# A Bayesian phylogeny of Patkaian (Northern Naga)

## statistical methods with large data on small languages

Kellen Parker van Dam

**Universität Zürich**  Switzerland
Department of Comparative Language Science
Center for the Interdisciplinary Study of Language Evolution

**La Trobe University**  Australia
Department of Languages & Cultures

**Universität Passau**  Germany
Multicultural Computational Linguistics

Universität Zürich UZH

LA TROBE UNIVERSITY

FNSNF
FONDS NATIONAL SUISSE
DE LA RECHERCHE SCIENTIFIQUE
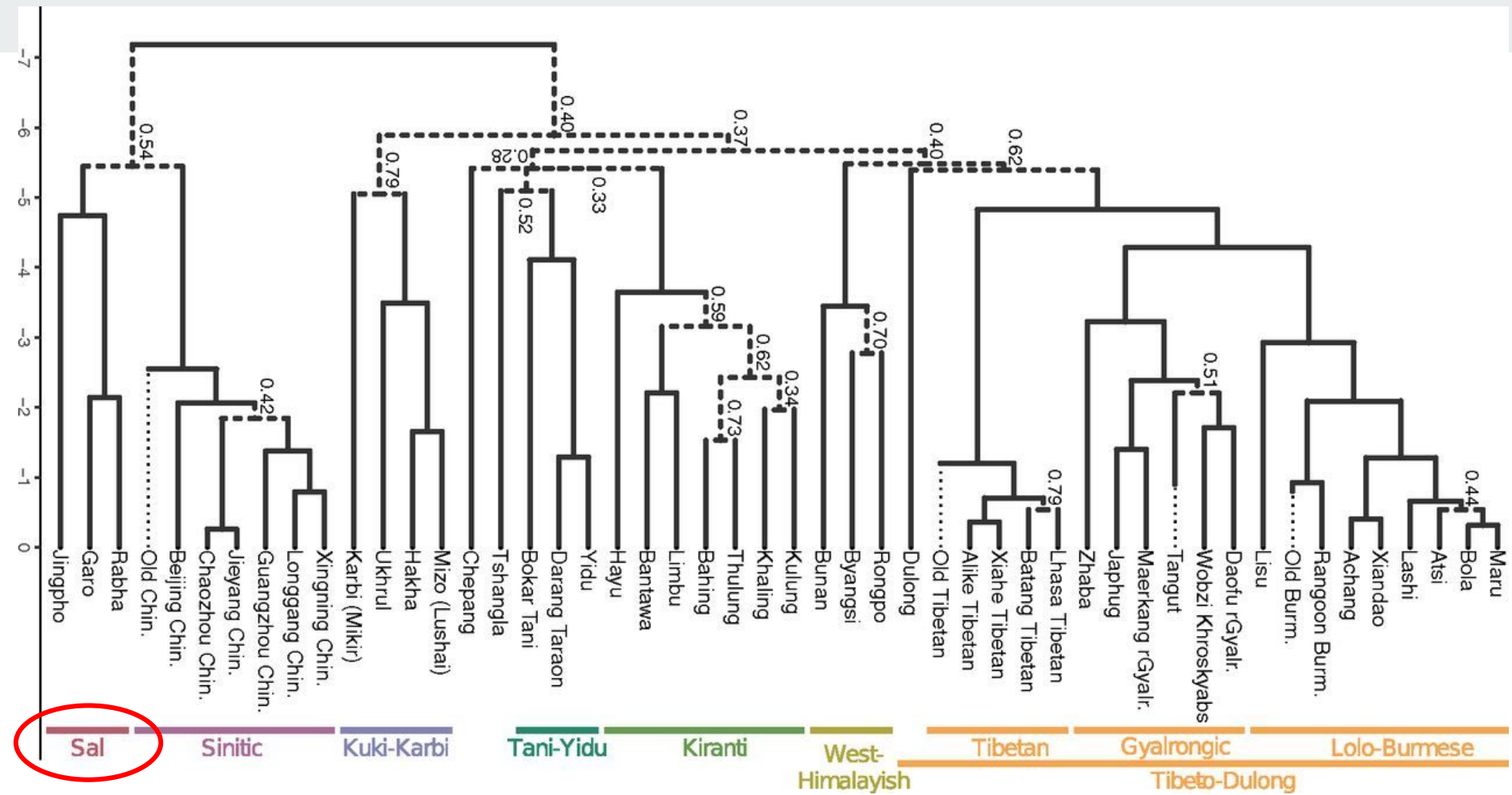
# Sino-Tibetan Phylogenies

