

A Morphological Transducer for the Limbu Language

Avyaya Singh, Jonathan North Washington

Swarthmore College
500 College Avenue, Swarthmore, PA 19081

{asingh5, jwashin1}@swarthmore.edu

Abstract

This paper presents the first ever morphological transducer for the Limbu language, also known by its endonym Yakthung Pan, an endangered Sino-Tibetan language primarily spoken in the area known as Limbuwan/Koshi Province in Eastern Nepal, with a minority population in the Sikkim state of India, where Limbu enjoys official status. Using a corpus of Limbu text produced by field interviews (Michailovsky, 1977), and a translation of the Holy Bible into the Limbu language, the paper presents various elements of the morphology of the language, how they were implemented into the transducer, and an evaluation of the transducer against identified corpora and a gold standard. With a relatively small lexicon, the transducer was found to have reasonable coverage (~60%), with high precision (88%) but low recall (32%). The paper discusses future expansion through further involvement with the community, which can help in the maintenance and revitalisation of the endangered language.

Keywords: Limbu, morphological transducers, morphological analysis, Tibeto-Burman languages, South Asia, Nepal

1. Introduction

This paper describes the creation of the first-ever morphological transducer for the Limbu language of Eastern Nepal, using tools from Apertium, a free and open source digital platform intended primarily for rule-based machine translation.

Linguistic preservation and revitalisation of Limbu is currently being undertaken by groups such as Yakthung Cho, an online Limbu Art collective, run by Subhas Thebe Limbu, a Limbu-origin film director, which seeks to produce Limbu language learning tools through online classes and produce content in the Limbu language aimed at informing the Limbu audience of issues pertaining primarily to their efforts at language revitalisation and autonomist demands in Eastern Nepal.¹ The first author of this paper themselves joined such a class conducted by Yakthung Cho in the English medium to learn the Limbu language.

A morphological transducer, here implemented as a finite-state transducer (FST), both analyses text in the language, providing tokenisation of the text and a morphological analysis of each token (morphological analysis), and creates tokens from input analyses (morphological generation). Morphological transducers can prove helpful in learning vocabulary (Katinskaia et al., 2018) and morphology (Arppe et al., 2022); as elements in word-form creators (Fernald et al., 2016; Kazantseva et al., 2018); and as the basis for spell checkers (Washington et al., 2021). They can further also be used in rule-based and hybrid machine translation systems (Khanna et al., 2021), and in a range of other NLP tasks and pipelines.

This paper is structured as follows. Section 2 provides background into the Limbu language and its general structure. Section 3 describes

the methodology of creating the transducer, and has details on encountered challenges. Section 4 presents an evaluation of the transducer, showing its coverage over various corpora as well as precision and recall over a hand-annotated gold standard. Section 5 concludes and discusses future work. Associated tools, including a keyboard and a disambiguator, are also described.

2. Background on Limbu/Yakthung

Limbu, known by its endonym Yakthung Pan (ཡཱུར་ཤུང་པཎ), is spoken by 350,436 people in Nepal (National Statistics Office, 2025), and by 40,835 speakers in North-Eastern India (Office of the Registrar General and Census Commissioner, 2018), mostly belonging to the Limbu ethnicity.² Figure 1 shows the geographic distribution of speakers.

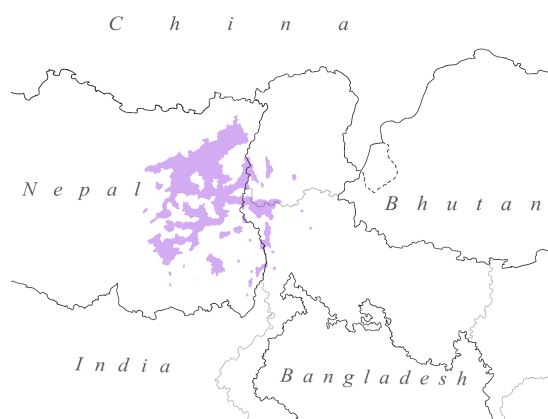


Figure 1 : A map highlighting the regions in which the Limbu language is spoken, namely, Eastern Nepal and the Sikkim state of India.³

