

Building an Endangered Language Resource in the Classroom: Universal Dependencies for Kakataibo

Roberto Zariquiey^ρ, Claudia Alvarado^ρ, Ximena Echevarria^ρ, Luisa Gomez^ρ,
Rosa Gonzales^ρ, Mariana Illescas^ρ, Sabina Oporto^ρ,
Frederic Blum^β, Arturo Oncevay^ε, Javier Vera^υ

^ρPontificia Universidad Católica del Perú, Peru

^βHumboldt-Universität zu Berlin and Leibniz-Zentrum Allgemeine Sprachwissenschaft

^εUniversity of Edinburgh, Scotland

^υPontificia Universidad Católica de Valparaíso, Chile

{rzariquiey, claudia.alvarado, ximena.echevarria, luisa.gomez, a20175617, m.illescasb, sabina.oporto}@pucp.edu.pe
frederic.blum@hu-berlin.de, a.oncevay@ed.ac.uk, javier.vera@pucv.cl

Abstract

In this paper, we launch a new Universal Dependencies treebank for an endangered language from Amazonia: Kakataibo, a Panoan language spoken in Peru. We first discuss the collaborative methodology implemented, which proved effective to create a treebank in the context of a Computational Linguistic course for undergraduates. Then, we describe the general details of the treebank and the language-specific considerations implemented for the proposed annotation. We finally conduct some experiments on part-of-speech tagging and syntactic dependency parsing. We focus on monolingual and transfer learning settings, where we study the impact of a Shipibo-Konibo treebank, another Panoan language resource.

Keywords: Universal Dependencies, Treebank, Collaborative Methodology, Kakataibo, Endangered Languages, Panoan, Amazonia, Peru

1. Introduction

Kakataibo is a language that belongs to the Panoan family spoken by around 3,000 native speakers in the Amazonian region of Peru. This paper describes the methodology implemented in the context of a regular undergraduate Computational Linguistics course to create a UD treebank for this language, as a strategy to develop significant learning settings and at the same time contribute to the computerization of this endangered language of Peruvian Amazon. The Kakataibo UD treebank would enhance the future development of an NLP toolkit for this language, since it is the main requirement to train a dependency parser. By taking advantage of the preexistence of a UD treebank for another Panoan language, Shipibo-Konibo (Vasquez et al., 2018), in this paper, we conduct some experiments in both monolingual and transfer learning settings.

The paper is organized as follows. First, §2 presents some background information on the Kakataibo language. Then, §3 describes the methodology implemented in the classroom. §4 introduces the Kakataibo UD treebank. §5 presents the experimentation conducted on part-of-speech tagging and syntactic dependency parsing in both monolingual and transfer learning settings (using the Shipibo-Konibo as a baseline). Finally, §6 summarises the conclusions of this paper.

2. The Kakataibo Language

Kakataibo (cbr) is a Panoan language spoken by approximately 3,000 people in the Peruvian departments of Huánuco and Ucayali. The Kakataibo people live in various communities along the Aguaytía, San Alejandro, Shamboyacu, Sungaroyacu and Pisqui

Rivers, where the language remains vital despite different degrees of contact between Kakataibo people and non-indigenous populations. Zariquiey (2011) distinguishes four living Kakataibo varieties: the Lower Aguaytía/Shamboyacu, Upper Aguaytía, Sungaroyacu and San Alejandro dialects. Nokamán, a variety named and minimally documented by Tessmann (1930), was a fifth variety, now extinct (Zariquiey, 2013). Among the living varieties, the most divergent is the San Alejandro one, with the Upper Aguaytía and Sungaroyacu varieties being highly similar to each other, and (to a lesser degree) to the Lower Aguaytía variety, which is the one represented in the treebank featured in this paper. The sentences belong to the first author’s database and were gathered in the field between 2007 and 2011. The Lower Aguaytía dialect, studied in this paper, exhibits the phonological inventory given in Tables 1 and 2.

	labial	alveolar	palatal	retroflex	velar	glottal
stop	p	t			k k ^w	ʔ
affricate		ts	tʃ			
fricative		s	ʃ	ʂ		
nasal	m	n	ɲ			
flap		r				
glide	β					

Table 1: Kakataibo consonant inventory

In terms of its typological profile, Kakataibo is an agglutinative language with synthetic verbal morphology (i.e., we find single verbal words composed of several morphemes). The language is both head and depen-

	front	central	back
high	i	i	u
mid	e		o
low		a	

Table 2: Kakataibo vowel inventory

dent marking, with a complex system of grammatical relations that combines ergative and tripartite (i.e., intransitive subject vs. transitive subject vs. transitive object) alignments in case marking with an accusative alignment in subject cross-referencing on both verbs and a closed set of second position clitics. Clausal constituent order is pragmatically determined, but there is a tendency towards verb-final clauses. Word order in the noun phrase is not fixed and most nominal modifiers can appear either before or after the nominal head. The language also exhibits a rich switch-reference system and pervasive use of nominalizations in discourse. Kakataibo verbs are inherently transitive or intransitive (with almost no labile verbs). The transitivity of the verb, which can only be altered by the use of valence changing markers, is encoded in various parts of the clause. Kakataibo exhibits a complex tense system with several past tense markers. There is a large set of verb morphemes and enclitics of different sorts encoding evidentiality, modality, mood, and a highly unusual typological category called speech genre in (Zariquiey, 2018). For a full reference grammar with an abundant discussion of each of these features, the readers are referred to Zariquiey (2018).