Recent Sino-Tibetan phylogenies offer interesting suggestions of large-scale relationships in the family.
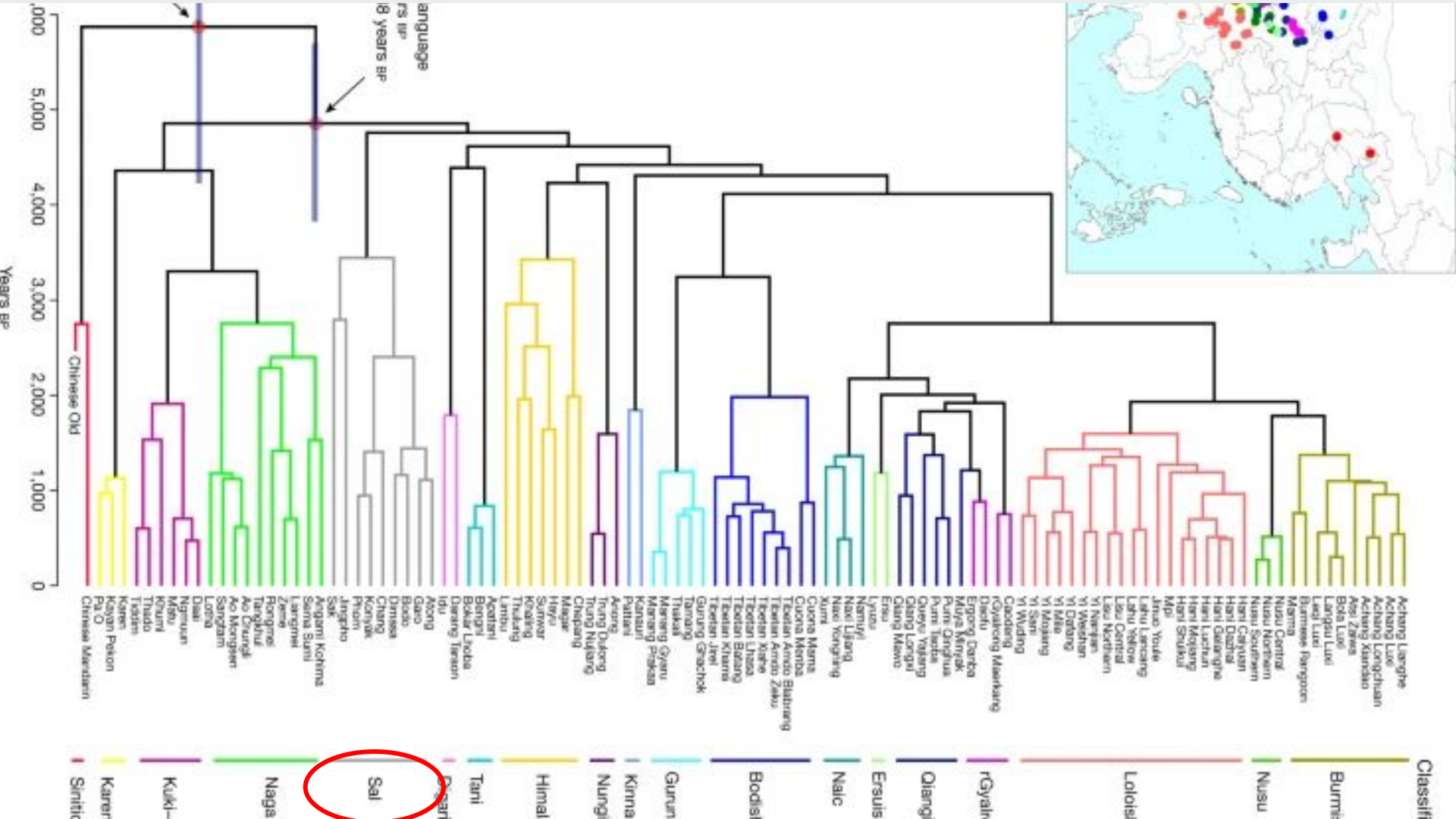
Methods and approaches vary, resulting in vastly different trees.

Posteriors are often quite low.

Sagart, L., Jacques, G., Lai, Y., Ryder, R.J., Thouzeau, V., Greenhill, S.J. and List, J.M., 2019. Dated language phylogenies shed light on the ancestry of Sino-Tibetan. Proceedings of the National Academy of Sciences, 116(21), pp.10317-10322.

Zhang, M., Yan, S., Pan, W. and Jin, L., 2019. Phylogenetic evidence for Sino-Tibetan origin in northern China in the Late Neolithic. Nature, 569(7754), pp.112-115.

Sal | Sinitic | Kuki-Karbi | Tani-Yidu | Kiranti | West-Himalayish | Tibetan | Gyalrongic | Lolo-Burmese | Tibeto-Dulong

# Sino-Tibetan Phylogenies

Recent Sino-Tibetan phylogenies offer interesting suggestions of large-scale relationships in the family.

However, significant differences in results and approach still leave us reliant on the "Fallen Leaves" model of van Driem (2012).
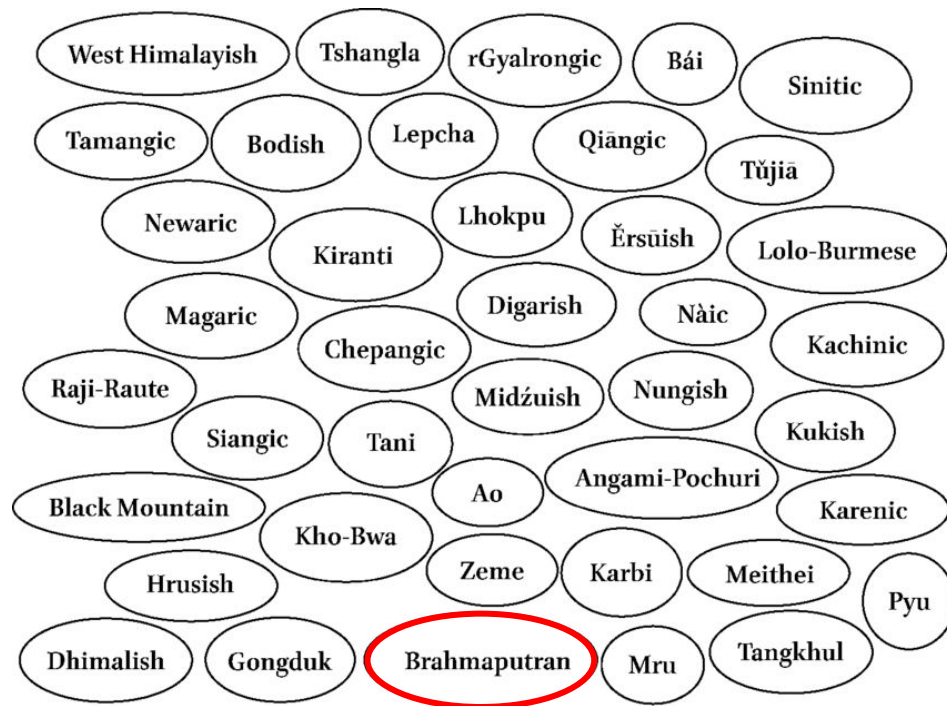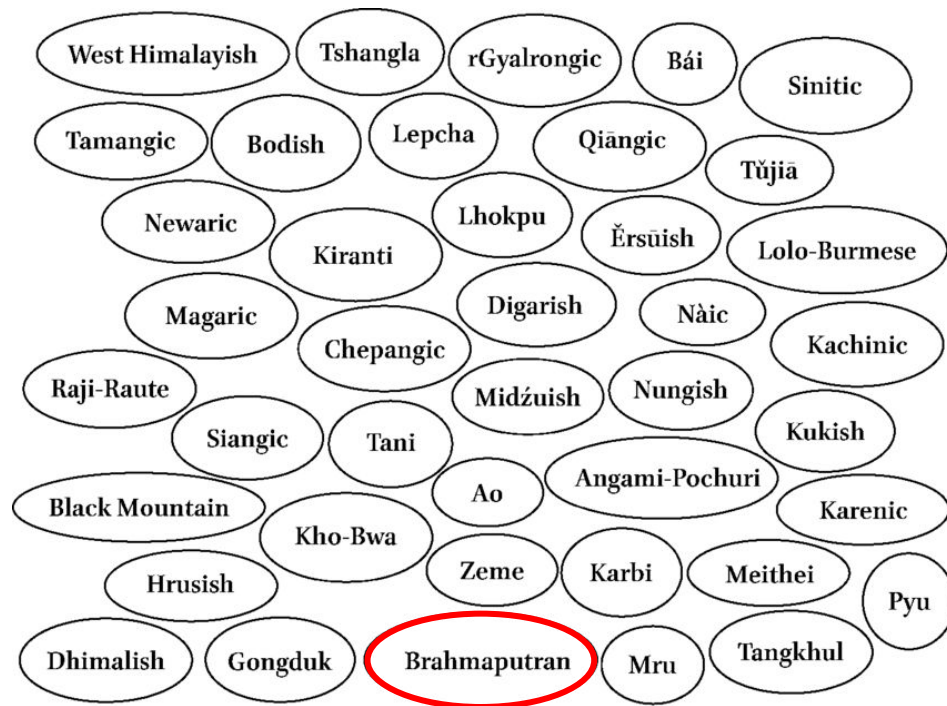


