

Rule-Based Clause-Level Morphology for Multiple Languages

Tillmann Dönicke

Göttingen Centre for Digital Humanities

University of Göttingen

tillmann.doenicke@uni-goettingen.de

Abstract

This paper describes an approach for the morphosyntactic analysis of clauses, including the analysis of composite verb forms and both overt and covert pronouns. The approach uses grammatical rules for verb inflection and clause-internal word agreement to compute a clause’s morphosyntactic features from the morphological features of the individual words. The approach is tested for eight typologically diverse languages in the 1st Shared Task on Multilingual Clause-Level Morphology, where it achieves F1 scores between 79% and 99% (94% in average).

1 Introduction

Until recently the prediction of clause-level morphological / morphosyntactic features has been approached for a few individual languages only (see Žáčková et al. (2000) for Czech, Choudhary et al. (2014) for Hindi, Faro and Pavone (2015) for Italian, Ramm et al. (2017) for English, French and German, Myers and Palmer (2019) for English revisited, and Dönicke (2020) for German revisited). Most of the approaches are rule-based, first of all because annotated training data barely exists. On the other hand, it seems intuitive to approach this task in a rule-based manner, since morphosyntax follows strict grammatical rules (as opposed to heuristics) that can be implemented by a linguist. The first work to our knowledge which considers multiple and typologically diverse languages at a time is that of Dönicke (2021), who presents a cross-linguistic algorithm for composite-verb analysis and implements it for 11 languages, but refrains from evaluating the approach due to the lack of annotated gold data. The 1st Shared Task on Multilingual Clause-Level Morphology tackles this lack of data and provides data sets for eight typologically diverse languages. We re-implement and extend Dönicke (2021)’s algorithm for the shared

task (Section 3), evaluate it (Section 4) and discuss its advantages and shortcomings (Section 5).

2 Shared Task and Data

The 1st Shared Task on Multilingual Clause-Level Morphology (Task 3 Analysis) provides data sets for eight languages. Training sets (10,000 samples each) and development sets (2,000 samples each) for six languages were released first, and test sets (1,000 samples each) as well as all sets for two surprise languages (Spanish and Swahili) were released two weeks before the system submission deadline. Each sample consists of a short sentence and a gold analysis. The sentence consists of a single clause and contains one verb form that can be simple (e.g. *he looks*) or composite (e.g. *he had not been looking*) as well as pronouns, adpositions and a sentence-final punctuation mark. The gold analysis consists of the main verb’s lemma, the analysis of the verb form and the analyses of all pronouns, both overtly expressed pronouns (as in *he looks*) and covertly expressed ones (as in *∅ look!*). The analyses are represented with UniMorph features (Sylak-Glassman, 2016). The task was to predict an analysis for an input sentence. Since the test sets were provided without gold analyses, the submission and evaluation of systems was performed via CodaLab.¹

3 Method

3.0 Motivation

Computing the morphosyntactic analysis of a clause can be modeled as a mapping from word-level morphological features to clause-level morphological features. This process follows grammatical rules, in particular (language-specific) rules for verb inflection and (language-independent) rules of agreement between words in a clause. Figure 1

¹https://codalab.lisn.upsaclay.fr/competitions/6830?secret_key=44e813c2-96c8-4889-b0fc-24dbe83ad2c6

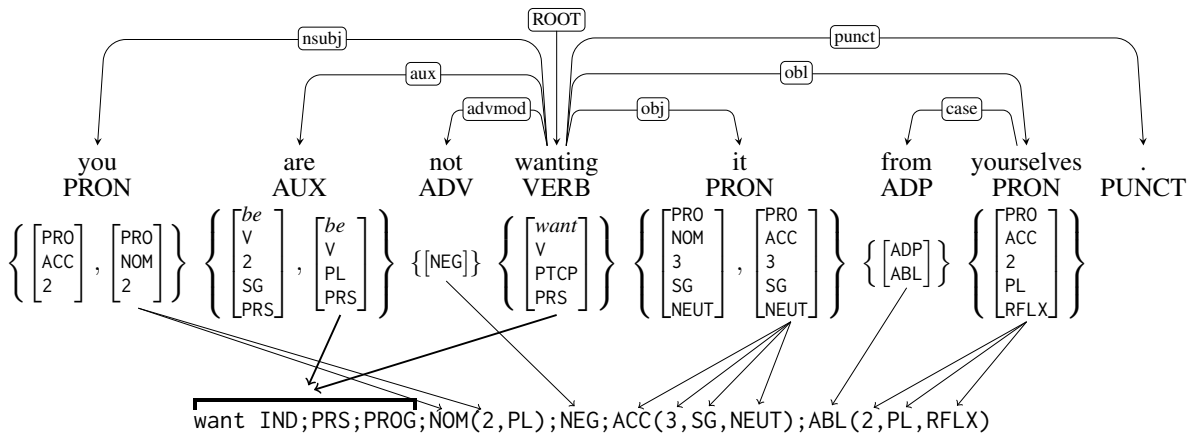


Figure 1: Mapping from word-level features to clause-level features for an English clause.

illustrates this for an English example sentence with a dependency tree on top and morphological analyses below each word. Some of the words are morphologically ambiguous and thus have more than one morphological analysis. The inflectional paradigm of English tells us that a finite present-tense (PRS) form of *be* and the present (PRS) participle (PTCP) of another verb expresses the indicative (IND) present (PRS) progressive (PROG) form of the latter verb, here *want*. To find the subject, it is to find a pronoun with nominative case (NOM), which could either be *you* or *it*. Since the subject has to agree with the finite verb *are* in person and number, the subject can only be *you*. In the consequence, *it* must be an object and cannot have nominative case, hence it receives the accusative (ACC) analysis. The third pronoun, *yourselves*, is reflexive (RFLX) and must therefore agree in person and number with another pronoun in the clause, where the only candidate is *you*. Because of the agreement of *you* and *yourselves*, *you* (which has no morphological number feature) has to be analyzed as plural (PL) and copies this features from *yourselves*. The adposition *from*, which is syntactically governed by *yourselves*, overrides the morphological case of the pronoun.