3. Methodology in the Classroom

Most NLP courses and textbooks focus only on algorithms and mathematical techniques, with much less attention to data collection and processing, among other more practical problems associated with the implementation of NLP projects (Vajjala, 2021). This is a fundamental issue when the audience of the course does not come from Computer Sciences and/or Engineering, as is often the case in the growing body of language technology techniques applied in humanities research around the world (Hinrichs et al., 2019; Hiippala, 2021). In this scenario, determining how much mathematics/programming/linguistics should be included in an NLP-focused course is neither a trivial nor an easy question (Vajjala, 2021). Taking into consideration the prototypical linguistics students' background, the most suitable approach would keep mathematics to the bare minimum, focusing on the basic programming elements (lists, dictionaries, if-else) that may provide the students with the necessary skills to deal with linguistic corpora, as well as accomplish basic NLP-related tasks (Vajjala, 2021). To produce a significant learning experience, such programming elements must be introduced in the solutions of concrete analytic problems, such as creating and processing real linguistic data. This perspective was taken in the implementa-

tion of a Computational Linguistics course for undergraduates, taught by two of the co-authors of this paper at the Humanities Department of the Pontificia Universidad Católica del Perú during the second semester of 2021 (duration: sixteen weeks, four teaching hours per week). The Kakataibo treebank launched in this contribution was one of the research outcomes of this course.

3.1. Background: Course goals and students

The course was designed for advanced undergraduate students with extensive knowledge in Linguistics (e.g., phonetics, phonology, morphology and syntax), as well as with training on the grammar of a sample Peruvian languages. It had no prerequisites as it is an optional course. Six linguistics undergraduate students, co-authors of this paper, were enrolled in the course. The class size was small due to the novelty of the course and the reduced Linguistics alumni at the university. The students had a minimal technology background, including some experience in web development, though most of the class had neither a deep knowledge of NLP resources nor a background in programming. Since all were undergraduate students, the materials provided were specially designed to cover topics that address the possible uses of specific NLP resources for conducting linguistic work on Amazonian languages. In line with Bender (2007), the pedagogical goals of the course were:

- to give students an introduction to Python programming;
- to provide hand-on experience in the analysis of linguistic data; and
- to explore the consequences of NLP systems and Computational Linguistic analysis, especially in minority languages.

3.2. Course content

This one-semester course was divided into three parts. In the first unit (five weeks, twenty teaching hours), the students participated on hand-on classes to increasing complex Python programming tasks. With this knowledge, in the second unit (six weeks, twenty-four teaching hours), the students analysed linguistic data of several sources: text documents (with the main goal to learn the basics of token/type frequency counts), *UniMorph* annotations (Sylak-Glassman et al., 2015) (with the main purpose to learn about Python dictionaries) and structured data from typological databases (like SAILS (Muysken et al., 2016)) (in order to practice with *csv* files and dictionaries). The third unit of the course (five weeks, twenty teaching hours) was focused on building a new Universal Dependencies (de Marneffe et al., 2021) treebank of a Peruvian minority language: Kakataibo.

3.3. Collaborative Methodology for the Development of the Language Resource

To accomplish the annotation task proposed in the third unit of the course, we developed a collaborative methodology and an annotation ecosystem that proved satisfactory to complete our project in a short time and with high-quality outcomes. The principle behind our methodology was to promote a bridge between the growing body of descriptive work on Peruvian languages and NLP initiatives. Although there is still a lot to be done in this respect, during the last 20 years a significant number of detailed typology-oriented reference grammars of Peruvian (and South American) languages have been published (Zariquiey et al., 2019). Such grammars often feature hundreds of fully annotated and parsed example sentences, accompanied by an analytical discussion that provides a sound basis for their interpretation. Therefore, developing treebanks based on these examples does not require advanced knowledge of the grammar of the language, but just a proper understanding of the examples to transform typology-oriented annotations into UD annotations.

We chose Kakataibo, a Panoan language spoken in Peru, since there was an available full reference grammar of the language (Zariquiey, 2018), written by the first author of this paper, who was also engaged in the regular teaching of the course. The grammar contains 1012 linguistic examples, three fully annotated complete narratives and a small dictionary. We expect to incorporate a larger set of fully glossed sentences into an extended version of the Kakataibo UD treebank launched here.

During the third unit of the course, each class consisted of discussions about UD annotations of some illustrative sentences. This served as guidance for the students regarding the parameters of the annotation, to ensure consistency across the annotated sentences. In addition, each student had a small group of test sentences that were reviewed by the expert on Kakataibo (the first author of this paper) to verify the quality of the work. The annotation quality was understood from two points of view: 1° the UD general guidelines; and 2° the particularities of the Kakataibo grammar. Then, each student randomly chooses a set of sentences to annotate and when the complete list of sentences was defined, they were manually annotated by the first author and the students. After that, each student annotated around 15 sentences and problematic cases would be discussed and resolved through Zoom meetings. Figure 1 features a caption of a manually annotated sentence, the yellow line was added by a student, and indicates that she identified a missing dependency. Once all the corpus was gathered in GitHub, a couple of students, with a supervisor professor, would oversee and correct any possible typing errors and make sure everything was in order.

Draft annotations were implemented in the annotation tool *UD Annotatrix* (Tyers et al., 2018), and each *CoNLL-U* file was also carefully revised to fix bugs and

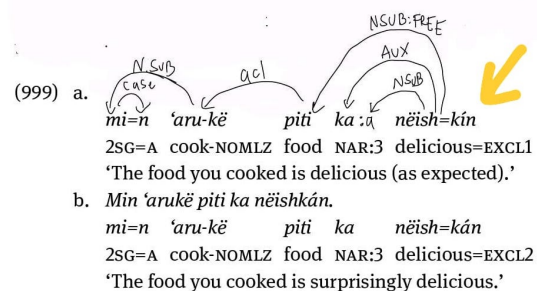


Figure 1: Annotation of a Kakataibo sentence in Zariquiey's grammar

other technical issues. We ended up with 130 fully annotated sentences from Zariquiey (2018), which constitute the first Kakataibo UD Treebank. Using the *conllu* Python library (Stenström, 2016), *CoNLL-U* formatted strings were transformed into Python dictionaries to conduct further in-class programming experiments. This process continued for four more weeks after the end of the academic semester, so the total time frame used in the creation of the treebank was nine weeks.

