu, Anna Korhonen, and Boris Katz. 2016. [Anchoring and agreement in syntactic annotations](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2215–2224, Austin, Texas. Association for Computational Linguistics.
- Verena Blaschke, Barbara Kovačić, Siyao Peng, and Barbara Plank. 2024. [MaiBaam annotation guidelines](#). Technical report, LMU Munich. ArXiv 2403.05902.
- Verena Blaschke, Hinrich Schütze, and Barbara Plank. 2023a. [Does manipulating tokenization aid cross-lingual transfer? A study on POS tagging for non-standardized languages](#). In *Tenth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2023)*, pages 40–54, Dubrovnik, Croatia. Association for Computational Linguistics.
- Verena Blaschke, Hinrich Schütze, and Barbara Plank. 2023b. [A survey of corpora for Germanic low-resource languages and dialects](#). In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 392–414, Tórshavn, Faroe Islands. University of Tartu Library.
- Emanuel Borges Völker, Maximilian Wendt, Felix Hennig, and Arne Köhn. 2019. [HDT-UD: A very large Universal Dependencies treebank for German](#). In *Proceedings of the Third Workshop on Universal Dependencies (UDW, SyntaxFest 2019)*, pages 46–57, Paris, France. Association for Computational Linguistics.
- Lars Bülow, Philip C. Vergeiner, and Stephan Elspaß. 2021. [Structures of adnominal possession in Austria's traditional dialects: Variation and change](#). *Journal of Linguistic Geography*, 9(2):69–85.
- Branden Chan, Stefan Schweter, and Timo Möller. 2020. [German's next language model](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6788–6796, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Silvia Dal Negro and Simone Ciccolone. 2020. [KONTATTO: A laboratory for the study of language contact in South Tyrol](#). *Sociolinguistica*, 34(1):241–247.
- Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021. [Universal Dependencies](#). *Computational Linguistics*, 47(2):255–308.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jürg Fleischer. 2019. [Vergleichende Aspekte der deutschen Regionalsprachen: Syntax](#). In Joachim Herrgen and Jürgen Erich Schmidt, editors, *Deutsch*, pages 635–664. De Gruyter Mouton.
- Jennifer-Carmen Frey, Aivars Glaznieks, and Egon W. Stemle. 2015. [The DiDi corpus of South Tyrolean CMC data](#). In *Proceedings of the 2nd Workshop on Natural Language Processing for Computer-Mediated Communication / Social Media*, pages 1–6. GSCL.
- Marta Ghilardi. 2019. [Eliciting comparable spoken data in minor languages: first observations from the corpus Kontatti](#). *Suvremena lingvistika*, 45(88):231–246.
- Stefan Grünewald and Annemarie Friedrich. 2020. [Unifying the treatment of preposition-determiner contractions in German Universal Dependencies treebanks](#). In *Proceedings of the Fourth Workshop on Universal Dependencies (UDW)*