Our entirely rule-based approach analyzes a clause in a very similar manner as in the example. The following subsections give an overview of the processing steps that an input sentence goes through to compute the output analysis. There are also two examples for French input sentences in the appendix. Further details can also be found in the documented source code.²

3.1 Preprocessing

All languages are preprocessed with spaCy.³ We use the pretrained models for French, Russian and Spanish, and trained new models on the Universal Dependencies (UD) treebanks (Zeman et al., 2022) for German (HDT), English (GUM), Hebrew (IAHLTwiki) and Turkish (Kenet). To improve the tokenization of spaCy, the raw text is preprocessed for some languages. For English, contractions are converted to full forms (e.g. *won't* \mapsto *will not*) using the Python package *contractions*⁴ and some additional conversions using regular expressions. Similarly, hyphenated contractions are converted to full forms for French by replacing - and -t- with a space (e.g. *regarde-t-il* \mapsto *regarde il*, *m'avaient-elles* \mapsto *m'avaient elles*). Since we could only train a spaCy model for unvocalized Hebrew, vocalized Hebrew is converted to unvocalized Hebrew using *unikud*⁵ before processing it with spaCy and afterwards replaced back with the original tokens.

Unfortunately, even for sentences as simple as in the shared task's data, spaCy makes errors in all processing steps: part-of-speech (POS) tagging, lemmatization and parsing. We fix the most errors with a mix of language-independent and language-specific rules. First, we look up the word-level analysis for every token in UniMorph (see Section 3.2 below) and overwrite the POS tag and/or lemma assigned by spaCy with that from UniMorph if it is unambiguous. Then, we apply some fixes to the parse tree according to the POS tags.

As there is no UD treebank for Swahili, it is

²<https://gitlab.gwdg.de/tillmann.doenicke/mrl2022-tmvm>

³<https://spacy.io/>

⁴<https://pypi.org/project/contractions/>

⁵<https://pypi.org/project/unikud/>

also not possible to train a spaCy model for the language. Here, we directly set the POS tags and lemmas according to the word-level analysis. As far as it concerns the shared task, parsing is not necessary for Swahili since we only need parses to connect adpositions or verbal particles with their heads, and Swahili has no such multi-word constructions.

3.2 Word-Level Analysis

We use UniMorph for word-level morphological analysis. UniMorph provides large word lists with POS and morphological analysis,⁶ however, only for verbs, nouns and adjectives. We therefore added analyses for pronouns, adpositions and in some languages also for auxiliary verbs (e.g. forms of *be* in English) if they are missing in the UniMorph files. Table 1 shows the number of word form analyses in the files from UniMorph and our extensions. Since UniMorph does not provide resources for Swahili,⁷ we only added analyses for the six personal pronouns and assume that every other input word is a verb, which we then analyze with the regular expression⁸

$$(Prefix)?(Subject)?(Tense)?(Object)? \\ (Stem)(Vowel),$$

where *Prefix*, *Subject*, *Tense* and *Object* can be any morpheme from an according predefined dictionary, e.g. $Subject \in \{ni : \left\{ \left[\begin{smallmatrix} 1 \\ SG \end{smallmatrix} \right] \right\}, u : \left\{ \left[\begin{smallmatrix} 2 \\ SG \end{smallmatrix} \right], \left[\begin{smallmatrix} 3 \\ SG \\ M_MI \end{smallmatrix} \right], \left[\begin{smallmatrix} 3 \\ SG \\ U \end{smallmatrix} \right] \right\}, \dots \}$, $Stem = .+?[aeiou]+[^\wedge aeiou]+$ and $Vowel = ([aeiou]?[aeiou])|((([aeiou]l)?(ia|ea)))$.

The word-level analyses are filtered and post-corrected in some cases depending on the context and using language-specific rules. For example, if the Spanish (usually reflexive) pronoun *se* precedes *la*, *las*, *lo* or *los*, it could also be a replacement for *le* or *les*, so the analyses of *le* and *les* are added.

3.3 Clause-Level Analysis

Composite verb forms are analyzed with the algorithm from Dönicke (2020, 2021), which maps the word-level features of the involved verbs to clause-level features. The algorithm itself is mostly language-independent and its application to dif-

⁶<https://github.com/unimorph/>

⁷UniMorph provides resources for Congo Swahili, another Swahili variant than that in the shared task’s data.

⁸The regular expression is mainly based on the *Swahili Cheat Sheet* which can be found at <https://www.swahilicheatsheet.com/>.

Language	UM	UM+	VF
English	652,482	43	25
French	367,732	123	10
German	519,143	93	15
Hebrew	33,177	190	6
Russian	473,481	109	6
Spanish	1,196,245	65	19
Swahili	–	6	24
Turkish	570,420	193	60

Table 1: Number of analyses in UniMorph (UM) and in our extension (UM+), and number of verb forms in the look-up table (VF).

ferent languages is possible by setting language-specific parameters, including the language’s basic word order (OV vs. VO) and a look-up table with the complete inflectional paradigm (i.e. all simple and composite forms, such as

$$\left\{ \left[\begin{smallmatrix} be \\ PRS \end{smallmatrix} \right], \left[\begin{smallmatrix} PTCPP \\ PRS \end{smallmatrix} \right] \right\} \mapsto \left[\begin{smallmatrix} IND \\ PRS \\ PROG \end{smallmatrix} \right], \text{ for English). Table 1}$$

shows the number of verb forms that are included in the look-up table for every language. Since every word may have several morphological analyses, there might also be several clause-level analyses, all of which we let return by the algorithm. The algorithm further identifies the finite verb in each composite analysis, which we return as well. This gives us tuples of the form (a, v) , where a is the analysis of the (possibly composite) verb form and v is the analysis of the finite verb in that form.