3.4. Discussion

An important lesson of this process was the enormous value that grammar's examples have for producing NLP resources. During the treebank creation, it became clear that implementing UD annotation based on typological categories like the ones used in reference grammars is a fairly straightforward process, due to the salient coincidences between linguistic typology and Universal Dependencies (Croft et al., 2017). We strongly believe that the implemented methodology can be efficiently replicated in the future collaborative creation of treebanks based on grammars' examples in the frame of NLP, programming or computational linguistics courses, or workshops for Linguistics students. Concerning this methodology's scalability, we propose to include inter-annotator agreement in the following way: bigger classrooms could be divided into groups, each of which would be assigned a set of sentences to be annotated. This dynamic would allow a finer annotation for each sentence. The selected sentences would have to be first analysed by each group member and later on discussed within the group. The final product would show a consensual annotation for that sentence, providing us with a more rigorous filter. In addition to this, we strongly believe that it would be important to maintain the establishment of parameters as a mandatory part of the course. It would help to ensure consistency across the corpus and to reduce the inquiries made to the language expert. Although we benefited from the first author's direct field experience on the language under study, this is not a requisite to successfully implementing this methodology. High-quality refer-

ence grammars are often self-explanatory if one has training in linguistic typology.

We believe that increasing the involvement of large groups of undergraduate students would not be an issue as they are generally keen to contribute and get engaged in research projects, particularly those which may contribute to the development and revitalization of endangered languages. The idea of submitting our joint work to an international academic conference was also highly appealing to the students and reinforced their commitment to the project.

4. The Kakataibo Treebank

4.1. Part-of-Speech

Universal Dependencies features a tagset of 17 Part of speech (POS) tags, mainly based on the Google universal part-of-speech tags (Petrov et al., 2012), and 15 of them were used to elaborate the Kakataibo treebank. The POS tags and frequencies in the treebank are shown in Table 6. The POS tags *X* and *SYM* were not included in this version of the treebank but were relevant for Shipibo-Konibo, another Panoan language (Vasquez et al., 2018; Pereira-Noriega et al., 2017). In the Shipibo-Konibo treebank, *X* was used for onomatopoeias, which is also a part of speech in Kakataibo, not yet attested in the annotated treebank. A future version of this treebank may incorporate those tags in its POS tags repertoire.

Following the Shipibo-Konibo treebank (Vasquez et al., 2018), Kakataibo enclitics are treated here as an independent closed POS, labelled as *PART*, which is one of the POS tags included in the UD POS tagset. According to Zariquiey (2018), there are three types of enclitics in Kakataibo: noun phrase (NP) enclitics, second position enclitics, and adverbial enclitics. Noun phrase enclitics appear at the right edge of NPs. Since some NP modifiers are post-nuclear, NP-enclitics do not necessarily attach to nouns, but also adjectives, nominalizations, determinants, and numerals, among other NP modifiers. Second position enclitics are positionally fixed: they always appear after the first constituent of the sentence independently of its syntactic nature. Adverbial enclitics are non-positional and they may appear attached to any constituent of the clause independently of its position. All enclitics in Kakataibo are phrase-level modifiers.

As can be seen in Table 6, *PART* is largely the most frequent POS in the Kakataibo treebank (freq = 0.40), this is mainly because second-position enclitics in Kakataibo are even more obligatory than verbs (there are verbless copula clauses in Kakataibo, but each independent sentence in the language must carry second position enclitics encoding speech genre and subject indexation).

4.2. Universal Dependency relations

UD defines a set of 37 dependency relations, mainly based on the Universal Stanford Dependencies (de

Marneffe et al., 2014), and 27 of these have been used for the annotation of the Kakataibo treebank, as specified in Table 8. Although UD aims to provide a universal set of syntactic dependencies as a strategy to facilitate consistent annotation across languages and cross-linguistic comparisons (Nivre et al., 2020), it also provides alternatives to code language-specific categories, using "subtype" relation labels. In the case of Kakataibo, it is required to acknowledge the distinction between auxiliary verbs and auxiliary particles. We coded this distinction employing various subtypes of the dependency *aux*, as discussed in 4.2.1. On the other hand, in addition to verbal morphology, subject encoding is accomplished through two independent constituents in the Kakataibo sentence: as part of the second-position enclitic complex and utilizing an independent noun phrase (only the former is obligatory). Based on this, we propose two subtypes for the *nsubj* dependency: *bound* and *free*.

4.2.1. Subtypes of *aux*

Second position enclitics clearly satisfy the definition of auxiliary provided in the UD protocol, but due to their different POS and their syntactic particularities, they need to be distinguished from verb auxiliaries in periphrastic verbal constructions. We implement this distinction by means of using *aux* without further subtype specification for auxiliary verbs and *aux:subtype* for the various types of categories encoded by means of second-position enclitics. Thus, we have the following *aux* subtypes: *aux:sgen* (related to the obligatory category of speech genre coded in each Kakataibo sentence, which encodes a pragmatic distinction between narrative and conversational genres); *aux:ev* (used for the reportative evidential); *aux:int* (used for the interrogative enclitic); and *aux:dub* (used for the dubitative enclitic). There are more second-position enclitics in Kakataibo, but they have not appear in our treebank yet. A detailed discussion on the syntax and semantic of these enclitics is offered somewhere else (Zariquiey, 2018). The Shipibo-Konibo treebank follows a similar approach, by including an *aux* subtype (*aux:val*) (Vasquez et al., 2018), which is more or less equivalent to *aux:ev* in Kakataibo (here we used *aux:ev*, since "evidential" is more widely used than "validator" in the contemporary typological literature).