**See also:** Orlandi, G., 2021.

# Sino-Tibetan Phylogenies

Recent Sino-Tibetan phylogenies offer interesting suggestions of large-scale relationships in the family.

However, significant differences in results and approach still leave us reliant on the "Fallen Leaves" model of van Driem (2012).

One solution: shifting to a intensive **bottom-up** approach to resolve some of these issues.



**See also:** Orlandi, G., 2021.

# The approach:

For each branch of Sal:

- collect data for all attested varieties / all published doculects based on the ~750 concept "SALIST" word list
- curate the data to account for biases / mistakes in elicitation, morphological features, semantic splits &c, omitting external borrowings and identifying internal borrowings
- produce trees following as closely as possible the methods of Sagart et al (2019)

# Stage 1 - Bodo-Garo

- *A Bayesian phylogeny of Bodo-Garo: Testing novel methods on established groupings.* North East Indian Linguistics Society conference. Guwahati, Assam, India. January 2023
- *Developing Bayesian language phylogenies from previously published data: A case study of Bodo-Garo.* Workshop on New Results and Methods in Reconstructing Population History. Universität Zürich, Zürich, Switzerland. 30 January – 1 February 2023
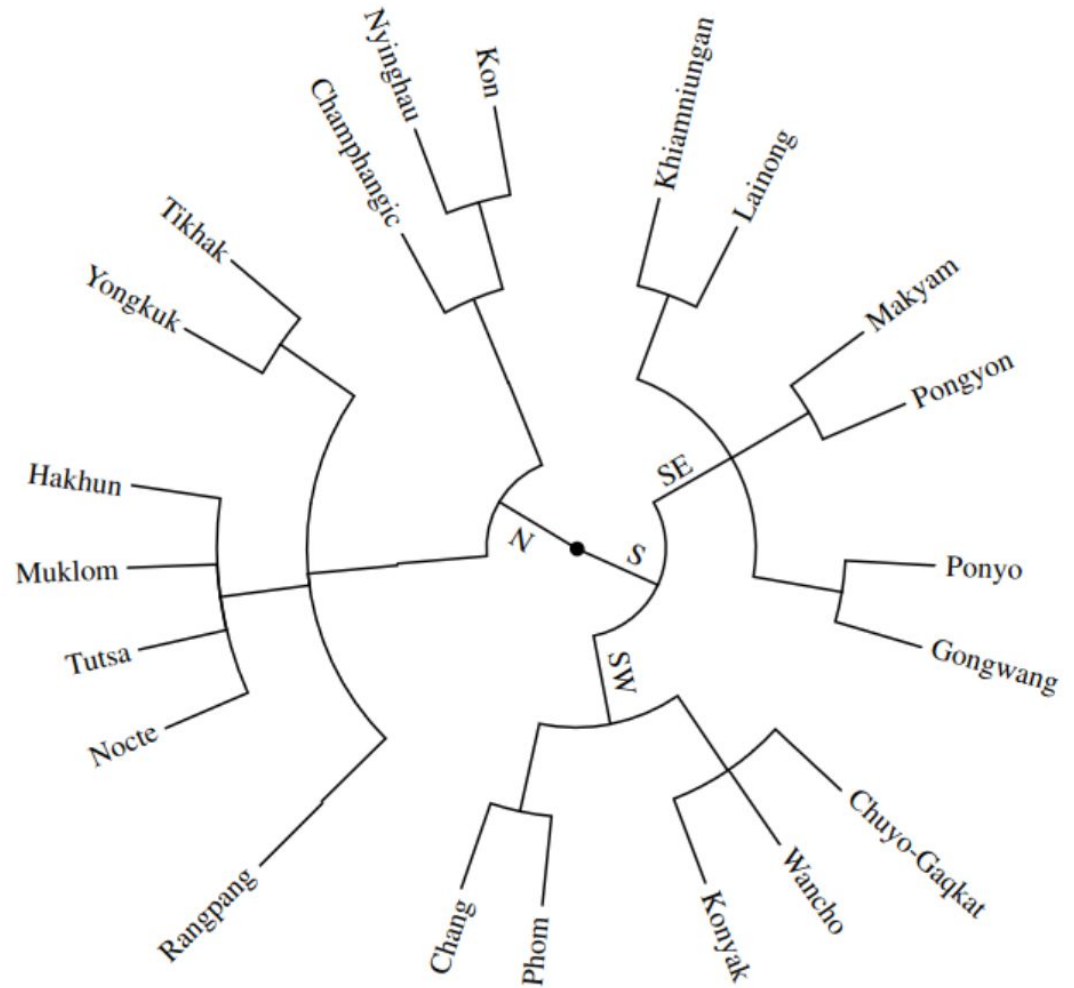
**Conclusion:**

Uncontrolled data sources / elicitation is an issue, but the nature of Bodo-Garo as a heavily creolised former lingua franca is a bigger issues.