¹ <https://www.instagram.com/yakthungcho/?hl=en>

² <https://www.ethnologue.com/language/lif/>

³ <https://commons.wikimedia.org/wiki/File:Limbu.map.png>

Limbu belongs to the Kiranti subfamily of the Tibeto-Burman languages within the wider Sino-Tibetan language family. The Kiranti subfamily, spoken in Eastern Nepal includes related languages such as Sunuwar, Rai, and Yakkha. Limbu is spoken primarily in the region known to the Limbu people as Limbuwan, called the Koshi Province by the Nepalese government, and it enjoys official status in the province.⁴ Although the province is autonomous, due to the federalisation of the Nepalese government following the republican revolution of 2006,⁵ the Limbu people have been demanding the reinstatement of Limbuwan through the renaming of the Koshi province, as these lands are seen as ancestral to them.⁶ The language is also spoken by a minority in the Indian state of Sikkim which borders this province to the east, where it also enjoys official status.

Limbu is a highly polysynthetic language: words are inflected using a series of affixes, and multiple words can fuse together. The language has four main dialects: Phedape, Chhathare, Tambarkhole, and Panthare (van Driem, 1987). Interesting linguistic features include dual number, split ergativity, and double (sometimes triple) negation, and primarily Subject-Object-Verb word order. Limbu uses its own indigenous orthography, called the Sirijonga script (§3.3.1).⁷

The phonemic inventory of Limbu is very similar to those of nearby Indo-Aryan languages, although it is itself Sino-Tibetan; this is because it is part of the Nepali Sprachbund. Thus, its phonological inventory and orthography follow the Indo-Aryan *Varga* system, where consonants are divided into 5 *Vargas* according to place of articulation. Thus, in terms of its phoneme inventory, the language and its orthography are comparable to Nepali, Bengali, Hindi, etc. (van Driem, 1987).

Only two reference grammars on the language were available to the authors, both primarily concerning the aforementioned Chhathare dialect: van Driem (1987) and Tumbahang (2007). These grammars were consulted in the creation of the transducer. Because of the more specific dialectal focus of Tumbahang (2007), words/phrases encountered in this source were included as synonyms/alternate spellings of words found in van Driem (2007), in an effort to make the transducer dialect-aware but generalisable for the Limbu language as a whole.

⁴ <https://www.ethnologue.com/language/lif/>

⁵ <https://interactive.aljazeera.com/aje/2016/nepal-maoi-st-dream/index.html>

⁶ <https://kathmandupost.com/province-no-1/2023/03/12/general-strike-throws-life-out-of-gear-in-koshi-province>

⁷ <https://www.ethnologue.com/language/lif/>

The authors do not speak the Limbu language, but this transducer was created through close study of the aforementioned reference grammars and deductions made through use of linguistic analysis. Future work will include native speakers.

3. Methodology & Implementation

This section goes over the details of the transducer, namely how it was constructed (§3.1), how its lexicon was expanded (§3.2), and solutions to challenges the authors encountered while creating the transducer (§3.3).

3.1 Tools used

The transducer described was created entirely by hand and is freely and publicly available under the GPL v3 Free/Open Source license.⁸ The intent is to increase access, as no other transducer is known to exist for the Limbu language, and resources on the language are scarce. The transducer is also available online for morphological analysis and generation via Apertium's web UI.⁹ Machine learning approaches were not used, due to their need for large corpora and multiple development cycles, the difficulty in correcting errors, and their lack of precision (Butt, 2021).

The lexd formalism (Swanson and Howell, 2021) was used to implement the morphotactics of Limbu. Previously published lexd transducers showcase its ability to handle non-suffixational morphology (Washington et al., 2021; Christopherson, 2023) and some discontinuous morphology (Swanson and Tyers, 2022), both of which are used extensively in Limbu. The transducer uses the Apertium framework (Forcada et al., 2011; Khanna et al., 2021) for compilation scripts and the HFST format and tools (Lindén et al., 2011) for storing and working with the compiled transducer. Apertium also offers a handful of tools for leveraging a morphological transducer for language maintenance, learning, and revitalisation, which are all future directions the authors intended the transducer for.

The morphophonology was implemented in the twol formalism, included with HFST. The analyser and generator are compiled by compose-intersecting the morphotactics transducer and the morphophonology transducer, as in sources cited above for other lexd transducers.