One interesting point about the distinction between the different types of *aux* dependencies has to do with the direction, auxiliary verbs appear to the right of the root, whereas second-position enclitics appear to the left of the root. A Kakataibo sentence featuring the dependencies *aux:sgen*, *aux:ev* and *aux* is presented in Figure 2.

4.2.2. Subtypes of *nsubj*

Noun phrases overtly encoding the subject of a clause are not obligatory in Kakataibo, but there is obligatory subject indexation in the verb and in the second position enclitics. While subject indexation in the verb can be considered as part of the verbal morphology, the de-

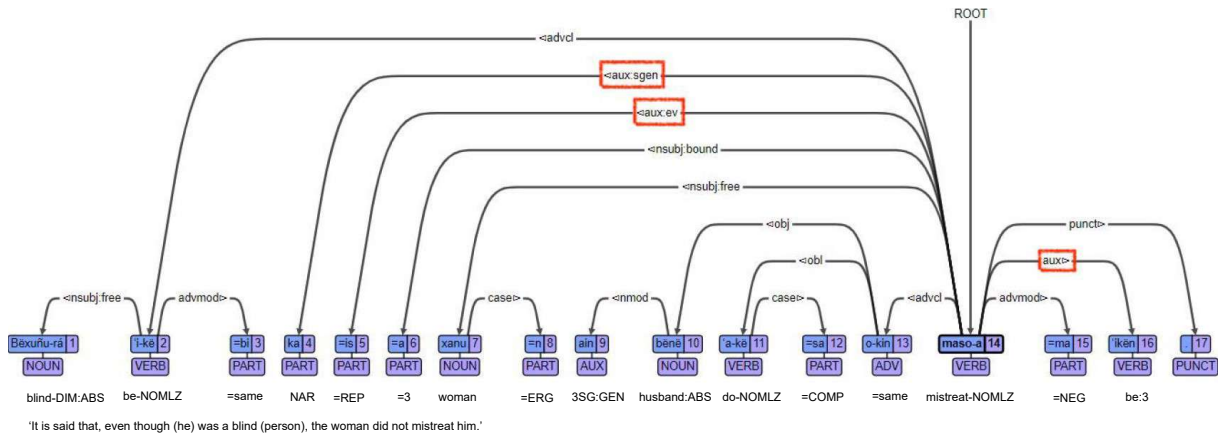


Figure 2: A Kakataibo sentence featuring the dependencies *aux:sgen*, *aux:ev* and *aux*. The featured sentence is *Bəxuñurá 'ikēbi kaisa xanun ain bënë 'akěsa okin masoama 'ikēn* 'It is said that, even though her husband was blind, the woman did not mistreat him.'

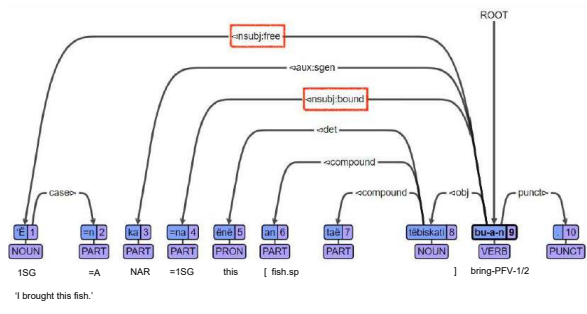


Figure 3: A Kakataibo sentence featuring the dependencies *nsbj:free* and *nsbj:bound*. The featured sentence is *'Ēn kana ənə an taě tēbiskati buan* 'I took this *an taě tēbiskati* (fish species)'

cision of treating enclitic as independent POS (PART) lead to an annotation in which the *nsbj* dependency goes from the root to both the second-position enclitic encoding subject indexation and to the head of the subject NP (if overtly expressed). To clarify that there are fundamental differences regarding the syntactic nature of the two ways of encoding subjects in Kakataibo (e.g., one is obligatory and the other is optional), we decided to identify two different subtypes of *nsbj*, that is *nsbj:free* and *nsbj:bound*. A Kakataibo sentence featuring the dependencies *nsbj:free* and *nsbj:bound* is presented in Figure 3.

5. Experimentation

Once we finished the current version of the Kakataibo treebank, we conducted experiments in POS tagging and dependency parsing in different monolingual and transfer settings for both the Shipibo-Konibo and the Kakataibo treebank. For this reason, we first compare the frequency of each tag in both treebanks (see Tables 6 and 8 in the Appendix) and the utterances length (see 4). Regarding the POS tags, we observe

that both treebanks contain similar proportions, although the Kakataibo treebank presents PART more frequently, and contains less punctuation marks. Besides, the utterances in the Kakataibo treebank tend to be longer than those in the Shipibo-Konibo treebank, with a mean length of 9.52 (±3.22) for Kakataibo, and 7.08 (±2.92) for Shipibo. Each language has one utterance that did not fit the limits of the graph, with a respective token length of 28 (Kakataibo) and 23 (Shipibo). These differences pose a limitation for the following transfer learning experiments.