# Stage 2 - Patkaian

Initial grouping based on speaker self-reporting & impressionistic descriptions in the literature.

Phonological reconstructions done for each major parent node, checked against neighbouring varieties or those for which some other connection may become apparent.

Through these reconstructions, regular correspondences have been established based on which cognacy can be judged.

# **Patkaian** (Northern Naga)

Not closely related to the other "Naga" languages (Angami, Ao, Sumi etc)

Typical community size is around 2'000 speakers for most varieties, some are much larger: 60'000 each for Khiamniungan & Wancho

The most diverse branch within the proposed Sal family.

# data collection & methodology
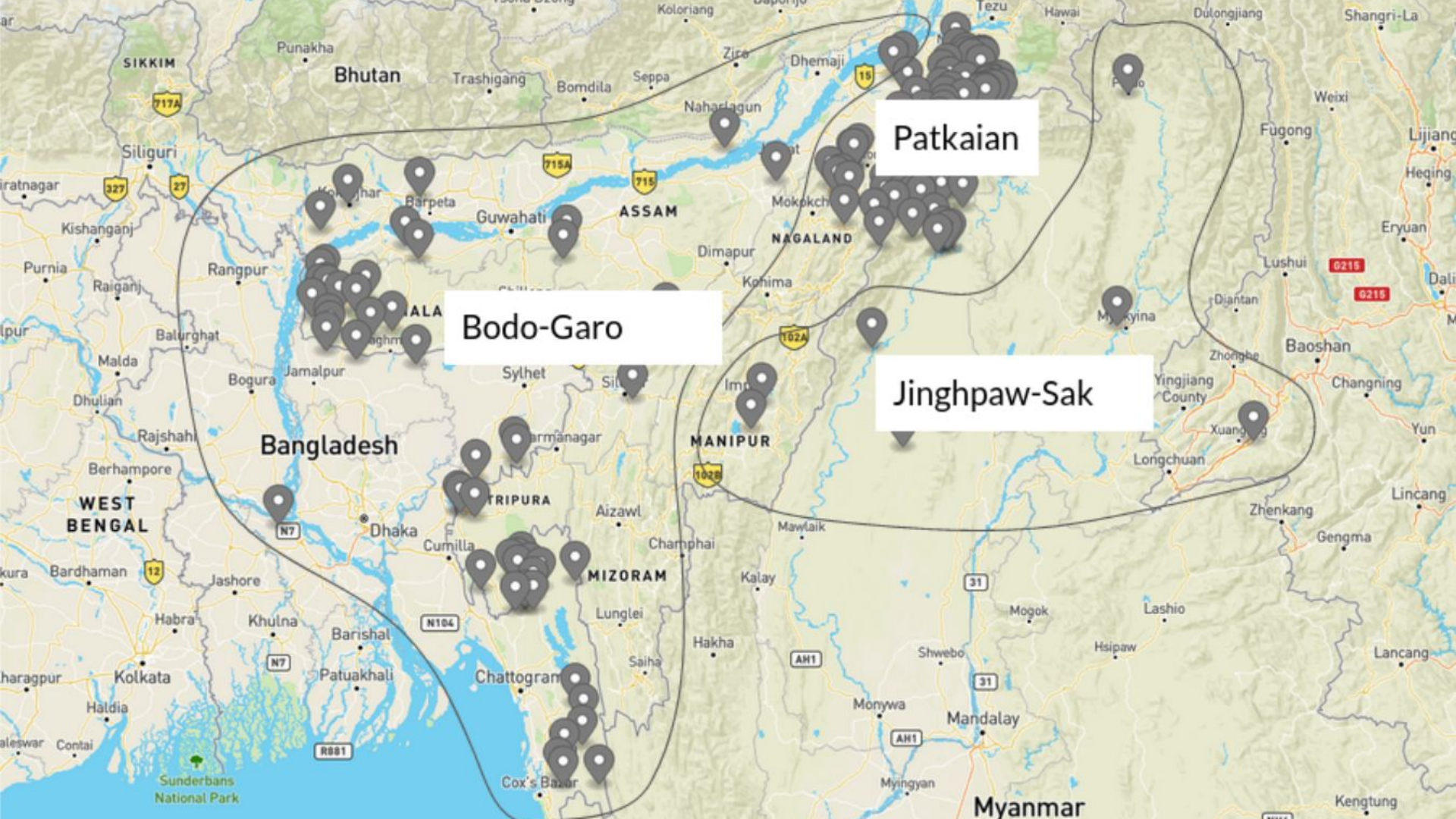
# Data collection

- **~750 concepts** covering **~175 doculects** of Patkaian
- Concepts derived from the CALMSEA word list (Matisoff 1978), the Intercontinental Dictionary Series word list (Key & Comrie 2023), and a large number of additional concepts which are widely attested in the doculects, forming the **SALIST** (Sal Area Lexical Inventory for Sino-Tibetan)

| concept_id | name | definition | hindi | assamese | mandarin | burmese |
|---|---|---|---|---|---|---|
| alcohol-brewed | brewed alcohol | Alcoholic beverages made through fermentation, such as beer or wine | सुराही दार शराब | জলকীয় দ্ৰব | 酿造酒 | အရက်ချက် |
| alcohol-distilled | distilled alcohol | Alcoholic beverages made through distillation, such as whiskey or vodka | अर्क | মদ | 蒸馏酒 | ပေါင်းခံအရက် |
| alive-living | to be alive | To have life and be living, not dead | जीवित | জিৱন থাকিব | 活着 | အသက်ရှင်ဖို့ |
| all | all | The whole quantity or extent of something; everyone or everything considered together. | सब | সমগ্ৰ | 所有 | အားလုံး |
| amber-glass | amber | A hard, translucent fossilized resin, typically yellowish-brown in color | अंबर | বাঁশফুলীয়া | 琥珀 | ပယင်း |
| | | A person from whom one is descended, | | | | |

# Data collection

- Only those terms which represent the most basic / typical word for a concept are included
- mini sketch grammars have been developed for each language regarding morphology, nominalisation, affixes in general
- Concepts with low coverage across branches were omitted in the end (< 8)

# Cognacy & mesolanguages

- Cognacy is enforced through regular sound correspondences.
    - However, irregular sound changes are common

        *$\mathbf{\gamma ap}$ → shoot, by extension kick/propel

        ⚹ $\mathbf{k^h i \eta_1}$ → *$\mathbf{k^h i \eta_1}$, of sky; *$\mathbf{k^h i \eta_3}$ → of water
- Such irregularities require some variability be allowed in cognate judgements

    re strictness of sound changes
- all terms have been coded for cognacy at the **morpheme** level

# Methodology

**LingPy** (List et al 2021) used for the creation of a MrBayes (Ronquist et al 2012) nexus file. Cognate assignment was still done manually.