- 2020), pages 94–98, Barcelona, Spain (Online). Association for Computational Linguistics.
- Bruno Guillaume. 2021. [Graph matching and graph rewriting: GREW tools for corpus exploration, maintenance and conversion](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 168–175, Online. Association for Computational Linguistics.
- Johannes Heinecke. 2019. [ConlluEditor: a fully graphical editor for Universal Dependencies treebank files](#). In *Proceedings of the Third Workshop on Universal Dependencies (UDW, SyntaxFest 2019)*, pages 87–93, Paris, France. Association for Computational Linguistics.
- Georg F. K. Höhn. 2021. [Towards a consistent annotation of nominal person in Universal Dependencies](#). In *Proceedings of the Fifth Workshop on Universal Dependencies (UDW, SyntaxFest 2021)*, pages 75–83, Sofia, Bulgaria. Association for Computational Linguistics.
- Daniel Hole. 2002. [Agentive selbst in German](#). In *Proceedings of Sinn und Bedeutung 6*, pages 133–150.
- Dirk Hovy, Barbara Plank, and Anders Søgaard. 2014. [When POS data sets don't add up: Combatting sample bias](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 4472–4475, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Alexandra Lenz, Timo Ahlers, and Martina Werner. 2014. [Zur Dynamik bairischer Dialektsyntax – eine Pilotstudie](#). *Zeitschrift für Dialektologie und Linguistik*, 81(1).
- Ryan McDonald, Joakim Nivre, Yvonne Quirnbach-Brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith Hall, Slav Petrov, Hao Zhang, Oscar Täckström, Claudia Bedini, Núria Bertomeu Castelló, and Jungmee Lee. 2013. [Universal Dependency annotation for multilingual parsing](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 92–97, Sofia, Bulgaria. Association for Computational Linguistics.
- Ludwig Merkle. 1993. *Bairische Grammatik*, 5th edition. Heinrich Hugendubel Verlag, Munich.
- Aleksandra Miletic, Myriam Bras, Marianne Vergez-Couret, Louise Esher, Clamença Poujade, and Jean Sibille. 2020. [Building a Universal Dependencies treebank for Occitan](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2932–2939, Marseille, France. European Language Resources Association.
- Ann-Marie Moser. 2023. [The ups and downs of relative particles in German diachrony: On loss, grammaticalization, and standardization](#). *Journal of Historical Linguistics*, 13(3).
- Max Müller-Eberstein, Rob van der Goot, and Barbara Plank. 2021. [How universal is genre in Universal Dependencies?](#) In *Proceedings of the 20th International Workshop on Treebanks and Linguistic Theories (TLT, SyntaxFest 2021)*, pages 69–85, Sofia, Bulgaria. Association for Computational Linguistics.
- Siyao Peng. 2023. [Cross-Paragraph Discourse Structure in Rhetorical Structure Theory Parsing and Treebanking for Chinese and English](#). Ph.D. thesis, Georgetown University.
- Siyao Peng, Zihang Sun, Huangyan Shan, Marie Kolm, Verena Blaschke, Ekaterina Artemova, and Barbara Plank. 2024. [Sebastian, Basti, Wastl?! Recognizing named entities in Bavarian dialectal data](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation*.
- Karin Pittner. 1996. [Attraktion, Tilgung und Verbposition: Zur diachronen und dialektalen Variation beim Relativpronomen im Deutschen](#). In Ellen Brandner and Gisella Ferraresi, editors, *Language Change and Generative Grammar*, volume 7 of *Linguistische Berichte Sonderhefte*, pages 120–153. Westdeutscher Verlag.
- Martin Popel, Zdeněk Žabokrtský, and Martin Vojtík. 2017. [Udapi: Universal API for Universal Dependencies](#). In *Proceedings of the NoDaLiDa 2017 Workshop on Universal Dependencies (UDW 2017)*, pages 96–101, Gothenburg, Sweden. Association for Computational Linguistics.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. [Stanza: A Python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 101–108, Online. Association for Computational Linguistics.
- Anthony R. Rowley. 2011. [Bavarian: Successful dialect or failed language?](#) In Joshua Fishman and Ofelia Garcia, editors, *Handbook of Language and Ethnic Identity*, volume 2 (The Success-Failure Continuum in Language and

- Ethnic Identity Efforts), pages 299–309. Oxford University Press.
- Alessio Salomoni. 2017. [Toward a treebank collecting German aesthetic writings of the late 18th century](#). In *Proceedings of the Fourth Italian Conference on Computational Linguistics CLiC-it 2017*. Accademia University Press.
- Manuela Schönenberger and Eric Haeberli. 2019. [Ein geparstes und grammatisch annotiertes Korpus schweizerdeutscher Spontansprachdaten](#). In *Germanistische Linguistik*, volume 241–243, pages 79–104. Georg Olms Verlag.
- Janine Siewert, Yves Scherrer, and Jörg Tiedemann. 2021. [Towards a balanced annotated Low Saxon dataset for diachronic investigation of dialectal variation](#). In *Proceedings of the 17th Conference on Natural Language Processing (KONVENS 2021)*, pages 242–246, Düsseldorf, Germany. KONVENS 2021 Organizers.
- Janine Siewert, Yves Scherrer, and Martijn Wieling. 2022. [Low Saxon dialect distances at the orthographic and syntactic level](#). In *Proceedings of the 3rd Workshop on Computational Approaches to Historical Language Change*, pages 119–124, Dublin, Ireland. Association for Computational Linguistics.
- Natalia Silveira, Timothy Dozat, Marie-Catherine de Marneffe, Samuel Bowman, Miriam Connor, John Bauer, and Chris Manning. 2014. [A gold standard dependency corpus for English](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 2897–2904, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Milan Straka. 2018. [UDPipe 2.0 prototype at CoNLL 2018 UD shared task](#). In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 197–207, Brussels, Belgium. Association for Computational Linguistics.
- Rob van der Goot, Ibrahim Sharaf, Aizhan Imankulova, Ahmet Üstün, Marija Stepanović, Alan Ramponi, Siti Oryza Khairunnisa, Mamoru Komachi, and Barbara Plank. 2021a. [From masked language modeling to translation: Non-English auxiliary tasks improve zero-shot spoken language understanding](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2479–2497, Online. Association for Computational Linguistics.
- Rob van der Goot, Ahmet Üstün, Alan Ramponi, Ibrahim Sharaf, and Barbara Plank. 2021b. [Massive choice, ample tasks \(MaChAmp\): A toolkit for multi-task learning in NLP](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 176–197, Online. Association for Computational Linguistics.
- Helmut Weiß. 1998. *Syntax des Bairischen*. Max Niemeyer Verlag.
- Peter Wiesinger. 1983. [Die Einteilung der deutschen Dialekte](#). In Werner Besch, Ulrich Knoop, Wolfgang Putschke, and Herbert Ernst Wiegand, editors, *Dialektologie: Ein Handbuch zur deutschen und allgemeinen Dialektforschung*, pages 807–900. Walter de Gruyter.
- Miriam Winkler, Virginija Juozapaityte, Rob van der Goot, and Barbara Plank. 2024. [Slot and intent detection resources for Bavarian and Lithuanian: Assessing translations vs natural queries to digital assistants](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation*.
- Guillaume Wisniewski and François Yvon. 2019. [How bad are PoS tagger in cross-corpora settings? Evaluating annotation divergence in the UD project](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 218–227, Minneapolis, Minnesota. Association for Computational Linguistics.
- Amir Zeldes. 2017. [The GUM corpus: Creating multilayer resources in the classroom](#). *Language Resources and Evaluation*, 51(3):581–612.
- Amir Zeldes and Nathan Schneider. 2023. [Are UD treebanks getting more consistent? A report card for English UD](#). In *Proceedings of the Sixth Workshop on Universal Dependencies (UDW, GURT/SyntaxFest 2023)*, pages 58–64, Washington, D.C. Association for Computational Linguistics.
- Daniel Zeman, Martin Popel, Milan Straka, Jan Hajič, Joakim Nivre, Filip Ginter, Juhani Luotolahti, Sampo Pyysalo, Slav Petrov, Martin Potthast, Francis Tyers, Elena Badmaeva, Memduh Gokirmak, Anna Nedoluzhko, Silvie Cinková, Jan Hajič jr., Jaroslava Hlaváčová, Václava Kettnerová, Zdeňka Uřešová, Jenna Kanerva,