In a subsequent step, we determine all possible morphological analyses for every pronoun in the input clause. If a pronoun has an adposition, we override the case of the pronoun with the case assigned by the adposition.⁹ Then, we construct all valid combinations of analyses (a, v, s, N) , where s is the analysis of the subject pronoun and $N \not\equiv s$ are the analyses of the other pronouns. A combination is valid iff s features nominative case and s agrees with v in all nominal features, i.e. number, person, formality and gender. If no valid combination is found, a covert subject pronoun with nominative case and the nominal features of v is introduced for s (this largely affects pro-drop languages like He-

⁹In some languages, the case assigned by an adposition can depend on the inherent case of the pronoun. For example, the German adposition *in* assigns IN+ALL to a dative pronoun and IN+ESS to an accusative pronoun. In these cases, we created case-specific entries for adpositions in our UniMorph extension and our algorithm selects the case for an adposition based on the case of the pronoun.

brew but also imperatives in some other languages). Covert object pronouns are also added to N if the verb form encodes these (this only affects Swahili).

In a last step, we search the clause for question marks and words of negation and add the corresponding features if applicable, yielding combinations of the form (a, v, s, N, c) with $c \sqsubseteq \begin{bmatrix} \text{NEG} \\ Q \end{bmatrix}$.

3.4 Filtering and Pooling

The number of analyses can be quite high but some analyses are more plausible than others. We therefore filter the analyses successively by the following steps:

1. If the clause contains an exclamation mark, only keep imperative analyses.

Motivation: In the shared task’s data, all clauses with an exclamation mark contain an imperative verb and vice versa.

2. For German only: If the clause contains a question mark and the clause is in V2 word order (i.e. it is not syntactically a question), remove the Q feature and only keep quotative analyses.¹⁰

Motivation: In the shared task’s data, all clauses with a question mark contain the Q feature and vice versa, except for German.

3. Only keep analyses where the subject pronoun features nominative case.

Motivation: In the shared task’s data, the subject always features nominative case. Generally, the nominative case marks the subject of a clause in many languages, although there are languages that also have non-nominative subjects (e.g. [Bejar, 2002](#), p. 313).¹¹

4. Only keep analyses where a minimal number of non-subject pronouns features nominative case.

Motivation: In the shared task’s data, non-subject pronouns never feature nominative case. Generally, nominative non-subjects only occur in specific linguistic constructions (e.g. to mark the predicate in copula constructions) or together with a non-nominative subject (cf. [Bejar, 2002](#), p. 313).

5. Only keep analyses with a non-reflexive subject pronoun.

¹⁰What is labeled as ‘quotative’ (QUOT) in the German data is usually called present subjunctive or subjunctive I in the literature and, unlike the labeling in the shared task suggests, not only used in quoted speech.

¹¹Not forgetting ergative languages, in which the subject’s case depends on the (transitivity of the) verb.

Motivation: In the shared task’s data, there are no reflexive subjects. Generally, there do not appear to be any languages with reflexive subjects ([Schachter, 1977](#)).¹²

6. Only keep analyses where every reflexive pronoun agrees with a non-reflexive pronoun. In case of agreement, the non-reflexive pronoun copies missing features from the reflexive pronoun.

Motivation: In the shared task’s data, every reflexive pronoun has an antecedent in the same clause. Generally, reflexive pronouns must have an antecedent in the same sentence (“Binding Principle A” of [Chomsky \(1981\)](#)).¹³

7. Only keep analyses where the pronouns feature a maximal number of different cases.

Motivation: In the shared task’s data, every case appears maximally once per clause. Generally, cases encode grammatical (and in a wider sense also semantic) roles and clauses typically contain every role only once (cf. [Jaworski and Przepiórkowski, 2014](#), p. 84).

8. For French only: Only keep analyses where every past participle agrees with the pronoun determined by the non-trivial French participle agreement rules (cf. [Past Participle Agreement in French, 2017](#)). In case of agreement, the pronoun copies missing features from the participle.

Motivation: In French and other Romance languages, past participles do not always agree with the subject (as it is usually the case) but sometimes with an object (cf. [Kayne, 1989](#)).

9. Only keep analyses where a maximal number of reflexive pronouns agrees with the subject pronoun. In case of agreement, the reflexive pronoun copies missing features from the subject pronoun.

Motivation: In the shared task’s data, reflexive pronouns in ambiguous sentences are sometimes annotated as having subjects and sometimes annotated as having non-subjects as antecedents. Generally, subjects are preferred over non-subjects as antecedents for reflexive pronouns in ambiguous sentences (cf. [White et al., 1997](#), p. 148).

If one of the steps would filter out all analyses, the step is skipped.

¹²English allows statistically rare exceptions (cf. [Song \(2017\)](#), or [Kirk and Kallen \(2006](#), p. 104) for the use of reflexive pronouns as subjects with a focus on Irish English).

¹³Again, English allows statistically rare exceptions (cf. [Kim et al., 2020](#), p. 296).