3.2 Lexicon

The transducer was developed against the only two available digital corpora of Limbu text of any substantial size: the Wycliffe translation of the Bible into the Limbu language (Wycliffe, 2009) and a collection of field interviews in the Limbu language compiled primarily by linguists Boyd

⁸ <https://github.com/apertium/apertium-lif.git>

⁹ <https://beta.apertium.org/#analysis?aLang=lif>

-3	0	+2	+4	+7	+8	+9	+12
𑄎- a- SUBJ.1PL.INCL-	𑄎𑄚 ab shoot	-∅ -NPST	-𑄚 -u -OBJ.3	-𑄚 -m -SUBJ.1PL	-𑄚 -si -SUBJ.NSG	-𑄚 -m -SUBJ.1PL	-∅ -PERF
𑄚- ke- SUBJ.2PL-	𑄚𑄚 baʔr- speak	-∅ -NPST	-𑄚 -i -OBJ/SUBJ.PL		-∅ -OBJ.SG		-∅ -PERF

Table 2: Morphological breakdown of two example verb forms: (a) 𑄎𑄚𑄚𑄚𑄚𑄚𑄚𑄚𑄚 *aabumsim* ‘We (incl) have shot them’ and (b) 𑄚𑄚𑄚𑄚 *kebaʔri* ‘You (pl) speak (obj)’. ∅ indicates an empty slot corresponding to a default.

"{Kk} becomes 𑄚 if succeeded by a vowel"
 %>{Kk%}:𑄚 <=> _ %>: (:𑄚) (:0) :Vowel ;

Figure 4: Twol rule for the {Kk} symbol.

𑄚𑄚𑄚𑄚𑄚𑄚𑄚𑄚<vbdo><inf>
 𑄚𑄚𑄚𑄚𑄚𑄚𑄚𑄚<vbdo><nonpast><s_sg3><o_3><o_sg><perf>

Figure 5: Example forms analysed by the transducer involving the twol rule from Figure 2, with a [k] in the first example a [g] in the second. The ‘:’ separates forms from analyses.

In order to fully implement the morphophonological alternations of the language, 45 twol rules were encoded.

In addition to this project, the first author also created a keyboard for the Sirijonga script, as currently available keyboards were not fully intuitive, as cited by members of the community (p.c., February, 2025). Existing keyboards were based on QWERTY,¹⁹ created by a non-Limbu Nepali artist,²⁰ or based on the consonant *Varga* system (available in Google's Gboard²¹). The keyboard, developed based on glyph-frequency, was made in consultation with members of the Limbu community, who commented on its effectiveness, and it is available under the MIT License on GitHub as a free download.²² It was first created using iBus-m17n,²³ and documentation of the layout is also available.²⁴ In order for seamless usage across platforms, the first author worked with Translations Commons, a language digitisation NGO, on making a mobile-friendly version of the author's produced keyboard. This Free/Open Source keyboard, made using the Keyman system,²⁵ is supported on iPadOS, Android, iOS, and Windows, and is available for free distribution by the author (60 downloads over the last 2 months) and on the Translations Commons

¹⁹ https://keyman.com/keyboards/sil_limbu_phonetic
²⁰ https://keyman.com/keyboards/sil_limbu_typewriter
²¹ <https://support.google.com/gboard/answer/6380730#zippy=%2Cfind-supported-languages>
²² https://github.com/SwatLangTech/Limbu_Keyboard
²³ <https://github.com/ibus/ibus-m17n>
²⁴ <https://wikis.swarthmore.edu/ling073/Limbu/Keyboard>
²⁵ <https://keyman.com/>

website,²⁶ as well as on the official Keyman website.

3.3.2 Verbal Morphology

Being a highly polysynthetic language, Limbu can agglutinate up to ~15 suffixes and ~5 prefixes (including cliticised adverbs) to a single verb root, and also incorporate other parts of speech into a verb form (van Driem, 1987). There are three negation affixes that can appear in Limbu also verbal morphology, amounting to the same meaning (the language does not possess double negation as a form of emphasis). For each verb, no affix is mandatory, and thus, a root can exist independently, or with the addition of up to all ~20 affixes. The absence of a filled affixal slot can either mean ‘default’ for that particular slot (singular for number, first person for person, active for voice, etc.) or simply be an open slot, not providing any additional meaning to the verb. Example verb forms are shown in Table 2, and one pattern for verbal morphology included in the transducer is shown in Figure 6.