Stina Ojala, Anna Missilä, Christopher D. Manning, Sebastian Schuster, Siva Reddy, Dima Taji, Nizar Habash, Herman Leung, Marie-Catherine de Marneffe, Manuela Sanguinetti, Maria Simi, Hiroshi Kanayama, Valeria de Paiva, Kira Droganova, Héctor Martínez Alonso, Çağrı Çöltekin, Umut Sulubacak, Hans Uszkoreit, Vivien Macketanz, Aljoscha Burchardt, Kim Harris, Katrin Marheinecke, Georg Rehm, Tolga Kayadelen, Mohammed Attia, Ali Elkahky, Zhuoran Yu, Emily Pitler, Saran Lertpradit, Michael Mandl, Jesse Kirchner, Hector Fernandez Alcalde, Jana Strnadová, Esha Banerjee, Ruli Manurung, Antonio Stella, Atsuko Shimada, Sookyong Kwak, Gustavo Mendonça, Tatiana Lando, Rattima Nitisaroj, and Josie Li. 2017. [CoNLL 2017 shared task: Multilingual parsing from raw text to Universal Dependencies](#). In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 1–19, Vancouver, Canada. Association for Computational Linguistics.

7. Language Resource References

IDS. n. d. *Deutsche Mundarten: Zwierner-Korpus*. Datenbank für gesprochenes Deutsch (DGD).

Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia Ackermann, Noëmi Aepli, Hamid Aghaei, Željko Agić, Amir Ahmadi, Lars Ahrenberg, Chika Kennedy Ajede, Salih Furkan Akkurt, Gabrielė Aleksandravičiūtė, Ika Alfina, Avner Algom, Khalid Alnajjar, Chiara Alzetta, Erik Andersen, Lene Antonsen, Tatsuya Aoyama, Katya Aplonova, Angelina Aquino, Carolina Aragon, Glyd Aranes, Maria Jesus Aranzabe, Bilge Nas Arcan, Pórunn Arnardóttir, Gashaw Arutie, Jessica Naraiswari Arwidarasti, Masayuki Asahara, Katla Ásgeirsdóttir, Deniz Baran Aslan, Cengiz Asmazoğlu, Luma Ateyah, Furkan Atmaca, Mohammed Attia, Aitziber Atutxa, Liesbeth Augustinus, Mariana Avelãs, Elena Badmaeva, Keerthana Balasubramani, Miguel Ballesteros, Esha Banerjee, Sebastian Bank, Verginica Barbu Mititelu, Starkaður Barkarson, Rodolfo Basile, Victoria Basmov, Colin Batchelor, John Bauer, Seyyit Talha Bedir, Shabnam Behzad, Juan Belieni, Kepa Bengoetxea, İbrahim Benli, Yifat Ben Moshe, Gözde Berk, Riyaz Ahmad Bhat, Erica Biagetti, Eckhard Bick, Agnė Bielinskienė, Kristín Bjarnadóttir, Rogier Blokland, Victoria Bobicev, Loïc Boizou, Emanuel Borges Völker, Carl Börstell, Cristina Bosco, Gosse Bouma, Sam Bowman, Adriane Boyd, Anouck Braggaar, António Branco, Kristina