Use of **MrBayes** Markov chain Monte Carlo method. No available means to calibrate a clock for Patkaian (or Sal more generally), so BEAST is not an option (as in Sagart et al 2019, Zhang et al 2019).

**Amri Karbi** (Konnerth, p.c; 2014) used as an outgroup, with additional Sal-internal outgroups for sub-branch trees (e.g. for determining three-way split within Patkaian)
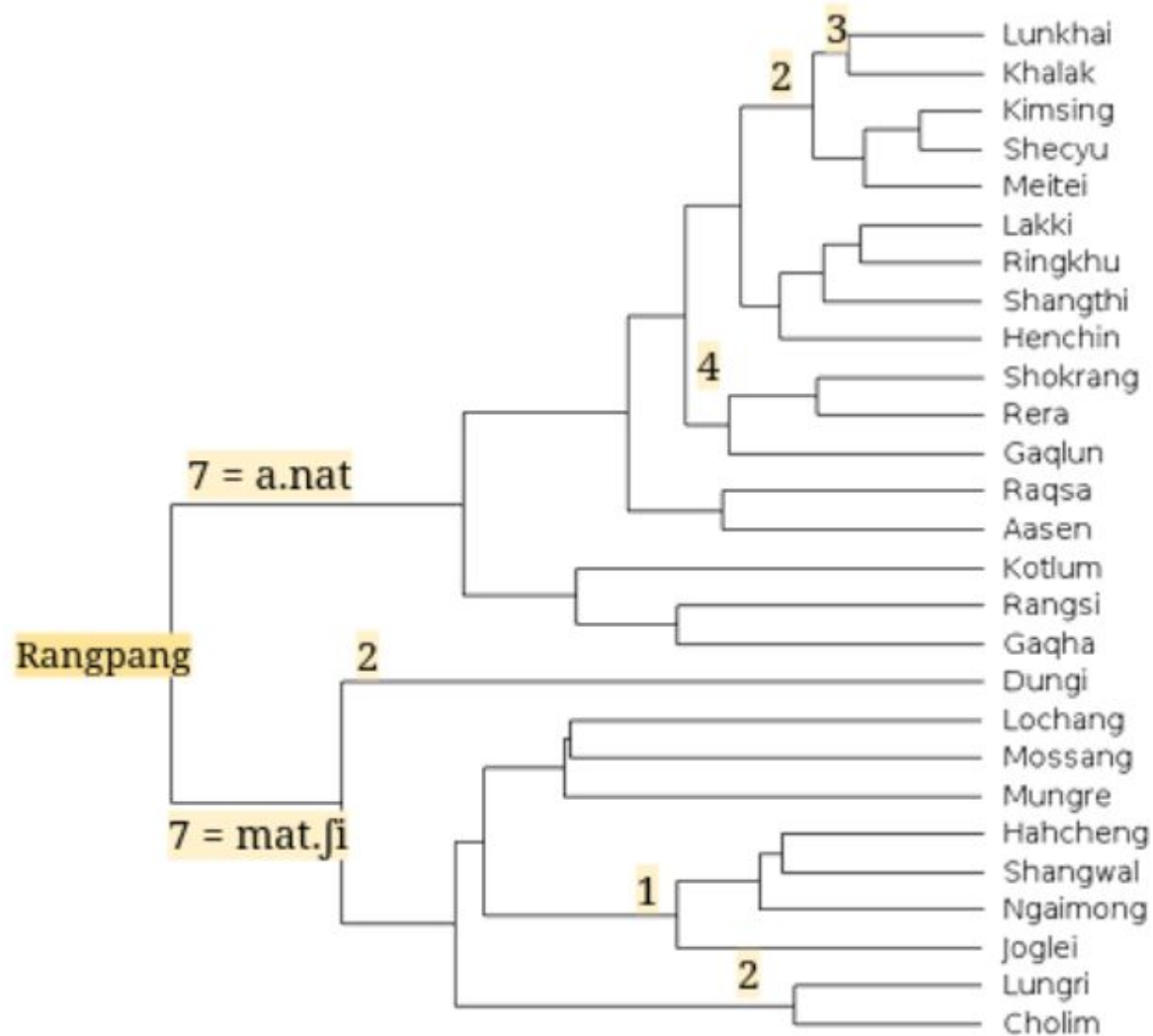
| concept | orthographic | phonetic | phonemic | full_segments | ipa | tokens | language_id | doculect_b | source |
|---|---|---|---|---|---|---|---|---|---|
| lick | | ɑdɑkdɑk | | a d a k + d a k | dak | d a k | KaisanNS | KaisanNS | statezni2021pc |
| lick | | tʰɑiɬoɬ | | tʰ a i | dai | ʒˡ a i | Karyaw | Karyaw | statezni2021pc |
| lick | | pʰwat kic | | pʰ u a t | pot | p o t | KhalakTK | Khalak | statezni2021pc |
| lick | ashi ao | | | a ʃ i | kʰoi | kʰ o i | KhiamNoklak | Khiamniungan | kumar1974khiam |
| lick | | su he | | s u + h e | kʰoi | kʰ o i | KhiamPasaung | Khiamniungan | statezni2021pc |
| lick | | ʃi | | ʃ i | kʰoi | kʰ o i | KhiamWolam | Khiamniungan | vandam2023wola |
| lick | | ɑɬ mɛɬ | | a m ɛ | mel | m e l | KimsingL | Kimsing | statezni2021pc |
| lick | | dək əɬ | | d ə k | ʒak | ʒ a k | KonChawang | Kon | statezni2021pc |
| lick | yai | | jai | j a i | lai | ʒˡ a i | KonyakM | KonyakTuensar | marrison1967cla |
| lick | yay; lay | | laj | l a j | lai | ʒˡ a i | KonyakN | KonyakWakchir | nagaraja1994kon |
| lick | | | jaj3 | j a j | lai | ʒˡ a i | KonyakTanhai | KonyakTanhai | jacques2010preli |
| lick | | ɑrɘkdɑyɩ | | r ɘ k | lik | l i k | KonYawngkon | Kon | statezni2021pc |
| lick | | lik sɘɬ m | | l i k | lik | l i k | Kotlum | Kotlum | statezni2021pc |
| lick | | jək ŋɑɬ | | j ə k | lik | l i k | Kotlum | Kotlum | statezni2021pc |
| lick | | ʔə mˡɛlɬ | | ə m j ə l | mel | m e l | KyahiP | Kyahi | statezni2021pc |
| lick | lɑipu | lɑi pu | | l a i | lai | ʒˡ a i | Kyan | Kyan | statezni2021pc |
| lick | kʰoiɬɑnɬ | | | kʰ o i | kʰoi | kʰ o i | LainongAnbaw | Lainong | statezni2021pc |
| lick | | kʰoiɬɑnɬ | | kʰ o i | kʰoi | kʰ o i | LainongHkamti | Lainong | statezni2021pc |
| lick | kʰoiɬ ɑnɬ | | | kʰ o i | kʰoi | kʰ o i | LainongHwiThaik | Lainong | statezni2021pc |
| lick | | kʰoiɬɑnɬ | | kʰ o i | kʰoi | kʰ o i | LainongLahe | Lainong | statezni2021pc |
| lick | | xoiɬɑnɬ | | x o i | kʰoi | kʰ o i | LainongLKNK | Lainong | statezni2021pc |