Adverbs? + SPN(1) + SPN(4) + Neg?(1) + **Verb-Stems** + Neg?(2) + Reflexivity? + Tense + SPN(5) + SPN(2) + ObjSlot + SPN(6) + NegS? + SPN(7) + ObjNum + SPN(8) + SPN(3) + NegS2? + Mode? + Aspect + Clitic? + Nom? + Subord? + *Adverbs?*

Figure 6: One verbal morphology pattern in the transducer, slightly simplified. *Italics* (added for demonstration) represent lexicons containing foreign roots, **Bold** (added for demonstration) represents the lexicon containing verbal roots, SPN represents the lexicon containing affixes related to the subject's person and number, Neg represents negation lexicons. Numbers in parentheses represent lexicon columns. Some columns include only tags that match prefix morphemes.

Although 13 other such patterns for verbs exist, as verbal morphology varies greatly in the language, the patterns are precise, meaning an error cannot occur in identifying which pattern a verb follows (and thus which slot an affix falls in). Patterns were identified based on van Driem (1987), Michailovsky (1977), and linguistic analysis on forms not encountered in these sources. The morphology of nouns is only somewhat less complex, with 7 patterns being

²⁶ <https://reco.to/JZazvh>, <https://translationcommons.org/keyboard-creation-projects/published-keyboards/>

included. The patterns included so far in the transducer cover the most common morphological patterns in the language.

3.3.3 Ambiguity

The aforementioned 15+5 affixes have many variants, due to phonological alternations, dialectal variation, and the fact that there appear to simply be synonymous morphemes. This is compounded by the fact that there exist 13-15 noun cases (van Driem, 1987) in this language, with each having between 2-5 lexical variants, with the genitive case having the most variants, at 11 (Figure 7), attributed to both orthographic variation and dialectal differences in the language. In the lexd file, the first line of every block of synonymous morphemes is used to denote the ‘standard form’, selected based on the reference grammars. The subsequent variants are either dialectal variants as denoted by Tumbahang (2007) or lexical variants encountered in the Pangloss collection of Limbu interviews. All forms besides the “standard form” include Dir/LR in the comment, which allows them to be included in the analysis transducer but excluded from the generation transducer.²⁷ Currently, the case lexicon itself contains 51 lines of code in order to account for these variations, 36 of which are considered variants.

```

<gen>:>{NG}{A}{n}
<gen>:>ꠔ # Dir/LR
<gen>:>ꠕ # Dir/LR
<gen>:>{i}ꠕ # Dir/LR
<gen>:>{ie}{i}ꠕ # Dir/LR
<gen>:>{a}{ie}{i}ꠕ # Dir/LR
<gen>:>ꠔꠕ # Dir/LR
<gen>:>ꠕꠕ # Dir/LR
<gen>:>{i}ꠕꠕ # Dir/LR
<gen>:>{a}{ie}{i}ꠕꠕ # Dir/LR

```

Figure 7: Variants of the genitive case based on dialectal or orthographical variation.

Additionally, some morphemes appear to be syncretic with other morphemes, despite being rarely used. For example, *-re/* is one of the alternative genitive case markers, but *-re/* can also mark vocative and ergative. Hence any token containing *-re/* has at least 3 analyses.

A disambiguator was created in VISL CG-3 (Bick and Didriksen, 2015) to reduce the ambiguity of analyses by context. An excerpt is provided in Figure 8—one heuristic which deals with the ambiguity of the genitive suffix, with an example of its application in Figure 9.

```

"<ꠔꠕ>"
: "ꠔ" n sg gen SELECT:85
: "ꠔ" n sg erg SELECT:85
: "ꠔ" n sg voc SELECT:85
"<ꠔꠕꠕ>"
: "ꠔꠕꠕ" n sg
: "ꠔꠕꠕ" n ser3sg
: "ꠔꠕꠕ" n attr REMOVE:293

```

²⁷ See Washington et al. (2021, 187) for details on this approach.

Figure 9: Disambiguation of forms in context (‘wheat-GEN bread’) employing the rule in Figure 8 (“85”). Analyses with preceding ; are removed by the rules from the list of possible analyses.

```

LIST NOUN = n np ;
LIST NGEN = (n gen) (np gen) ;
SELECT NGEN IF
(0 NGEN)
(1 NOUN) ;

```

Figure 8: A CG-3 rule from the disambiguator.