Brokaitė, Aljoscha Burchardt, Marisa Campos, Marie Candito, Bernard Caron, Gauthier Caron, Catarina Carvalheiro, Rita Carvalho, Lauren Cassidy, Maria Clara Castro, Sérgio Castro, Tatiana Cavalcanti, Gülşen Cebiroğlu Eryiğit, Flavio Massimiliano Cecchini, Giuseppe G. A. Celano, Slavomír Čéplö, Neslihan Cesur, Savas Cetin, Özlem Çetinoğlu, Fabricio Chalub, Liyanage Chamila, Shweta Chauhan, Ethan Chi, Taishi Chika, Yongseok Cho, Jinho Choi, Jayeol Chun, Juyeon Chung, Alessandra T. Cignarella, Silvie Cinková, Aurélie Collomb, Çağrı Çöltekin, Miriam Connor, Claudia Corbetta, Daniela Corbetta, Francisco Costa, Marine Courtin, Benoît Crabbé, Mihaela Cristescu, Vladimir Cvetkoski, Ingerid Løyning Dale, Philemon Daniel, Elizabeth Davidson, Leonel Figueiredo de Alencar, Mathieu Dehouck, Martina de Laurentiis, Marie-Catherine de Marneffe, Valeria de Paiva, Mehmet Oguz Derin, Elvis de Souza, Arantza Diaz de Ilarraza, Carly Dickerson, Arawinda Dinakaramani, Elisa Di Nuovo, Bamba Dione, Peter Dirix, Kaja Dobrovoljc, Adrian Doyle, Timothy Dozat, Kira Droganova, Magali Sanches Duran, Puneet Dwivedi, Christian Ebert, Hanne Eckhoff, Masaki Eguchi, Sandra Eiche, Marhaba Eli, Ali Elkahky, Binyam Ephrem, Olga Erina, Tomaž Erjavec, Farah Essaidi, Aline Etienne, Wograine Evelyn, Sidney Facundes, Richárd Farkas, Federica Favero, Jannatul Ferdaousi, Marília Fernanda, Hector Fernandez Alcalde, Amal Fethi, Jennifer Foster, Theodorus Fransen, Cláudia Freitas, Kazunori Fujita, Katarína Gajdošová, Daniel Galbraith, Federica Gamba, Marcos Garcia, Moa Gärdenfors, Fabrício Ferraz Gerardi, Kim Gerdes, Luke Gessler, Filip Ginter, Gustavo Godoy, Iakes Goenaga, Koldo Gojenola, Memduh Gökırmak, Yoav Goldberg, Xavier Gómez Guinovart, Berta González Saavedra, Bernadeta Griciūtė, Matias Gironi, Loïc Grobol, Normunds Grūzītis, Bruno Guillaume, Kirian Guiller, Céline Guillot-Barbance, Tunga Güngör, Nizar Habash, Hinrik Hafsteinsson, Jan Hajič, Jan Hajič jr., Mika Hämäläinen, Linh Hà Mỹ, Na-Rae Han, Muhammad Yudistira Hanifmuti, Takahiro Harada, Sam Hardwick, Kim Harris, Dag Haug, Johannes Heinecke, Oliver Hellwig, Felix Hennig, Barbora Hladká, Jaroslava Hlaváčová, Florinel Hociung, Petter Hohle, Yidi Huang, Marivel Huerta Mendez, Jena Hwang, Takumi Ikeda, Anton Karl Ingason, Radu Ion, Elena Irimia, Qlájídé Ishola, Artan Islamaj, Kaoru Ito, Sandra Jagodzińska, Siratun Jannat, Tomáš Jelínek, Apoorva Jha, Katharine Jiang, Anders Johannsen, Hildur Jónsdóttir, Fredrik Jørgensen, Markus Juutinen, Hüner Kaşıkara, Nadezhda Kabaeva, Sylvain