In a pooling step, redundant features are removed from the analyses, which may result in some of the analyses becoming identical and hence collapsing into one. For example, if there are three analyses for a French input that differ in the analysis of the pronoun *leur*,

$$\begin{bmatrix} \text{PRO} \\ \text{DAT} \\ 3 \\ \text{PL} \\ \text{MASC} \end{bmatrix} \text{ vs. } \begin{bmatrix} \text{PRO} \\ \text{DAT} \\ 3 \\ \text{PL} \\ \text{FEM} \end{bmatrix} \text{ vs. } \begin{bmatrix} \text{PRO} \\ \text{DAT} \\ 3 \\ \text{PL} \\ \text{NEUT} \end{bmatrix}, \text{ then } \begin{bmatrix} \text{PRO} \\ \text{DAT} \\ 3 \\ \text{PL} \end{bmatrix} \text{ be-}$$

comes the reduced analysis of *leur* in each of the analyses. If the three analyses are now completely identical, they are combined into one analysis. On the contrary, if there are two analyses for a German input that differ in the analysis of the pronoun

$$ihm, \begin{bmatrix} \text{PRO} \\ \text{DAT} \\ 3 \\ \text{SG} \\ \text{MASC} \end{bmatrix} \text{ vs. } \begin{bmatrix} \text{PRO} \\ \text{DAT} \\ 3 \\ \text{SG} \\ \text{NEUT} \end{bmatrix}, \text{ the gender feature is not}$$

redundant (since *ihm* cannot be feminine) and can therefore not be removed.

3.5 Ranking

The analyses that are not filtered out are assumed to be correct by the program and can all be output. For the shared task, we rank the analyses according to the following sorting procedures and choose the first one as final result:

1. Choose verb analyses in this order: *lemma of non-auxiliary verb* > *lemma of auxiliary verb*.

Motivation: Analyses with a lemma of an auxiliary verb usually result from errors in the word- or clause-level analysis steps, so we prefer analyses with a lemma of a non-auxiliary verb.

2. Choose pronoun analyses in this order: MASC > FEM > NEUT > no gender.

Motivation: We did not find a general preference for any grammatical gender of ambiguous pronouns in the shared task's data, but we wanted our system to not arbitrarily choose one and this is the order in which many grammars name the genders.

3. Choose pronoun analyses in this order: no class > any class (this only affects Swahili).

Motivation: We experimented with both variants on the training and development set for Swahili and matched the gold analysis in more cases by preferring analyses without class feature over analyses with class feature.

4. Choose pronoun analyses in this order: not LGSPEC3 > LGSPEC3 (this only affects Spanish).

Motivation: We experimented with both variants on the training and development set for Spanish and matched the gold analysis in more cases by preferring analyses without LGSPEC3 feature over analyses with LGSPEC3 feature.

5. Choose pronoun analyses in this order: NOM > ACC > DAT > other case (this effectively prefers analyses where the cases of pronouns appear in this word order).

Motivation: We observed that sentences with ambiguous pronouns always receive cases in this order in the shared task's gold analyses.

6. Choose pronoun analyses in this order: RFLX > not RFLX (except for Spanish, where the sorting is reversed).

Motivation: We experimented with both variants on the training and development set for every language and (for all languages but Spanish) matched the gold analysis in more cases by preferring reflexive readings over non-reflexive readings for ambiguous pronouns.¹⁴

Note that later sorting procedures ignore the previous ones and are therefore more effective.

3.6 Postprocessing

Sometimes, UniMorph contains incorrect lemmas with a trailing *e* for English (e.g. *answere* instead of *answer*). We fix this using NLTK's WordNetLemmatizer¹⁵ and the Python package `pyspellchecker`.¹⁶

The result analysis is then converted to a string in the output format of the shared task.

4 Evaluation and Results

For the evaluation, the gold analysis and the predicted analysis are decomposed into features. For example, the analysis

IND; PST; PFV; NOM(3, PL, MASC); ACC(1, PL, MASC); NEG; Q

is decomposed into the features

$$\Phi = \{\text{IND, PST, PFV, NOM-3, NOM-PL, NOM-MASC, ACC-1, ACC-PL, ACC-MASC, NEG, Q}\}.$$

Given the features for the gold analysis Φ_g and for the predicted analysis Φ_p , the F1 score for one sample is calculated as follows:

$$P = \frac{|\Phi_p \cap \Phi_g| + s_\ell}{|\Phi_p| + w_\ell} \quad R = \frac{|\Phi_p \cap \Phi_g| + s_\ell}{|\Phi_g| + w_\ell}$$

¹⁴An example for an ambiguous pronoun is German *mich*, which can mean 'me' or 'myself' (cf. Hole, 2005, p. 65).

¹⁵<https://www.nltk.org/api/nltk.stem.wordnet.html>

¹⁶<https://pypi.org/project/pyspellchecker/>

Language	Train	Dev	Test
English	.994	.995	.993
French	.973	.974	.977
German	.946	.952	.974
Hebrew (unvoc)	.959	.955	.965
Hebrew (voc)	.966	.970	.955
Russian	.908	.917	.931
Spanish	.931	.920	.943
Swahili	.730	.760	.789
Turkish	.934	.928	.929
Average	.927	.930	.940

Table 2: F1 scores for all languages on the respective training, development and test sets.

$$F_1 = \frac{2 \cdot P \cdot R}{P + R}$$

Hereby, $s_\ell = 3$ if the predicted lemma matches the gold lemma and $s_\ell = 0$ otherwise, and $w_\ell = 3$. While the development and training sets always contain only one gold analysis per sample, the test sets contain multiple gold analyses for samples with an ambiguous sentence. In case of such ambiguous sentences, the predicted analysis is compared to each gold analysis and the highest F1 score is chosen. The F1 for a data set (e.g. for the English test set) is the average F1 over all samples in that set.

The results using our method are shown in Table 2. We achieve F1 scores over 92% on the test sets for each language except Swahili (79%), and an average F1 score of 94% (96% without Swahili). Since our approach is not based on machine learning methods, we observe a relatively stable performance on all data splits (training, development, test) and, in particular, no decrease on the test set.