The Limbu disambiguator (apertium-lif.lif.rlx) consists of 88 rules. The performance of the disambiguator is discussed in §4.1.

4. Evaluation

For the sake of evaluating this transducer, several corpora were prepared using the previously described only extant digital corpora of Limbu text in the Sirijonga orthography: the Bible corpus. Genesis 1 & 2 as well as 10 sentences (of 324) from the *Paddy Dancing* interview were used for development and are hence treated separately.

4.1 Naive Coverage and Ambiguity

The transducer²⁸ was evaluated in terms of naïve coverage, which is defined as the number of forms in a corpus which return an analysis, irrespective of whether this analysis is correct or not. The ambiguity of the transducer on these corpora, defined as the average number of analyses returned by the analyser per analysed token, was also calculated with and without the disambiguator.²⁹ The results are in Table 2.

corpus	tokens	coverage	ambiguity before / after
10 sents	177	100.00 %	3.82 / 1.62
Gen. 1 & 2	1253	63.05 %	3.03 / 1.69
Full Bible	759K	60.62 %	2.28 / 1.44

Table 2: Naive coverage and ambiguity on development and test corpora.

It’s seen that coverage is reasonable for a young transducer, at over 60%. Ambiguity is reduced greatly by the disambiguator, from ~3 analyses per token to around 1.5.

4.2 Accuracy

Precision and recall were measured to understand the accuracy of the transducer. Precision is defined as the percentage of analyses returned by the transducer for each form that are correct. Recall is defined as the percentage of the correct analyses for each form that are returned as opposed to not being returned at all.

²⁸ The transducer was last evaluated on revision 4b0ebbb in March, 2026.

²⁹ Scripts used for calculating coverage, ambiguity, precision, and recall are available in the dev/ directory in the transducer’s git repository.

In order to measure accuracy, the first author randomly selected 100 unique tokens (of ~4k) from the *Paddy Dancing* interview, hand-transliterated them, and annotated them manually to create a gold standard. The words in the gold standard appeared to be generally representative of the lexicon and morphology of texts the authors have worked with. The gold standard attempted to capture all possible analyses for each token, by using the glosses provided by Mikhailovsky (1977) and leveraging the first author's knowledge of Limbu morphology. The gold standard was compared to the list of analyses returned by the analyser.

Precision was measured at ~88% (108/122 analyses), while recall was only ~32% (108/340 analyses), with an F-score of 0.4675; i.e., most analyses returned by the transducer were correct (except in two instances, noted below), but many correct analyses were not returned due to the fact that the transducer has low coverage on the interview.

A preliminary qualitative evaluation of forms that did not have an analysis returned by the transducer yields several patterns. Besides missing lexical items (including verb stems, adverbs, nouns, classifiers, etc.), there were a few uncommon morphotactic patterns not implemented in the transducer, mostly for verbs, such as polar question and reduplicated forms. Additionally, there are some mistakes in pattern definitions that cause some patterns not to function fully as intended; e.g., while the language allows up to three negation morphemes in a verb, and the morphotactics were designed to support this, currently the transducer only recognises up to two negation morphemes in a verb form.

Given the fact that the majority of analyses not returned by the transducer were caused by stems missing from the lexicon, the low recall result substantiates the fact that the stems in the lexicon are indeed more common than those missing. Assuming similar numbers of possible analyses per form, and given the random unique sampling approach used to build the gold standard, if included and missing stems were evenly distributed in frequency, then recall would approximate coverage.

There were very few forms with incorrect analyses returned from the transducer; all involve hypothetical but unattested verb forms which are synonymous with another word. For example, the word चिऱ *char* refers to the numeral "4", however, it is also the predicted verb form involving the optative suffix -ऱ *-r* added to the verb चिऱि *chama*. However, चिऱ *char* is always used to reference the numeral, and an alternative optative form, चिऱि *chal*, is instead used for चिऱि *chama*.³⁰ Thus, the optative analyses returned by the transducer in this

³⁰ For other verbs, the two optative forms are used interchangeably.

instance are incorrect. Similarly, the word वाऱि *wa?i* refers to a polar question verb form of the stem वाऱि *wama* "to be" in all contexts, rather than any of the other analyses returned by the transducer, despite these being predicted by the morphotactics. One of these analyses involves object marking, which is currently incorrectly allowed for verbs which are not transitive.