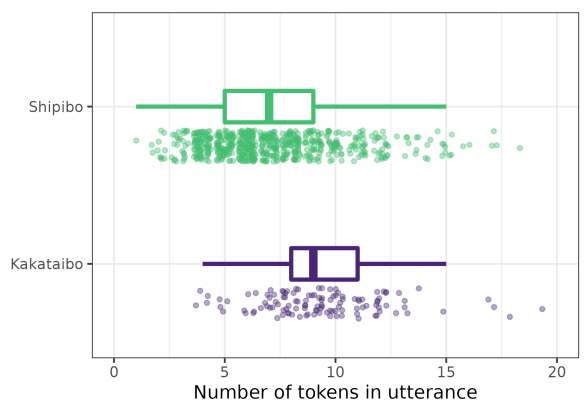


Figure 4: Utterance length in the Panoan treebanks

5.1. POS tagging

For the POS-tagging experiment, Shipibo-Konibo was split into an 80/10/10 training/dev/test set. In order to have sufficient utterances for the dev- and the test-set, Kakataibo was split into partitions of 60/20/20. This ensured that the evaluation was done on more than 20 utterances, even though the training set now only consisted of 72 sentences. While this lead to slightly lower accuracy and f1-scores, it significantly improved the stability of the results and should be considered more

	train	fine-tune	cbr accuracy	cbr f1	shp accuracy	shp f1
1	cbr		84.5±1	46.9±2.4	35.3±1.2	10.6±0.2
2	shp		61.0±1.4	21.1±2.1	93.4±0.7	84.6±1.6
3	shp	cbr	76.9±2.8	34.1±4	93.2±0.2	85.6±0.5

Table 3: Results of the POS tagging experiment (cbr = Kakataibo, shp = Shipibo-Konibo)

POS	precision	recall	f1-score	n
PART	0.9681	1.0000	0.9838	91
NOUN	0.7317	0.8108	0.7692	37
VERB	0.7941	0.8710	0.8308	31
PUNCT	1.0000	0.9200	0.9583	25
PRON	0.8889	0.7273	0.8000	22
DET	0.2500	0.2500	0.2500	4
ADJ	0.0000	0.0000	0.0000	3
ADV	0.0000	0.0000	0.0000	4
AUX	0.6667	0.6667	0.6667	3
PROPN	1.0000	0.6667	0.8000	3
NUM	0.0000	0.0000	0.0000	2
CCONJ	0.0000	0.0000	0.0000	1
micro avg	0.8496	0.8496	0.8496	226
macro avg	0.5250	0.4927	0.5049	226

Table 4: f1-scores for POS-tagging

reliable¹. The distribution of POS-tags across the sets for Kakataibo are given in Table 7. We tested three different experiment settings for POS tagging.

1. Monolingual training with the Kakataibo treebank.
2. Monolingual training with the Shipibo-Konibo treebank and zero-shot transfer to Kakataibo.
3. Monolingual training with the Shipibo-Konibo treebank and fine-tuning for Kakataibo.

As model architecture, we used a BiLSTM-CRF dependency parser implemented in *flair* (Akbik et al., 2019).² The contextual string embeddings (Akbik et al., 2018) were based on the JW300-corpus (Agić and Vulić, 2019), which was specifically trained on typologically diverse low-resource languages, and showed significantly better results than transformer-based embeddings for our experiments. The overall results can be found in Table 3 and the f1-scores per POS tag is given in Table 4.

Despite the small training set, the accuracy in both the fine-tuning setting as well as the monolingual Kakataibo training showed good results. Especially the monolingual training was very successful, with an accuracy over 84% on average. The low f1-score partially has its origin in the fact that not all tags are equally

¹Another option is to perform a leave-one-out analysis, but given our limited resources, we stick to the partition split.

²<https://github.com/flairNLP/flair>, Version 0.10, MIT License

present in the three different data partitions used for training and testing.

It is also noteworthy that even the second setting, a fully lexicalized zero-transfer of the POS-tagger, achieved an accuracy over 60%. For a semi-automated workflow of annotating a new treebank, it could prove worthwhile to train a zero-shot model on a closely related language and then correct at least 100 utterances manually. From this point, it would then be recommended to start building a monolingual tagger for further annotations.

5.2. Dependency parsing

For dependency parsing, we use a deep bi-affine neural dependency parser (Dozat and Manning, 2017) that is implemented in *supar* (Zhang et al., 2020).³ The following settings were used in the experiment:

1. Delexicalized transfer from Kazakh to (lexical) Kakataibo.
2. Delexicalized transfer from Shipibo-Konibo to (lexical) Kakataibo.
3. Delexicalized transfer from Kazakh to delexicalized Kakataibo.
4. Delexicalized transfer from Shipibo-Konibo to delexicalized Kakataibo.
5. Monolingual training for Shipibo-Konibo and zero-shot transfer to Kakataibo.
6. Monolingual training of the reduced Kakataibo set (60/20/20).
7. Monolingual training with the full Kakataibo set (80/10/10).

Experiment 1 and 3 are motivated through previous work on Shipibo-Konibo, showed that the typological proximity of Kazakh (Tyers and Washington, 2015) provides good results for delexicalized transfer of a dependency parser to that language (Vasquez et al., 2018). The goal of this setting is to confirm these results for a second Panoan language, and see whether the results are stable, or only a mix of typological proximity and shared random patterns present in both datasets.