a mini-sketch grammar has been worked out for each variety covering word formation & basic morphology

# Results

|         | Varieties                                              | Contrastive Features                                                                                                                                                                                                      |
| ------- | ------------------------------------------------------ | ----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------- |
| Group 1 | Ngaimong, Joglei, Muishaung, Mungre, Maitai            | stop finals in past / negative; postverbal only marking in past, negative and (mostly) in the future                                                                                                                      |
| Group 2 | Cholim, Longri, Chamchang, Shecyü, Louchäng            | open finals carrying tone mostly 3 in past / negative (except some 3[rd] persons) preverbal mV- + open syllables, carrying tone 2, in the future                                                                           |
| Group 3 | Lungkhi, Khalak                                        | open finals in past (in k-), negative (in b-) and future (except some 3[rd] persons) no preverbal elements in combination with agreement marking                                                                           |
| Group 4 | Yvngban Wvng (Rangsi), Shangti, Gaqlun, Rinkhu, Rera   | perverbal marking in the negative, with postverbal agreement markers usually bare preverbal marking in the future in some varieties tone marking of agreement markers mostly tones 1 and 2                                  |

| English | Ngai-mong | Mui-shaung | Mungre | Louchäng | Cham-chang | Shecyü | Cholim | Rinkhu | Song Language |
|---|---|---|---|---|---|---|---|---|---|
| | 1 | 1 | 1 | 1 | 2 | 2 | ? | stop | stop |
| blow | $mul_1$ | $\text{əmui}_1$ | $moj_1$ | $maɯ_1$ | $meɪ_2$ | $me_2$ | mɔ? | ($p^hɯt$) | |
| fall | $dəl_1$ | $dəi_1$ | $daj_1$ | $de_1$ | $dɛə_2$ | $dia_2$ | djɤ? | dit | dət |
| ill | $ða_1$ | $ʈuu_1$ | $tsa_1$ | $di_1$ | $tsi_2$ | $dzi_2$ | de? | rak | |
| cloth | $k^həl_1$ | $k^həi_1$ | $k^haj_1$ | $khe_1$ | $k^hɛə_2$ | $khia_2$ | $k^hjɤ?$ | $k^het$ | $k^hət$ |
| trample | $na_1$ | $nɯɯ_1$ | $na_1$ | | $ŋi_2$ / $ni_2$ | $ni_2$ | ne? | | nak |
| hear | $tal_1$ | $tai_1$ | $təj_1$ | $ti_1$ | $təi_2$ | $tai_2$ | $te_1$ | (i)tat | tat |
| open up | dəp | $dau_1$ | $dəj_1$ | | $di_2$ | $di_2$ | $de_1$ | | dep |
| fear | $hil_1$ | $hi_1$ | $xaj_1$ | $hai_1$ | $hai_2$ | $hai_2$ | hjɤ? | ($p^hap$) | |