For completeness, we also show the accuracies in terms of exact matches, i.e. the percentage of predicted analyses that exactly match a gold analysis (including ordering of the features), in Table 3, although we consider this metric to be inadequate for the evaluation of the task since the elements of a feature structure are naturally unordered. Since our rule-based approach cannot learn the ordering of the features from examples, we hard-coded the order of the features in the output string so that it roughly complies with the ordering in the shared task’s data for most languages. After the final system submission, we noticed a mistake in the ordering of the features NEG and Q. Therefore, numbers

Language	Train	Dev	Test
English	.976 (.975)	.977 (.977)	.974
French	.637 (.845)	.676 (.870)	.693
German	.452 (.590)	.465 (.619)	.550
Hebrew (unvoc)	.765 (.765)	.744 (.739)	.827
Hebrew (voc)	.794 (.794)	.807 (.815)	.748
Russian	.459 (.452)	.456 (.472)	.609
Spanish	.492 (.537)	.473 (.553)	.637
Swahili	.041 (.048)	.048 (.066)	.067
Turkish	.841 (.842)	.806 (.808)	.816
Average	.606 (.650)	.606 (.658)	.658

Table 3: Exact matches for all languages on the respective training, development and test sets.

in brackets in Table 3 show exact-match performance after fixing their ordering, while the other numbers are the performances of the system as submitted.¹⁷ The high differences that result from this small change in some languages (e.g. +15% in German) further illustrate the inadequateness of the metric.

5 Discussion

The main advantage of the presented method is probably the performance, although there is naturally some room for improvement. The second major advantage of the method is that it does not require any training data. This makes it a promising option for analyzing every language where manually annotated gold data is not available. No training also means that no training bias can be induced by the data, which arguably makes the method’s performance more stable across text domains. In terms of languages, the algorithm is relatively universally applicable since the underlying mechanisms of inflection and agreement are the same across natural languages. This is also indicated by the performance that is very similar across languages and language families.¹⁸

However, the method is not without shortcomings, all of which are clearly visible in the case of

¹⁷Since gold analyses for the test data have not been released, yet, we cannot re-evaluate our system on the test sets, but we can assume that the performance is nearly the same as on the other splits, or even a bit higher since the test sets can contain more than one gold analysis per sample.

¹⁸The languages in the shared task belong to the following families: Indo-European (English, French, German, Russian, Spanish), Afro-Asiatic (Hebrew), Niger-Congo (Swahili), Turkic (Turkish). [Dönicke \(2021\)](#) also implements the composite-verb analysis for languages from other families.

Swahili. First of all, the method requires a (word-level) morphological analysis and a parser for the language to analyze. We decided to use UniMorph because the output format in the shared task also uses UniMorph features. Dönicke (2021), on the other hand, does not only use the parser but also the morphological analyzer that can be trained by spaCy on a UD treebank.¹⁹ The current version of the UD treebanks includes treebanks for 130 languages and 61 languages are listed as possible future extensions—Swahili being one of them—, and UniMorph currently provides resources for 167 languages. Nonetheless, the current lack of both a treebank and morphological resources for Swahili forced us to implement a workaround resulting in a much lower performance compared to the other languages. Another drawback of our method is that knowledge about the grammar of the language to analyze is required to set-up the language-specific inflection table, the list of auxiliary verbs, the word-order parameter (OV vs. VO), and in the current implementation also a list of words of negation as well as UniMorph-style entries for pronouns and adpositions. Dönicke (2021) already mentions that the study of composite verb forms in a foreign language can be extensive, but it is also prone to errors. It may be a coincidence that the languages with the best performance (English, French, German) are those languages which the author of this paper has the profoundest knowledge of, but it may also be due to the incomplete knowledge about the other languages acquired in the short term. Although the algorithm is designed to be language-independent (with language-specific operations being controlled through the aforementioned parameters), its performance can be sometimes improved by language-specific special rules (e.g. the rules for participle agreement in French), which again can only be implemented by someone who has the according knowledge of the language. Table 4 shows how many of these rules are hard-coded in our implementation. It should be added, however, that some of these rules are only implemented to meet the output format of the shared task and are not related to the morphosyntactic nature of the language. For example, there is no apparent reason why all gold analyses for Swahili have the feature V (verb) while the analyses for the other languages do not; but for the shared task there had to be a special

¹⁹McCarthy et al. (2018) compare UD features and UniMorph features and also provide a tool to convert the former into the latter.

Language	P1	A1	A2	F	R	P2	Σ
English	1	2	1	0	0	1	5
French	1	2	1	1	0	0	5
German	0	1	1	1	0	0	3
Hebrew	2	0	1	1	0	0	4
Russian	0	0	1	0	0	0	1
Spanish	0	1	0	0	1	0	2
Swahili	1	1	2	1	0	1	6
Turkish	0	0	1	0	0	1	2
Σ	5	7	8	4	1	3	28

Table 4: Number of hard-coded language-specific rules in the code. P1: preprocessing, A1: word-level analysis, A2: clause-level analysis, F: filtering and pooling, R: ranking, P2: postprocessing.

rule that adds this feature to every output analysis for Swahili. Probably, the requirement of linguistic knowledge is not that much of a disadvantage, since research teams working on a language usually include some speakers of that language.

6 Conclusion

We presented a method to predict clause-level morphological / morphosyntactic features. The main advantages are its performance (94% F1 on average), that it does not require training data and that it is applicable for multiple languages. The disadvantages are that it requires a preceding word-level morphological analysis, linguistic knowledge about the language to analyze and some time to set-up the method for a new language. While the implementation within the frame of the shared task is not applicable for general use (mainly because of the pre- and postprocessing), interested readers may want to have a look at the implementation from Dönicke (2021).

Acknowledgements

This work is funded by Volkswagen Foundation.