³<https://github.com/yzhangcs/parser>, Version 1.01, MIT License

	model	train	UAS cbr	LAS cbr	UAS shp	LAS shp
1	delex to lex	ktb	32.4	17.7	26.5	10.4
2	delex to lex	shp	9.5	3.1		
3	delex to delex	ktb	51.3	28.2	64.9	40.5
4	delex to delex	shp	60.4	39.3		
5	mono	shp	20.3±4.7	2±0.5	87.7±0.9	80.1±1
6	mono	cbr	73.1±3.3	60.4±3		
7	mono full	cbr	77.4±3.7	67.4±1.6		

Table 5: Results of the Dependency parsing experiment

The Kazakh data is taken from the current UD release (Makazhanov et al., 2015; Tyers and Washington, 2015).⁴

We extend the experiment of delexicalized transfer by adding experiment 2 and 4, and test the delexicalized transfer of two closely-related languages, albeit one, Kakataibo, having less annotated data available (Zeman and Resnik, 2008). The splits for Experiment 6 and 7 are presented in Table 9. The results of those experiments are presented in Table 5. The Unlabelled Attachment Score (UAS) refers to the correct assignment of a head for any element, without taking the UPOS into account. The Labelled Attachment Score (LAS) calculates the score for the combination of dependency relation and UPOS. The UAS for Kakataibo in the different experimental settings is presented in Figure 5.

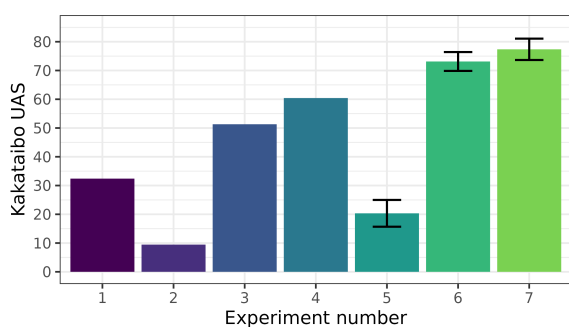


Figure 5: UAS results for Kakataibo

With respect to experiment 2, it is important to note that the results are not an error. Using the Kazakh embeddings together with the Shipibo-Konibo data yielded results that regularly surpassed 20% during training, but then collapsed back into <10%-results for the test data as well. It seems as if the Shipibo-Konibo dependency-relations data is actually surprisingly unfit for transfer to Kakataibo. This holds for the delexical-to-lexical setting, as well as for the lexical zero-shot transfer model. On the other hand, the fully delexicalized results was better than the corresponding Kazakh model. We hypothesise that the origin of this problem lies in the different utterance lengths of the treebanks described earlier.

⁴https://github.com/UniversalDependencies/UD_Kazakh-KTB/tree/master, Version 2.9, CC BY-SA license.

The findings that delexicalized transfer from Kazakh to Panoan languages can temptatively be confirmed, but not to the same extent as in the previous findings. This suggests that even though the typological proximity has a strong effect on transferability of delexicalized dependency parsers, it may just have been a coincidence that it was Kazakh out of all available typologically similar languages that showed the best results in previous work on Shipibo-Konibo. The exact factors for leveraging typological similarity for sharing NLP resources remain unclear, but further studies on this topic are pressing (Bender, 2009).

Comparing experiment 6 and 7, we were again surprised by the small difference between the reduced dataset and the full treebank. This shows that even a small training set of around 100 utterances (74 train + 24 test) shows results that can be implemented in annotation workflows or further experimentation settings. Delexicalized transfer from a closely related language can boost initial annotation steps, once POS tags are already available.

6. Conclusions

We introduced here a new NLP resource for a Peruvian endangered language: a Kakataibo Universal Dependencies treebank. The resource comprises 130 annotated sentences, and features 15 POS tags and 27 dependency relations. Two of the Kakataibo dependency relations feature further subtypes: *aux* (*aux*, *aux:sgen*, *aux:ev*, *aux:int* and *aux:dub*) and *nsubj* (*nsubj:free* and *nsubj:bound*). This treebank is the first one produced for Kakataibo, but the second one for a Panoan language, since there is also a UD treebank for Shipibo-Konibo (Vasquez et al., 2018). The existence of a treebank for two Panoan languages allowed us to conduct some experiments on automatic part-of-speech tagging and syntactic dependency parsing in monolingual and transfer learning settings. We did not find a consistently positive impact of transfer learning from the Shipibo-Konibo treebank. However, the results strongly suggest that annotating a small preliminary version of a UD treebank for a minority language can be helpful for reducing annotation efforts in further iterations.

We also discussed here the collaborative methodology implemented for the creation of the Kakataibo treebank, which was conceived as part a regular Computational Linguistic course for linguistics undergraduates

in Peru. The methodology proposed here proved efficient to promote collaborative work among researchers and students in order to produce a full treebank of an endangered language in a limited time frame and in a formative setting. The idea behind the methodology implemented was to promote a bridge between descriptive linguistics and NLP developments, by means of building a UD treebank based on example sentences in published high-quality grammars. We are optimistic about the possibility of replicating our methodology for future similar projects and envisage a near future with larger numbers of endangered languages annotated in the UD framework.

Finally, the resources and experimentation details for reproducibility are published in: <https://github.com/Tarotis/Building-an-Endangered-Language-Resource-in-the-Classroom>

Acknowledgements

The first author acknowledges the support of CONCYTEC-ProCiencia, Peru, under the contract 183-2018-FONDECYT-BM-IADT-MU from the funding call E041-2018-01-BM.