Kahane, Hiroshi Kanayama, Jenna Kanerva, Neslihan Kara, Ritván Karahóga, Andre Kåsen, Tolga Kayadelen, Sarveswaran Kengatharayer, Václava Kettnerová, Lilit Kharatyan, Jesse Kirchner, Elena Klementieva, Elena Klyachko, Petr Kocharov, Arne Köhn, Abdullatif Köksal, Kamil Kopacewicz, Timo Korkiakangas, Mehmet Köse, Alexey Koshevoy, Natalia Kotsyba, Jolanta Kovalevskaitė, Simon Krek, Parameswari Krishnamurthy, Sandra Kübler, Adrian Kuqi, Oğuzhan Kuyrukçu, Aslı Kuzgun, Sookyong Kwak, Kris Kyle, Käbi Laan, Veronika Laippala, Lorenzo Lambertino, Tatiana Lando, Septina Dian Larasati, Alexei Lavrentiev, John Lee, Phương Lê Hồng, Alessandro Lenci, Saran Lertpradit, Herman Leung, Maria Levina, Lauren Levine, Cheuk Ying Li, Josie Li, Keying Li, Yixuan Li, Yuan Li, KyungTae Lim, Bruna Lima Padovani, Yi-Ju Jessica Lin, Krister Lindén, Yang Janet Liu, Nikola Ljubešić, Irina Lobzhanidze, Olga Logina, Lucelene Lopes, Stefano Lusito, Andry Luthfi, Mikko Luukko, Olga Lyashevskaya, Teresa Lynn, Vivien Mackentanz, Menel Mahamdi, Jean Maillard, Ilya Makarchuk, Aibek Makazhanov, Michael Mandl, Christopher Manning, Ruli Manurung, Büşra Marşan, Cătălina Mărănduc, David Mareček, Katrin Marheinecke, Stella Markantonatou, Héctor Martínez Alonso, Lorena Martín Rodríguez, André Martins, Cláudia Martins, Jan Mašek, Hiroshi Matsuda, Yuji Matsumoto, Alessandro Mazzei, Ryan McDonald, Sarah McGuinness, Gustavo Mendonça, Tatiana Merzhevich, Niko Miekka, Aaron Miller, Karina Mischenkova, Anna Missilä, Cătălin Mititelu, Maria Mitrofan, Yusuke Miyao, AmirHossein Mojiri Froushani, Judit Molnár, Amirsaeid Moloodi, Simonetta Montemagni, Amir More, Laura Moreno Romero, Giovanni Moretti, Shinsuke Mori, Tomohiko Morioka, Shigeki Moro, Bjartur Mortensen, Bohdan Moskalevskyi, Kadri Muischnek, Robert Munro, Yugo Murawaki, Kaili Müürisep, Pinkey Nainwani, Mariam Nakhlé, Juan Ignacio Navarro Horñiacek, Anna Nedoluzhko, Gunta Nešpore-Bērzkalne, Manuela Nevaci, Lương Nguyễn Thị, Huyền Nguyễn Thị Minh, Yoshihiro Nikaido, Vitaly Nikolaev, Rattima Nitisaroj, Alireza Nourian, Maria das Graças Volpe Nunes, Hanna Nurmi, Stina Ojala, Atul Kr. Ojha, Hulda Óladóttir, Adédayò Olúòkun, Mai Omura, Emeka Onwuegbuzia, Noam Ordan, Petya Osenova, Robert Östling, Lilja Øvrelid, Şaziye Betül Özateş, Merve Özçelik, Arzucan Özgür, Balkız Öztürk Başaran, Teresa Paccosi, Alessio Palmero Aprosio, Anastasia Panova, Thiago Alexandre Salgueiro Pardo, Hyunji Hayley Park, Niko Partanen, Elena Pascual, Marco Passarotti, Agnieszka

Patejuk, Guilherme Paulino-Passos, Giulia Pedonese, Angelika Peljak-Łapińska, Siyao Peng, Siyao Logan Peng, Rita Pereira, Sílvia Pereira, Cenel-Augusto Perez, Natalia Perkova, Guy Perrier, Slav Petrov, Daria Petrova, Andrea Peverelli, Jason Phelan, Claudel Pierre-Louis, Jussi Piitulainen, Yuval Pinter, Clara Pinto, Rodrigo Pintucci, Tommi A Pirinen, Emily Pitler, Magdalena Plamada, Barbara Plank, Thierry Poibeau, Larisa Ponomareva, Martin Popel, Lauma Pretkalniņa, Sophie Prévost, Prokopis Prokopidis, Adam Przepiórkowski, Robert Pugh, Tiina Puolakainen, Sampo Pyysalo, Peng Qi, Andreia Querido, Andriela Rääbis, Alexandre Rademaker, Mizanur Rahoman, Taraka Rama, Loganathan Ramasamy, Carlos Ramisch, Joana Ramos, Fam Rashel, Mohammad Sadegh Rasooli, Vinit Ravishankar, Livy Real, Petru Rebeja, Siva Reddy, Mathilde Regnault, Georg Rehm, Arij Riabi, Ivan Riabov, Michael Rießler, Erika Rimkutė, Larissa Rinaldi, Laura Rituma, Putri Rizqiyah, Luisa Rocha, Eiríkur Rögnvaldsson, Ivan Roksandic, Mykhailo Romanenko, Rudolf Rosa, Valentin Roşca, Davide Rovati, Ben Rozonoyer, Olga Rudina, Jack Rueter, Kristján Rúnarsson, Shoal Sadde, Pegah Safari, Aleks Sahala, Shadi Saleh, Alessio Salomoni, Tanja Samardžić, Stephanie Samson, Manuela Sanguinetti, Ezgi Saniyar, Dage Särg, Marta Sartor, Mitsuya Sasaki, Baiba Saulīte, Agata Savary, Yanin Sawanakunanon, Shefali Saxena, Kevin Scannell, Salvatore Scarlata, Emmanuel Schang, Nathan Schneider, Sebastian Schuster, Lane Schwartz, Djamé Seddah, Wolfgang Seeker, Mojgan Seraji, Syeda Shahzadi, Mo Shen, Atsuko Shimada, Hiroyuki Shirasu, Yana Shishkina, Muh Shohibussirri, Maria Shvedova, Janine Siewert, Einar Freyr Sigurðsson, João Silva, Aline Silveira, Natalia Silveira, Sara Silveira, Maria Simi, Radu Simionescu, Katalin Simkó, Mária Šimková, Haukur Barri Símónarson, Kiril Simov, Dmitri Sitchinava, Ted Sither, Maria Skachedubova, Aaron Smith, Isabela Soares-Bastos, Per Erik Solberg, Barbara Sonnenhauser, Shafi Sourov, Rachele Sprugnoli, Vivian Stamou, Steinþór Steingrímsson, Antonio Stella, Abishek Stephen, Milan Straka, Emmett Strickland, Jana Strnadová, Alane Suhr, Yogi Lesmana Sulestio, Umut Sulubacak, Shingo Suzuki, Daniel Swanson, Zsolt Szántó, Chihiro Taguchi, Dima Taji, Fabio Tamburini, Mary Ann C. Tan, Takaaki Tanaka, Dipta Tanaya, Mirko Tavoni, Samson Tella, Isabelle Tellier, Marinella Testori, Guillaume Thomas, Sara Tonelli, Liisi Torga, Marsida Toska, Trond Trosterud, Anna Trukhina, Reut Tsarfaty, Utku Türk, Francis Tyers, Sveinbjörn Þórðarson,