References

- Susana Bejar. 2002. Movement, morphology and learnability. *Syntactic Effects of Morphological Change*, pages 307–325.
- Noam Chomsky. 1981. Lectures on government and binding.
- Narayan Choudhary, Pramod Pandey, and Girish Nath Jha. 2014. A rule based method for the identification of TAM features in a PoS tagged corpus. In

- Human Language Technology Challenges for Computer Science and Linguistics*, pages 178–188, Cham. Springer International Publishing.
- Tillmann Dönicke. 2020. [Clause-level tense, mood, voice and modality tagging for German](#). In *Proceedings of the 19th International Workshop on Treebanks and Linguistic Theories*, pages 1–17, Düsseldorf, Germany. Association for Computational Linguistics.
- Tillmann Dönicke. 2021. [Delexicalised multilingual discourse segmentation for DISRPT 2021 and tense, mood, voice and modality tagging for 11 languages](#). In *Proceedings of the 2nd Shared Task on Discourse Relation Parsing and Treebanking (DISRPT 2021)*, pages 33–45, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Simone Faro and Arianna Pavone. 2015. [Refined tagging of complex verbal phrases for the Italian language](#). In *Proceedings of the Prague Stringology Conference 2015*, pages 132–145, Czech Technical University in Prague, Czech Republic.
- Daniel Hole. 2005. Zur Sprachgeschichte einiger deutscher Pronomina. *Sprachwissenschaft*, 430(1):49–75.
- Wojciech Jaworski and Adam Przepiórkowski. 2014. [Semantic roles in grammar engineering](#). In *Proceedings of the Third Joint Conference on Lexical and Computational Semantics (*SEM 2014)*, pages 81–86, Dublin, Ireland. Association for Computational Linguistics and Dublin City University.
- Richard S. Kayne. 1989. [Facets of Romance past participle agreement](#). In Paola Benincá, editor, *Dialect Variation and the Theory of Grammar*, pages 85–104. De Gruyter Mouton, Berlin, Boston.
- Ji-Hye Kim, Soojin An, and Ahreum Jung. 2020. [Binding conditions of English reflexives and pronouns in the ICE-USA](#). *Lanaguage Research*, 56(3):287–307.
- John M. Kirk and Jeffrey L. Kallen. 2006. Irish Standard English: How Celticised? How Standardised? *The Celtic Englishes IV*, pages 88–113.
- Arya D. McCarthy, Miikka Silfverberg, Ryan Cotterell, Mans Hulden, and David Yarowsky. 2018. [Marrying Universal Dependencies and Universal Morphology](#). In *Proceedings of the Second Workshop on Universal Dependencies (UDW 2018)*, pages 91–101, Brussels, Belgium. Association for Computational Linguistics.
- Skatje Myers and Martha Palmer. 2019. [ClearTAC: Verb tense, aspect, and form classification using neural nets](#). In *Proceedings of the First International Workshop on Designing Meaning Representations*, pages 136–140, Florence, Italy. Association for Computational Linguistics.
- Past Participle Agreement in French. 2017. [Study.com](#). March 11.
- Anita Ramm, Sharid Loáiciga, Annemarie Friedrich, and Alexander Fraser. 2017. [Annotating tense, mood and voice for English, French and German](#). In *Proceedings of ACL 2017, System Demonstrations*, pages 1–6, Vancouver, Canada. Association for Computational Linguistics.
- Paul Schachter. 1977. [Reference-related and role-related properties of subjects](#). In Peter Cole and Jerrold M. Sadock, editors, *Grammatical Relations, Syntax and Semantics*, pages 279 – 306. Brill, Leiden, The Netherlands.
- Sanghoun Song. 2017. A corpus study of unbound reflexive pronouns in English. *영어학*, 17(2):275–305.
- John Sylak-Glassman. 2016. [The composition and use of the universal morphological feature schema \(UniMorph schema\)](#). Johns Hopkins University.
- Lydia White, Joyce Bruhn-Garavito, Takako Kawasaki, Joe Pater, and Philippe Prévost. 1997. The researcher gave the subject a test about himself: Problems of ambiguity and preference in the investigation of reflexive binding. *Language Learning*, 47(1):145–172.
- Eva Žáčková, Luboš Popelínský, and Miloslav Nepil. 2000. [Automatic tagging of compound verb groups in Czech corpora](#). In *Text, Speech and Dialogue*, pages 115–120, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia Ackermann, Noëmi Aepli, Hamid Aghaei, Željko Agić, Amir Ahmadi, Lars Ahrenberg, Chika Kennedy Ajede, Gabrielè Aleksandravičiūtė, Ika Alfina, Avner Algom, Erik Andersen, Lene Antonsen, Katya Aplonova, Angelina Aquino, Carolina Aragon, Glyd Aranes, Maria Jesus Aranzabe, Bilge Nas Arican, Hórunn Arnardóttir, Gashawe Arutie, Jessica Naraiswari Arwidarasti, Masayuki Asahara, Deniz Baran Aslan, Cengiz Asmazoğlu, Luma Ateyah, Furkan Atmaca, Mohammed Attia, Aitziber Atutxa, Liesbeth Augustinus, 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, Kepa Bengoetxea, Yifat Ben Moshe, Gözde Berk, Yevgeni Berzak, Irshad Ahmad Bhat, Riyaz Ahmad Bhat, Erica Biagetti, Eckhard Bick, Agnè Bielinskienė, Kristín Bjarnadóttir, Rogier Blokland, Victoria Bobicev, Loïc Boizou, Emanuel Borges Völker, Carl Börstell, Cristina Bosco, Gosse Bouma, Sam Bowman, Adriane Boyd, Anouck Braggaa, Kristina Brokaitė, Aljoscha Burchardt, Marie Candito, Bernard Caron, Gauthier Caron, Lauren Cassidy, Tatiana Cavalcanti, Gülşen Cebiroğlu Eryiğit, Flavio Massimiliano Cecchini, Giuseppe G. A. Celano, Slavomír Čéplö, Neslihan Cesur, Savas Cetin, Özlem Çetinoğlu, Fabricio Chalub, 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, Daniela Corbetta, Marine Courtin, Mihaela Cristescu, Philemon Daniel, Elizabeth Davidson, 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, Timothy Dozat, Kira Droganova, Puneet Dwivedi, Hanne Eckhoff, Sandra Eiche, Marhaba Eli, Ali Elkahky, Binyam Ephrem, Olga Erina, Tomaž Erjavec, Aline Etienne, Wograine Evelyn, Sidney Facundes, Richárd Farkas, Federica Favero, Jannatul Ferdaousi, Marília Fernanda, Hector Fernandez Alcalde, Jennifer Foster, Cláudia Freitas, Kazunori Fujita, Katarína Gajdošová, Daniel Galbraith, Federica Gamba, Marcos Garcia, Moa Gärdenfors, Sebastian Garza, Fabrício Ferraz Gerardi, Kim Gerdes, Filip Ginter, Gustavo Godoy, Iakes Goenaga, Koldo Gojenola, Memduh Gökirmak, Yoav Goldberg, Xavier Gómez Ginovart, Berta González Saavedra, Bernadeta Griciūtė, Matias Grioni, Loïc Grobol, Normunds Grūzītis, Bruno Guillaume, 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, Jena Hwang, Takumi Ikeda, Anton Karl Ingason, Radu Ion, Elena Irimia, Olájidé Ishola, Kaoru Ito, Siratun Jannat, Tomáš Jelínek, Apoorva Jha, Anders Johannsen, Hildur Jónsdóttir, Fredrik Jørgensen, Markus Juutinen, Sarveswaran K, Hüner Kaşıkara, Andre Kaasen, Nadezhda Kabaeva, Sylvain Kahane, Hiroshi Kanayama, Jenna Kanerva, Neslihan Kara, Ritván Karahóga, Boris Katz, Tolga Kayadelen, Jessica Kenney, Václava Kettnerová, Jesse Kirchner, Elena Klementieva, Elena Klyachko, Arne Köhn, Abdullatif Köksal, Kamil Kopacewicz, Timo Korkiakangas, Mehmet Köse, Natalia Kotsyba, Jolanta Kovalevskaitė, Simon Krek, Parameswari Krishnamurthy, Sandra Kübler, Oğuzhan Kuyrukçu, Aslı Kuzgun, Sookyong Kwak, Veronika Laippala, Lucia Lam, Lorenzo Lambertino, Tatiana Lando, Septina Dian Larasati, Alexei Lavrentiev, John Lee, Phương Lê Hồng, Alessandro Lenci, Saran Lertpradit, Herman Leung, Maria Levina, Cheuk Ying Li, Josie Li, Keying Li, Yuan Li, KyungTae Lim, Bruna Lima Padovani, Krister Lindén, Nikola Ljubešić, Olga Loginova, Stefano Lusito, Andry Luthfi, Mikko Luukko, Olga Lyashevskaya, Teresa Lynn, Vivien Macketanz, Menel Mahamdi, Jean Maillard, Aibek Makazhanov, Michael Mandl, Christopher Manning, Ruli Manurung, Büşra Marşan, Cătălina Mărânduc, David Mareček, Katrin Marheinecke, Stella Markantonatou, Héctor Martínez Alonso, Lorena Martín Rodríguez, André Martins, Jan Mašek, Hiroshi Matsuda, Yuji Matsumoto, Alessandro Mazzei, Ryan McDonald, Sarah McGuinness, Gustavo Mendonça, Tatiana Merzhevich, Niko Miekka, Karina Mischenkova, Margarita Misirpashayeva, Anna Missilä, Cătălin Mititelu, Maria Mitrofan, Yusuke Miyao, AmirHos-