5. Conclusion & Future Work

This paper presents the first ever morphological transducer for the endangered Limbu language, to the best of the authors' knowledge. With significant coverage on the Limbu Bible, and with high precision, this paper contributes to the linguistic description of the language, and introduces the ability to perform various text-processing tasks in the language, enhancing the digital presence of the language.

In terms of future work on the language, the authors believe that coverage can be expanded in two main ways: (1) by hand-transcribing Limbu print books, the Universal Declaration of Human Rights, and the interviews from Mikhailovsky (1977), automating transliteration of other sources, and expanding the lexicon based on their contents, and (2) working more closely with native speakers to make sense of forms encountered in these sources that are not able to be understood using published resources. The transducer is easily scalable through adding to the lexicon, tweaking the existing morphological patterns, and adding new ones.

As mentioned before, this transducer can also be used as a component in a rule-based machine translation system. Given that there are currently no MT systems for the Limbu language, the first author did create a basic Limbu-English MT system as a proof-of-concept. However, it is still quite barebones, and is only able to handle a few simple phrase types and words of the language. The authors seek to expand this MT system in the future.

This transducer and its associated projects (e.g., the Limbu keyboard) have already been introduced to the Limbu community. There appears to be some initial uptake, and community members have commented on the efficiency and usefulness of these projects.

Acknowledgments

The authors would like to thank Subhas Thebe Limbu^{31,32} and Sumi Limbu³³ for their support and for their help in looking over our work. We would also like to thank all members of the Limbu community for keeping this language alive and thriving, and Yakthung Cho for their contribution