Vilhjálmur Þorsteinsson, Sumire Uematsu, Roman Untilov, Zdeňka Urešová, Larraitz Uria, Hans Uszkoreit, Andrius Utka, Elena Vagnoni, Sowmya Vajjala, Socrates Vak, Rob van der Goot, Martine Vanhove, Daniel van Niekerk, Gertjan van Noord, Viktor Varga, Uliana Vedenina, Giulia Venturi, Eric Villemonte de la Clergerie, Veronika Vincze, Natalia Vlasova, Aya Wakasa, Joel C. Wallenberg, Lars Wallin, Abigail Walsh, Jonathan North Washington, Maximilian Wendt, Paul Widmer, Shira Wigderson, Sri Hartati Wijono, Vanessa Berwanger Wille, Seyi Williams, Mats Wirén, Christian Wittern, Tsegay Woldemariam, Tak-sum Wong, Alina Wróblewska, Qishen Wu, Mary Yako, Kayo Yamashita, Naoki Yamazaki, Chunxiao Yan, Koichi Yasuoka, Marat M. Yavrumyan, Arife Betül Yenice, Olcay Taner Yıldız, Zhuran Yu, Arlisa Yuliawati, Zdeněk Žabokrtský, Shorouq Zahra, Amir Zeldes, He Zhou, Hanzhi Zhu, Yilun Zhu, Anna Zhuravleva, and Rayan Ziane. 2023. [Universal Dependencies 2.13](#). LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

A. Data Statement

Header

- *Dataset title:* MaiBaam
- *Dataset curator(s):* Verena Blaschke, Barbara Kovačić, Siyao Peng, Barbara Plank
- *Dataset version:* 1.0 (UD version 2.14)
- *Dataset citation:* MaiBaam should be cited by citing this article.
- *Data statement authors:* Verena Blaschke, Barbara Kovačić, Siyao Peng, Hinrich Schütze, Barbara Plank
- *Data statement version:* 1.0
- *Data statement citation and DOI:* To cite this data statement, please cite this publication.
- *Links to versions of this data statement in other languages:* —

Executive Summary MaiBaam is a manually annotated dependency treebank for Bavarian. It contains 15k tokens and is annotated with part-of-speech tags and syntactic dependencies according to Universal Dependency guidelines.²⁴

²⁴<https://universaldependencies.org/bar/>

MaiBaam encompasses diverse text genres (wiki articles and discussions, grammar examples, fiction, and commands for virtual assistants) and dialects from the North, Central and South Bavarian areas as well as the dialectal transition areas in between. In addition, it provides sentence-level genre and dialect metadata.

Curation Rationale The purpose of MaiBaam is to allow research on computational methods for processing non-standardized language data, including the evaluation of cross-lingual transfer setups (given the large amount of German and other data also annotated according to UD guidelines). Furthermore, it allows researching syntactic structures of Bavarian (on their own, and in contrast to the other Germanic languages covered by UD). Our goal is to represent as many Bavarian dialects and as many text genres as possible given the availability and licensing of such data. Each data instance is a sentence, annotated with part-of-speech tags and syntactic dependencies.