sein Mojiri Foroushani, Judit Molnár, Amirsaeid Moloodi, Simonetta Montemagni, Amir More, Laura Moreno Romero, Giovanni Moretti, Keiko Sophie Mori, 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 Nikaïdo, Vitaly Nikolaev, Rattima Nitisaroj, Alireza Nourian, Hanna Nurmi, Stina Ojala, Atul Kr. Ojha, Adédayo 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, Hyunji Hayley Park, Niko Partanen, Elena Pascual, Marco Passarotti, Agnieszka Patejuk, Guilherme Paulino-Passos, Giulia Pedonese, Angelika Peljak-Lapińska, Siyao Peng, Cene-Augusto Perez, Natalia Perkova, Guy Perrier, Slav Petrov, Daria Petrova, Andrea Peverelli, Jason Phelan, Jussi Piitulainen, Tommi A Pirinen, Emily Pitler, Barbara Plank, Thierry Poibeau, Larisa Ponomareva, Martin Popel, Lauma Pretkalniņa, Sophie Prévost, Prokopis Prokopidis, Adam Przepiórkowski, Tiina Puolakainen, Sampo Pyysalo, Peng Qi, Andriela Rääbis, Alexandre Rademaker, Mizanur Rahoman, Taraka Rama, Loganathan Ramasamy, Carlos Ramisch, Fam Rashel, Mohammad Sadegh Rasooli, Vinit Ravishankar, Livy Real, Petru Rebeja, Siva Reddy, Mathilde Regnault, Georg Rehm, Ivan Riabov, Michael Rießler, Erika Rimkutė, Larissa Rinaldi, Laura Rituma, Putri Rizqiyah, Luisa Rocha, Eiríkur Rögnvaldsson, Mykhailo Romanenko, Rudolf Rosa, Valentin Roşca, Davide Rovati, Ben Rozonoyer, Olga Rudina, Jack Rueter, Kristján Rúnarsson, Shoval Sadde, Pegah Safari, Benoît Sagot, Aleksí Sahala, Shadi Saleh, Alessio Salomoni, Tanja Samardžić, Stephanie Samson, Manuela Sanguinetti, Ezgi Sanyar, Dage Särg, Baiba Saulīte, Yanin Sawanakunanon, Shefali Saxena, Kevin Scannell, Salvatore Scarlata, Nathan Schneider, Sebastian Schuster, Lane Schwartz, Djameé Seddah, Wolfgang Seeker, Mojgan Seraji, Syeda Shahzadi, Mo Shen, Atsuko Shimada, Hiroyuki Shirasu, Yana Shishkina, Muh Shohibussirri, Dmitry Sichinava, Janine Siewert, Einar Freyr Sigurðsson, Aline Silveira, Natalia Silveira, Maria Simi, Radu Simionescu, Katalin Simkó, Mária Šimková, Kiril Simov, Maria Skachedubova, Aaron Smith, Isabela Soares-Bastos, Shafi Sourov, Carolyn Spadine, Rachele Sprugnoli, Vivian Stamou, Steinhórfur Steingrímsson, Antonio Stella, Milan Straka, Emmett Strickland, Jana Strnadová, Alane Suhr, Yogi Lesmana Sulestio, Umot Sulubacac, Shingo Suzuki, Daniel Swanson, Zsolt Szántó, Chihiro Taguchi, Dima Taji, Yuta Takahashi, 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, Sumire Uematsu, Roman Untilov, Zdeňka Urešová,