Tree diagram with the following tip labels (top to bottom):
Lunkhai, Khalak, Kimsing, Shecyu, Meitei, Lakki, Ringkhu, Shangthi, Henchin, Shokrang, Rera, Gaqlun, Raqsa, Aasen, Kotlum, Rangsi, Gaqha, Dungi, Lochang, Mossang, Mungre, Hahcheng, Shangwal, Ngaimong, Joglei, Lungri, Cholim
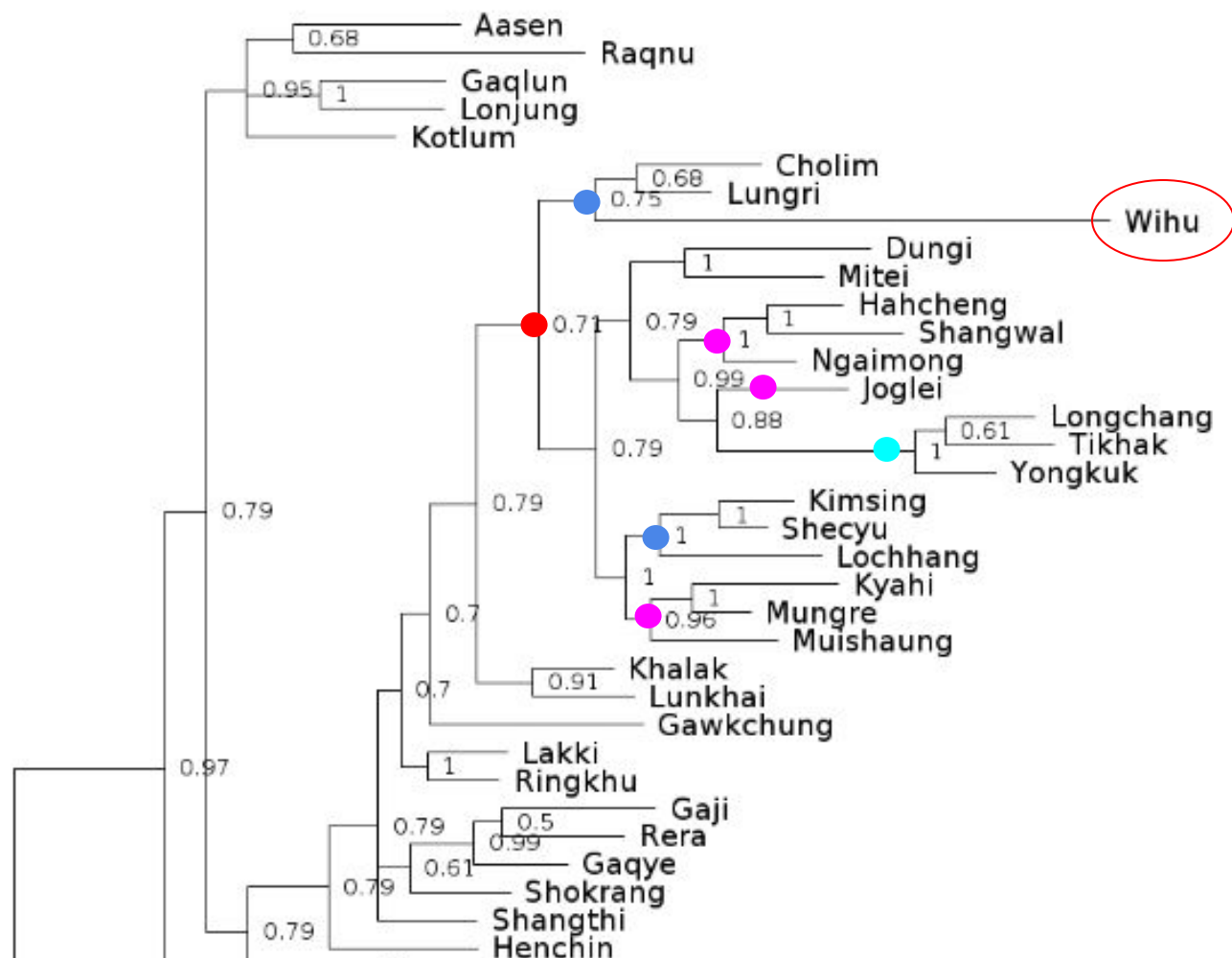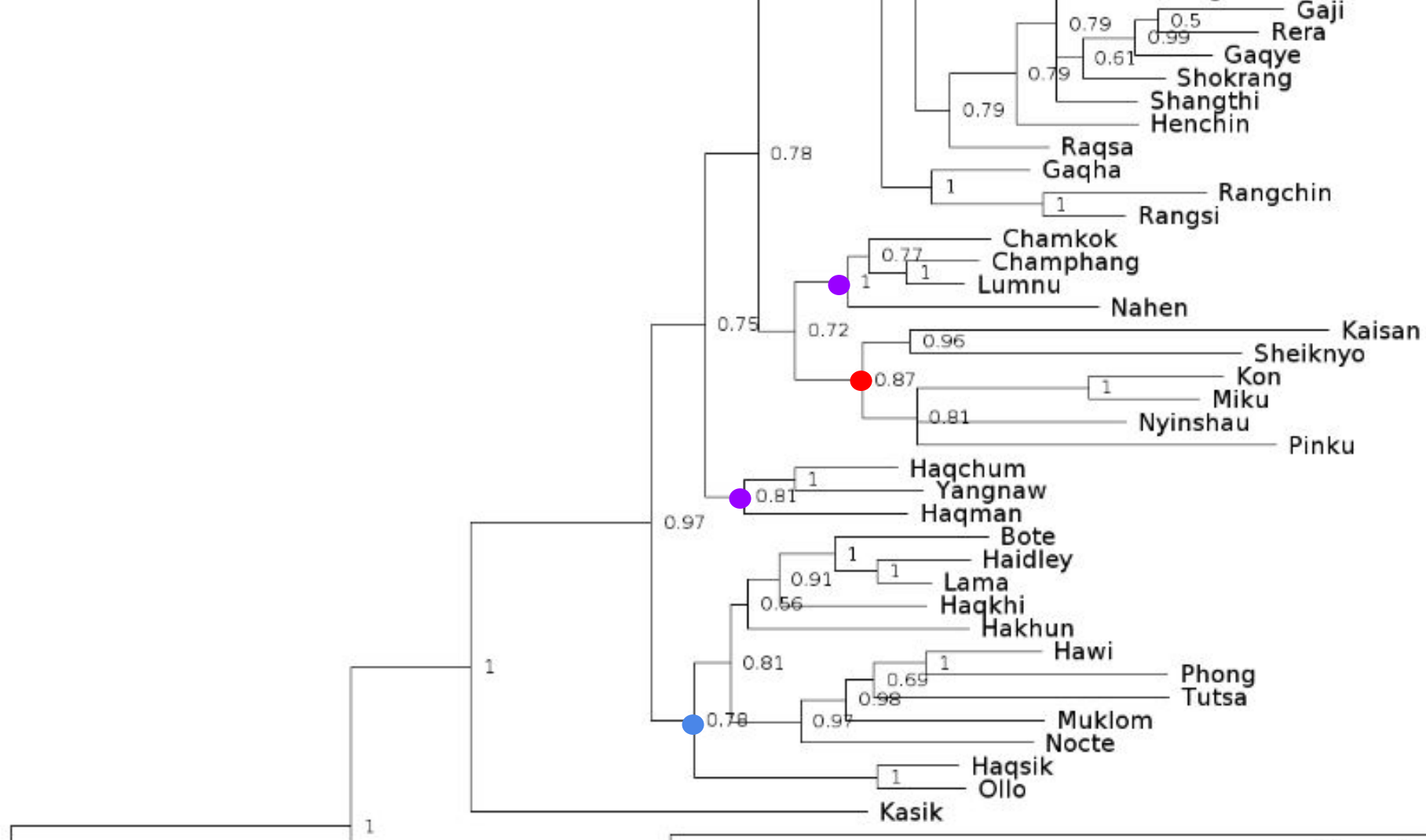
Node labels: 3, 2, 4, 7 = a.nat, Rangpang, 2, 7 = mat.ʃi, 1, 2

Morey's groups correspond to phonological & morphological features.