Documentation for Source Datasets

MaiBaam contains sentences from the South Tyrolean (Winkler et al., 2024) and Central Bavarian (Winkler et al., 2024) versions of xSID (CC BY-SA 4.0 International) as well as additional data from Winkler et al. (2024) (to be released soon, likely under the same license), Tatoeba (CC-BY 2.0 FR), Wikipedia (CC BY-SA 4.0 International), and UD’s Cairo CILing Corpus (no license).

Language Varieties MaiBaam contains Bavarian data (ISO 639-3: [bar](#), Glottocode: [bava1246](#), BCP-47: [bar-DE](#), [bar-AT](#), [bar-IT](#)) from the North, Central and South Bavarian areas as well as the transition areas between North and Central Bavarian and Central and South Bavarian. Details are in Table 3.

Speaker Demographic We do not have detailed demographic information on the speakers whose data we include in MaiBaam, with the exception of differently granular geographic and/or dialectal information (this applies to 80% of the utterances). About 17% of the tokens (13% of the sentences) belong to dialects spoken in larger cities (Munich, Vienna, Salzburg). Details on the geographic and dialectal distribution can be found in Table 3.

Annotator Demographic All involved in providing Bavarian translations and helping us with questions about Bavarian words/sentences or linguistic structures are native speakers of Bavarian (two Central Bavarian speakers from Bavaria, one South Bavarian speaker from South Tyrol).

The annotator is a native speaker of German and a (non-Bavarian) Upper German dialect. The annotation guidelines were created and refined by a native speaker of German familiar with German dialectology, a learner of German with experience in annotating UD treebanks, and the annotator. They were also reviewed by a native speaker of a Bavarian dialect.

Everybody involved in this project, except for one of the Bavarian informants, has a background in (computational) linguistics. Additional details are in Section 3.2.

Speech Situation and Text Characteristics

- *Time of linguistic activity*: Wikipedia articles and discussions: unknown between 2006–2024; Tatoeba sentences: 2013–2022; Fairy tales: 2018; xSID: translated in 2021 (South Tyrolean) and 2023 (data from Bavaria); other virtual assistant data: 2023; CICLing: translated in 2023
- *Date of data collection*: 2023–2024
- *Modality*: Written
- *Synchronous vs. asynchronous interaction*: Asynchronous
- *Scripted/edited vs. spontaneous*: Presumably edited in most cases. The CICLing and xSID sentences are translations, as are many of the Tatoeba sentences. Some of the sentences from Wikipedia articles or from fairy tales might be translations.
- *Speakers' intended audience*: The sentences in xSID are queries for a hypothetical digital assistant. The linguistic example sentences from CICLing and Wikipedia articles are for people interested in linguistics. All other data are for an audience of internet users who are interested in reading Bavarian content, be they themselves speakers of Bavarian or not.
- *Genre*: See Section 3.1 and Table 2.
- *Topic*: Various. The wiki articles include locations, traditions/customs, food and entertainment/media, among other topics.
- *Non-linguistic context*: —

Preprocessing and Data Formatting We manually replace usernames mentioned in Wikipedia discussions with `USERNAME`. We ignore the original text formatting choices (italics, boldface, typeface). We do not include the raw, unannotated data in the dataset. The dataset adheres to CoNLL-U formatting.

Capture Quality No known issues.

Limitations We cannot verify that all sentences were written by competent speakers of Bavarian.

Metadata

- *License*: [CC BY-SA 4.0 International](#)
- *Annotation guidelines*: [Blaschke et al. \(2024\)](#)
- *Annotation process*: See Section 3.2.
- *Dataset quality metrics*: —
- *Errata*: None so far. Please report errors by contacting the authors or opening an issue at github.com/UniversalDependencies/UD_Bavarian-MaiBaam/issues.

Disclosures and Ethical Review We only collected and annotated data that were shared under licences that explicitly permit adapting and re-sharing the data. Everyone involved in annotating and translating data and everyone we consulted with questions about Bavarian was hired and compensated according to local standards. An institutional ethics review process was not accessible at the time of dataset creation.

There are no conflicts of interest. This research is supported by the ERC Consolidator Grant DIALECT 101043235. We also gratefully acknowledge partial funding by the European Research Council (ERC #740516).

Other —

Glossary —

About this document A data statement is a characterization of a dataset that provides context to allow developers and users to better understand how experimental results might generalize, how software might be appropriately deployed, and what biases might be reflected in systems built on the software.

This data statement was written based on the template for the Data Statements Version 2 Schema. The template was prepared by Angelina McMillan-Major, Emily M. Bender, and Batya Friedman, and can be found at <http://techpolicylab.uw.edu/data-statements>.

B. POS Tag Distributions

Table 6 provides an overview of the part-of-speech tags in MaiBaam, and Table 7 shows the distribution of dependency relations.