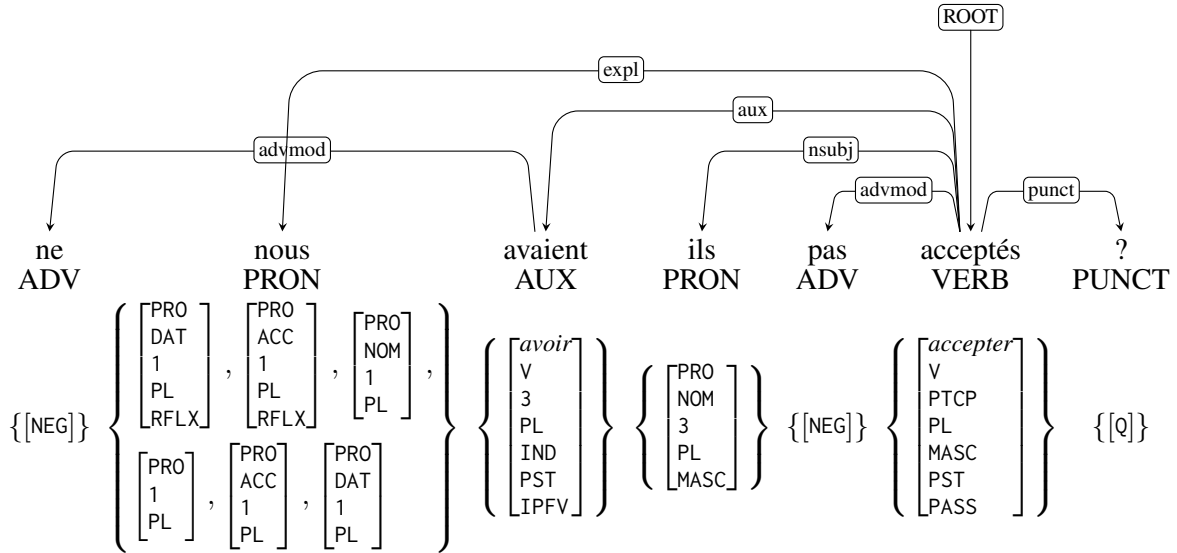
Larraitz Uria, Hans Uszkoreit, Andrius Utkas, Elena Vagnoni, Sowmya Vajjala, Rob van der Goot, Martine Vanhove, Daniel van Niekerk, Gertjan van Noord, Viktor Varga, Uliana Vedenina, Eric Villemonte de la Clergerie, Veronika Vincze, Natalia Vlasova, Aya Wakasa, Joel C. Wallenberg, Lars Wallin, Abigail Walsh, Jing Xian Wang, Jonathan North Washington, Maximilian Wendt, Paul Widmer, Shira Wigderson, Sri Hartati Wijono, Seyi Williams, Mats Wirén, Christian Wittern, Tsegay Woldemariam, Tak-sum Wong, Alina Wróblewska, Mary Yako, Kayo Yamashita, Naoki Yamazaki, Chunxiao Yan, Koichi Yasuoka, Marat M. Yavrumyan, Arife Betül Yenice, Olcay Taner Yıldız, Zhuoran Yu, Arlisa Yuliawati, Zdeněk Žabokrtský, Shorouq Zahra, Amir Zeldes, He Zhou, Hanzhi Zhu, Anna Zhuravleva, and Rayan Ziane. 2022. [Universal dependencies 2.10](#). LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

A French Example 1

Input: ne nous avaient-ils pas acceptés?

Gold Output: accepter IND;PST;PFV;NOM(3,PL,MASC);ACC(1,PL,MASC);NEG;Q

After preprocessing and word-level analysis:



After ...	<i>a</i>	<i>v</i>	<i>s</i>	<i>N</i>	<i>c</i>
add <i>a, v</i>	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]			
add <i>s</i>	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]		
add <i>N</i>	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL RFLX] }	
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL RFLX] }	
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO NOM 1 PL] }	
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO 1 PL] }	
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL] }	
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL] }	
add <i>c</i>	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL RFLX] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL RFLX] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO NOM 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL] }	[NEG Q]
filter 4	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL RFLX] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL RFLX] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL] }	[NEG Q]
filter 6	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL] }	[NEG Q]
filter 7	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL] }	[NEG Q]
ranking	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO ACC 1 PL] }	[NEG Q]
	[accepter IND PST PFV]	[avoir V 3 PL IND PST IPFV]	[PRO NOM 3 PL MASC]	{ [PRO DAT 1 PL] }	[NEG Q]

Pred. Output: accepter IND;PST;PFV;NOM(3,PL,MASC);ACC(1,PL,MASC);NEG;Q