³¹ subashthebe.com/

³² eyebeam.org/artists/subash-thebe-limbu/

³³ scholarworks.uark.edu/etd/5325/

to Limbu linguistic preservation through online classes.³⁴

References

- Arppe, Antti & Poulin, Jolene & Neitsch, Andrew & Harrigan, Atticus & Hieber, Daniel (2022). *itwêwina: Towards a morphologically intelligent and user-friendly online dictionary of Plains Cree*. 10.13140/RG.2.2.33073.08804.
- Bick, Eckhard & Didriksen, Tino. (2015). CG-3 - Beyond Classical Constraint Grammar. Proceedings of the 20th Nordic Conference of Computational Linguistics (NODALIDA 2015), Vilnius, Lithuania.
- Butt, Miriam (2020). Building Resources: Language Comparison and Analysis. At the Second Workshop on Computational Research in Linguistic Typology (SIGTYP 2020), invited to talk. https://youtu.be/D2AzyAY_3Mc.
- Christopherson, Cody Scott (2023). A finite-state morphological analyzer for Q'eqchi' using Helsinki Finite-State Technology (HFST) and the Giellatekno infrastructure. Master's thesis, Brigham Young University.
- van Driem, George (1987). *A Grammar of Limbu*, Berlin, New York: De Gruyter Mouton. <https://doi.org/10.1515/9783110846812>
- Fernald, Theodore B., Kashyap, Nabil, and Fahringer, Jeremy (2016). Navajo verb generator. Development version.
- Forcada, Mikel L., Ginestí-Rosell, Mireia, Nordfalk, Jacob, O'Regan Jim, Ortiz-Rojas, Sergio, Pérez-Ortiz, Juan Antonio, Sánchez-Martínez, Felipe, Ramírez-Sánchez, Gema, and Tyers, Francis M. (2011). Apertium: a free/open-source platform for rule-based machine translation. *Machine Translation*, 25(2):127–144.
- Gaenszle, Martin (2021). The Limbu Script and the Production of Religious Books in Nepal. *Philological Encounters*. 6 (1–2): 43–69. doi:10.1163/24519197-bja10014. ISSN 2451-9197.
- Katinskaia, A., Nouri, J., and Yangarber, R. (2018). Revita: a language-learning platform at the intersection of ITS and CALL. In Proceedings of LREC: 11th International Conference on Language Resources and Evaluation, Miyazaki, Japan.
- Kazantseva, Anna, Maracle, Owennatekha Brian, Maracle, Ronkwe'tiyóhstha Josiah, and Pine, Aidan (2018). Kawennón:nis: the wordmaker for Kanyen'kéha. In Proceedings of the Workshop on Computational Modeling of Polysynthetic Languages, pages 53–64, Santa Fe, New Mexico, USA. Association for Computational Linguistics
- Khanna, T., Washington, J., Tyers, Francis M., Bayatli, Sevilay, Swanson, Daniel G., Pirinen, Tommi A., Tang, Irene, and Alòs i Font, Hèctor (2021). Recent advances in Apertium, a free / open-source rule-based machine translation platform for lowresource languages. *Machine Translation*.
- Lindén, K., Silfverberg, Miikka, Axelson, Erik, Hardwick, Sam, and Pirinen, Tommi (2011). HFST— Framework for Compiling and Applying Morphologies, volume 100 of Communications in Computer and Information Science, pages 67–85. Springer.
- Michailovsky, Boyd; Everson, Michael (2002). L2/02-055: Revised proposal to encode the Limbu script in the UCS (PDF). Retrieved 2026.
- Michailovsky, Boyd, Mazaudon Martine (1977). Pangloss Collection | Limbu Corpus. URL: <https://pangloss.cnrs.fr/corpus/Limbu?lang=en&mode=norma>
- National Statistics Office (2025). Languages in Nepal. Kathmandu: National Statistics Office. (National Population and Housing Census 2021). URL: <https://censusnepal.cbs.gov.np/results/files/result-folder/Language%20in%20Nepal.pdf>
- Office of the Registrar General & Census Commissioner (2018). Language: India, States and Union Territories (C-16_25062018). In Census of India 2011. URL: <https://www.censusindia.gov.in/nada/index.php/catalog/42458>
- Swanson, Daniel and Howell Nick (2021). Lexd: A finite-state lexicon compiler for non-suffixational morphologies In Mika Hämmäläinen, et al., editors, *Multilingual Facilitation*, pages 133–146. Helsingin yliopisto.
- Swanson, Daniel G. and Francis M. Tyers (2022). Handling Stress in Finite-State Morphological Analyzers for Ancient Greek and Ancient Hebrew. In *Proceedings of the Second Workshop on Language Technologies for Historical and Ancient Languages*, pages 108–113, Marseille, France. European Language Resources Association. <https://aclanthology.org/2022.lt4hala-1.15/>
- Tumbahang, Govinda Bahadur (2007). *A Descriptive Grammar Of Chhatthare Limbu*.

³⁴ <https://docs.google.com/forms/d/e/1FAIpQLSctveBqsNhZIIzgUFA-4cneHCV1ghaDCOIBGCBdN04EzrDBg/viewform>

- Kathmandu: Central Department of English Kirtipur. hdl:123456789/3017.
- Tumbahang, Govinda Bahadur (2007). *A Descriptive Grammar Of Chhatthare Limbu*. Kathmandu: Central Department of English Kirtipur. hdl:123456789/3017.
- Tyers, Francis, Pirinen, Tommi, Washington, Jonathan (2015) Finite-state morphologies and text corpora as resources for improving morphological descriptions. *Workshop on Computational Phonology and Morphology*. University of Chicago. Poster presentation.
- Washington Jonathan, Lopez, Felipe, and Lillehaugen, Brook (2021). Towards a morphological transducer and orthography converter for Western Tlacolula Valley Zapotec. In Proceedings of the First Workshop on Natural Language Processing for Indigenous Languages of the Americas, pages 185–193, Online. Association for Computational Linguistics.
- Washington, Jonathan, Salimzianov, Ilnar, Tyers, Francis M., Gökırmak, Memduh, Ivanova, Sardana, and Kuyrukçu, Oğuzhan (2021). Free/open-source technologies for Turkic languages developed in the Apertium project. In Proceedings of the Seventh International Conference on Computer Processing of Turkic Languages (TurkLang 2019).
- Wycliffe Bible Translators, Inc. & Isai Limbu Literature Association (2009). *Limbu New Testament (Limbu Script)*. Digital Bible Society. URL: <https://www.scriptureearth.org/00eng.php?idx=1480&LN=Limbu>. Accessed: February, 2026. Currently only available as a PDF: <https://find.bible/bibles/LIFWBT/>