7. Bibliographical References

- Agić, Ž. and Vulić, I. (2019). JW300: A wide-coverage parallel corpus for low-resource languages. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3204–3210, Florence, Italy, July. Association for Computational Linguistics.
- Akbik, A., Blythe, D., and Vollgraf, R. (2018). Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649.
- Akbik, A., Bergmann, T., Blythe, D., Rasul, K., Schweter, S., and Vollgraf, R. (2019). Flair: An easy-to-use framework for state-of-the-art NLP. In *NAACL 2019, 2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 54–59.
- Bender, E. M. (2007). Combining research and pedagogy in the development of a crosslinguistic grammar resource. In Tracy Holloway King et al., editors, *Proceedings of the GEAF07 Workshop*, pages 26–45, Stanford, CA. CSLI.
- Bender, E. M. (2009). Linguistically naïve!= language independent: Why nlp needs linguistic typology. In *Proceedings of the EACL 2009 Workshop on the Interaction between Linguistics and Computational Linguistics: Virtuous, Vicious or Vacuous?*, pages 26–32.
- Croft, W., Nordquist, D., Looney, K., and Regan, M. (2017). Linguistic typology meets universal dependencies. In *TLT*, pages 63–75.
- de Marneffe, M.-C., Dozat, T., Silveira, N., Haverinen, K., Ginter, F., Nivre, J., and Manning, C. D. (2014). Universal Stanford dependencies: A cross-linguistic typology. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 4585–4592, Reykjavik, Iceland, May. European Language Resources Association (ELRA).
- de Marneffe, M.-C., Manning, C. D., Nivre, J., and Zeman, D. (2021). Universal Dependencies. *Computational Linguistics*, 47(2):255–308, 07.
- Dozat, T. and Manning, C. D. (2017). Deep biaffine attention for neural dependency parsing. In *ICLR 2017*.
- Hiippala, T. (2021). Applied language technology: NLP for the humanities. In *Proceedings of the Fifth Workshop on Teaching NLP*, pages 46–48, Online, June. Association for Computational Linguistics.
- Hinrichs, E., Hinrichs, M., Kübler, S., and Trippel, T. (2019). Language technology for digital humanities: introduction to the special issue. *Language Resources and Evaluation*, 53(4):559–563, Dec.
- Makazhanov, A., Sultangazina, A., Makhambetov, O., and Yessenbayev, Z. (2015). Syntactic Annotation of Kazakh: Following the Universal Dependencies Guidelines. A report. In *3rd International Conference on Turkic Languages Processing, (TurkLang 2015)*, pages 338–350.
- Pieter Muysken, et al., editors. (2016). *South American Indigenous Language Structures (SAILS) Online*. Max Planck Institute for Evolutionary Anthropology, Jena.
- Nivre, J., de Marneffe, M.-C., Ginter, F., Hajič, J., Manning, C. D., Pyysalo, S., Schuster, S., Tyers, F., and Zeman, D. (2020). Universal Dependencies v2: An evergrowing multilingual treebank collection. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4034–4043, Marseille, France, May. European Language Resources Association.
- Pereira-Noriega, J., Mercado-Gonzales, R., Melgar, A., Sobrevilla-Cabezudo, M., and Oncevay-Marcos, A. (2017). Ship-lemmatagger: Building an nlp toolkit for a peruvian native language. In Kamil Ekštejn et al., editors, *Text, Speech, and Dialogue*, pages 473–481, Cham. Springer International Publishing.
- Petrov, S., Das, D., and McDonald, R. (2012). A universal part-of-speech tagset. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2089–2096, Istanbul, Turkey, May. European Language Resources Association (ELRA).
- Stenström, E. (2016). conllu. <https://github.com/EmilStenstrom/conllu/>.
- Sylak-Glassman, J., Kirov, C., Post, M., Que, R., and Yarowsky, D. (2015). A universal feature schema for rich morphological annotation and fine-grained cross-lingual part-of-speech tagging. In Cerstin Mahlow et al., editors, *Systems and Frameworks for Computational Morphology*, pages 72–93, Cham. Springer International Publishing.

- Tessmann, G. (1930). *Die Indianer Nordost-Perus: grundlegende Forschungen für eine systematische Kulturkunde*, volume 2 of *Veröffentlichung der Harvey-Bassler-Stiftung*. Hamburg, Hamburg.
- Tyers, F. M. and Washington, J. N. (2015). Towards a Free/Open-source Universal-dependency Treebank for Kazakh. In *3rd International Conference on Turkic Languages Processing, (TurkLang 2015)*, pages 276–289.
- Tyers, F. M., Sheyanova, M., and Washington, J. N. (2018). Ud annotatrix: An annotation tool for universal dependencies. In *Proceedings of the 16th International Workshop on Treebanks and Linguistic Theories (TLT16)*, pages 10–17.
- Vajjala, S. (2021). Teaching NLP outside linguistics and computer science classrooms: Some challenges and some opportunities. In *Proceedings of the Fifth Workshop on Teaching NLP*, pages 149–159, Online, June. Association for Computational Linguistics.
- Vasquez, A., Ego Aguirre, R., Angulo, C., Miller, J., Villanueva, C., Agić, Ž., Zariquiey, R., and Oncevay, A. (2018). Toward Universal Dependencies for Shipibo-konibo. In *Proceedings of the Second Workshop on Universal Dependencies (UDW 2018)*, pages 151–161, Brussels, Belgium, November. Association for Computational Linguistics.
- Zariquiey, R., Hammarström, H., Arakaki, M., Oncevay, A., Miller, J., García, A., and Ingunza, A. (2019). Obsolescencia lingüística, descripción gramatical y documentación de lenguas en el Perú: hacia un estado de la cuestión. *Lexis*, 43(2):271–337.
- Zariquiey, R. (2011). Aproximación dialectológica a la lengua cashibo-cacataibo (pano). *Lexis*, 35(1):5–46.
- Zariquiey, R. (2013). Tessmann’s nokamán: a linguistic investigation of a mysterious panoan group. *Cadernos de Etnolingüística*, 5(2):1–48.
- Zariquiey, R. (2018). *A grammar of Kakataibo*, volume 75 of *Mouton Grammar Library*. Mouton, Berlin.
- Zeman, D. and Resnik, P. (2008). Cross-language parser adaptation between related languages. In *Proceedings of the IJCNLP-08 Workshop on NLP for Less Privileged Languages*.
- Zhang, Y., Li, Z., and Min, Z. (2020). Efficient Second-Order TreeCRF for Neural Dependency Parsing. In *Proceedings of ACL*, pages 3295–3305.