Part-of-speech tag		# Tokens	%	TTR	Top 5 most frequent tokens
NOUN	Noun	2 269	15.1	0.65	<i>Wecka</i> ‘alarm’, <i>Kafää</i> ‘coffee’, <i>Mingara</i> ‘Munich citizen’, <i>Leid</i> ‘people’, <i>See</i> ‘lake’
PUNCT	Punctuation	2 105	14.0	0.01	. , ” ? (
DET	Determiner	1 946	13.0	0.11	<i>m</i> , <i>da</i> , <i>a</i> , <i>de</i> , <i>an</i>
VERB	Verb	1 458	9.7	0.65	<i>gibt</i> ‘[there] is’, <i>hod</i> ‘has’, <i>Erinner</i> ‘remind’, <i>sogt</i> ‘says’, <i>gsogt</i> ‘said’
ADP	Adposition	1 417	9.4	0.14	<i>in</i> , <i>vo</i> , <i>i</i> , <i>auf</i> , <i>mit</i>
ADV	Adverb	1 206	8.0	0.39	<i>aa</i> ‘also’, <i>so</i> ‘so’, <i>a</i> ‘also’, <i>no</i> ‘still’, <i>do</i> ‘there’
PRON	Pronoun	1 127	7.5	0.11	<i>ma</i> ‘we, me’, <i>i</i> ‘I’, <i>s</i> ‘it, she/her, you.PL’, <i>I</i> ‘I’, <i>des</i> ‘this/that’
AUX	Auxiliary	926	6.2	0.21	<i>is</i> ‘is’, <i>hod</i> ‘has’, <i>san</i> ‘[we/they] are’, <i>hob</i> ‘[I] have’, <i>hom</i> ‘[we] have’
ADJ	Adjective	799	5.3	0.83	<i>neie</i> ‘new’, <i>guat</i> ‘good’, <i>Soizburga</i> ‘of Salzburg’, <i>guad</i> ‘good’, <i>gaunzn</i> ‘entire’,
PROPN	Proper noun	550	3.7	0.66	<i>Minga</i> ‘Munich’, <i>Thomas</i> , <i>Tom</i> , <i>Gretl</i> , <i>Hansl</i>
CCONJ	Coordinating conjunction	380	2.5	0.08	<i>und</i> ‘and’, <i>oda</i> ‘or’, <i>owa</i> ‘but’, <i>oder</i> ‘or’, <i>Und</i> ‘and’
SCONJ	Subordinating conjunction	341	2.3	0.23	<i>dass</i> ‘that’, <i>das</i> ‘that’, <i>wia</i> ‘than’, <i>wej</i> ‘REL’, <i>wos</i> ‘REL’
NUM	Numeral	240	1.6	0.58	<i>5</i> , <i>zwoa</i> ‘two’, <i>4</i> , <i>6 drei</i> ‘three’
PART	Particle	165	1.1	0.13	<i>ned</i> ‘not’, <i>net</i> ‘not’, <i>zu</i> ‘INF’, <i>niat</i> ‘not’, <i>nim</i> ‘not [anymore]’
X	Other	64	0.4	0.89	<i>e</i> , <i>d</i> , <i>Schuhplattler</i> * ‘[dance]’, <i>München</i> * ‘Munich’, <i>miŋ(:)e</i>
INTJ	Interjection	23	0.2	0.70	<i>Gö</i> , <i>gö</i> , <i>Ja</i> ‘yes’, <i>Bfiade</i> ‘bye’, <i>Seavas</i> ‘hi/bye’
SYM	Symbol	7	0.0	1.00	*, †, %, <, :-)

Table 6: **POS tag statistics.** For each POS tag, we provide the absolute and relative (% , in percent) number of tokens, the type-token ratio (TTR) and the most frequent tokens. The asterisk* denotes German words that are clearly presented as non-Bavarian material in a given sentence.

Relation	Abs.	%	Relation	Abs.	%	Relation	Abs.	%
punct	2105	14.0	appos	176	1.2	nsubj:pass	66	0.4
det	1746	11.6	advcl	174	1.2	fixed	41	0.3
advmod	1403	9.3	nummod	140	0.9	acl	31	0.2
case	1329	8.8	flat	126	0.8	orphan	30	0.2
nsubj	1128	7.5	expl	110	0.7	discourse	23	0.2
root	1070	7.1	expl:pv	107	0.7	advcl:relcl	21	0.1
obl	798	5.3	obl:arg	105	0.7	vocative	16	0.1
obj	670	4.5	compound:prt	103	0.7	compound	15	0.1
aux	575	3.8	ccomp	99	0.7	obl:agent	11	0.1
amod	468	3.1	det:poss	90	0.6	goeswith	8	0.1
conj	467	3.1	acl:relcl	88	0.6	csubj	7	0.0
nmod	446	3.0	parataxis	84	0.6	dislocated	7	0.0
cc	377	2.5	aux:pass	81	0.5	reparandum	3	0.0
mark	342	2.3	xcomp	70	0.5	dep	2	0.0
cop	265	1.8						

Table 7: **Dependency relation statistics.** For each dependency relation, we provide the absolute (#) and relative (% , in percent) number of occurrences.