Appendix A: Part-of-speech tags used in the datasets

	upos	n _{Shipibo}	n _{Kakataibo}	freq _{Shipibo}	freq _{Kakataibo}
1	ADJ	153	29	0.03	0.03
2	ADP	36	6	0.01	0.01
3	ADV	144	30	0.03	0.03
4	AUX	204	17	0.04	0.02
5	CCONJ	91	4	0.02	0.00
6	DET	133	33	0.03	0.03
7	INTJ	35	1	0.01	0.00
8	NOUN	646	216	0.14	0.21
9	NUM	26	5	0.01	0.00
10	PART	956	415	0.20	0.40
11	PINT	5		0.00	
12	PRON	451	78	0.10	0.07
13	PROPN	58	15	0.01	0.01
14	PUNCT	867	127	0.19	0.12
15	SCONJ	1	2	0.00	0.00
16	SUFN	2		0.00	
17	SUFV	5		0.00	
18	SYM	4		0.00	
19	VERB	855	164	0.18	0.16
20	VERB_AUX	1		0.00	
21	X	7		0.00	
	total	4680	1142	1	1

Table 6: Part-of-speech tags in the datasets. For both the Shipibo-Konibo and Kakataibo datasets, we report the number of annotated tags (n) and the proportion of each tag with respect to all the tags in the dataset (freq).

	upos	Training	Dev	Test
1	ADJ	19	7	3
2	ADP	6	0	0
3	ADV	17	9	4
4	AUX	9	5	3
5	CCONJ	3	0	1
6	DET	15	14	4
7	INTJ	1	0	0
8	NOUN	128	51	37
9	NUM	3	0	2
10	PART	235	89	91
11	PRON	41	15	22
12	PROPN	11	1	3
13	PUNCT	76	26	25
14	SCONJ	2	0	0
15	VERB	97	36	31
	total	663	253	253

Table 7: 60/20/20 Split for the UPOS-tags of the Kakataibo dataset

Appendix B: Dependency relations used in the dataset

	deprel	n _{Shipibo}	n _{Kakataibo}	freq _{Shipibo}	freq _{Kakataibo}
1	acl	5		0.00	
2	advcl	114	40	0.02	0.04
3	advmod	144	48	0.03	0.05
4	amod	130	25	0.03	0.02
5	appos	6	8	0.00	0.01
6	aux	213	14	0.05	0.01
7	aux:val	234		0.05	
8	case	573	131	0.12	0.13
9	cc	92	1	0.02	0.00
10	ccomp	1	4	0.00	0.00
11	compound	79	9	0.02	0.01
12	conj	41	2	0.01	0.00
13	cop	137	2	0.03	0.00
14	det	139	20	0.03	0.02
15	discourse	7	5	0.00	0.00
16	flat	9	1	0.00	0.00
17	iobj	22	2	0.00	0.00
18	Lfcl	183		0.04	
19	marker	1		0.00	
20	nmod	69	76	0.01	0.07
21	nsubj	538		0.11	
22	nummod	17	21	0.00	0.02
23	obj	256	66	0.05	0.06
24	obl	123	48	0.03	0.05
25	punct	865	127	0.18	0.12
26	root	667	120	0.14	0.11
27	vocative	2	1	0.00	0.00
28	x	1		0.00	
29	xcomp	12		0.00	
30	aux:dub		1		0.00
31	aux:ev		22		0.02
32	aux:int		1		0.00
33	aux:sgen		121		0.12
34	csubj		1		0.00
35	dislocated		1		0.00
36	list		2		0.00
37	nsubj:bound		116		0.11
38	nsubj:free		97		0.09
39	parataxis		9		0.01
	total	4680	1142	1	1

Table 8: Dependency relations used in the datasets. For both the Shipibo-Konibo and Kakataibo datasets, we report the number of annotated tags (n) and the proportion of each tag with respect to all the tags in the dataset (freq).

deprel	Experiment 6 (60/20/20)			Experiment 7 (80/10/10)		
	Training	Dev	Test	Training	Dev	Test
1 advcl	25	8	7	30	5	5
2 advmod	27	13	8	41	2	5
3 amod	15	7	3	22	2	1
4 appos	6	1	1	8	0	0
5 aux	6	5	3	11	2	1
6 aux:dub	1	0	0	1	0	0
7 aux:ev	10	8	4	15	3	4
8 aux:int	0	0	1	1	0	0
9 aux:sgen	68	27	26	94	13	14
10 case	76	27	28	103	16	12
11 cc	0	0	1	0	0	1
12 ccomp	4	0	0	4	0	0
13 compound	9	0	0	9	0	0
14 conj	1	0	1	1	0	1
15 cop	2	0	0	2	0	0
16 csubj	1	0	0	1	0	0
17 det	10	9	1	13	5	2
18 discourse	3	1	1	5	0	0
19 dislocated	1	0	0	1	0	0
20 flat	1	0	0	1	0	0
21 iobj	1	1	0	2	0	0
22 list	0	2	0	0	2	0
23 nmod	42	18	16	59	6	11
24 nsubj:bound	67	25	24	92	13	11
25 nsubj:free	59	17	21	77	11	9
26 nummod	13	4	4	17	2	2
27 obj	40	15	11	53	5	8
28 obl	23	12	13	35	7	6
29 parataxis	4	3	2	7	0	2
30 punct	76	26	25	101	12	14
31 root	72	24	24	96	12	12
32 vocative	0	0	1	1	0	0
total	663	253	253	903	118	121

Table 9: Split for the dependency relations for the Kakataibo dataset