Many of these may have other explanations, e.g.:

- common sound changes for phonological features
- loss of productive *-ʔ nominaliser for tonal differences
- speculatively: Jesperson's cycle for negation placement, or IRR (vs FUT) < NEG, or just well attested innovation in TB negation marking (DeLancey 2015)
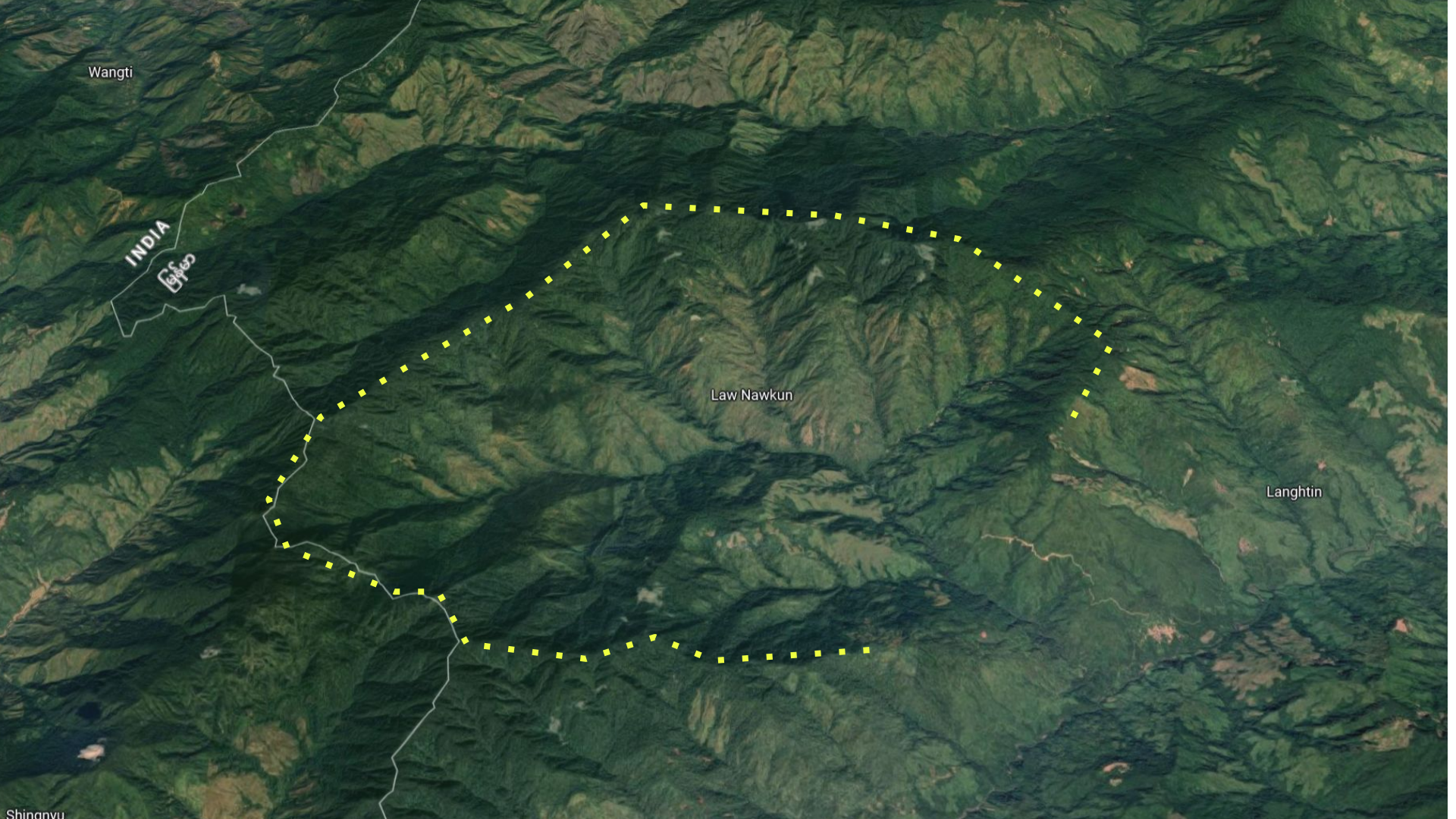- esoterogeny!

etc.

Kuku
Makyam
Santhong
1
1

Aasen
Raqnu
0.68
Gaqlun
Lonjung
0.95  1
Kotlum

Cholim
Lungri
0.68
0.75
Wihu

Dungi
Mitei
1
Hahcheng
Shangwal
0.79  1
Ngaimong
1
Joglei
0.99
0.71
0.88
Longchang
Tikhak
0.61
Yongkuk
1
0.79

Kimsing
Shecyu
1
Lochhang
1
Kyahi
Mungre
1
Muishaung
0.96

Khalak
Lunkhai
0.91
Gawkchung
0.79
0.7
0.7

Lakki
Ringkhu
1
Gaji
Rera
0.5
Gaqye
0.99
Shokrang
0.61
0.79
Shangthi
Henchin
0.79
0.79
0.97

# Southern branches

Law and Ponyiu-Gongwan placement has
interesting complications

# **Southern branches** - Law

# Conclusions

# General conclusions

Many conventional groupings often do now hold up to scrutiny, rely more on administrative boundaries etc.

Often the result of modern (not historical) geographic proximity

Bayesian methods aren't always well applied, and may not always give the answers, but can often point us in important and unconsidered directions.

By combination of phonological **reconstruction**, elicitation **consistency**, incorporation of traditional migration **narratives** to reconstruct potential historical **contact** with particular focus on **geography**, and an understanding of word formation to control for elicitation inconsistencies, such methods can be of great value.

Much more care is required than is often taken.

# General conclusions

Bayesian methods aren't always well applied, and may not always give the answers even if done well.

The methods can still point us in directions we may otherwise miss, while also helping reduce **some** forms of researcher bias.

Ideally:

- phonological **reconstruction**, at all major branches
- consistent **elicitation** methods
- incorporation of migration **narratives** & historical **contact**
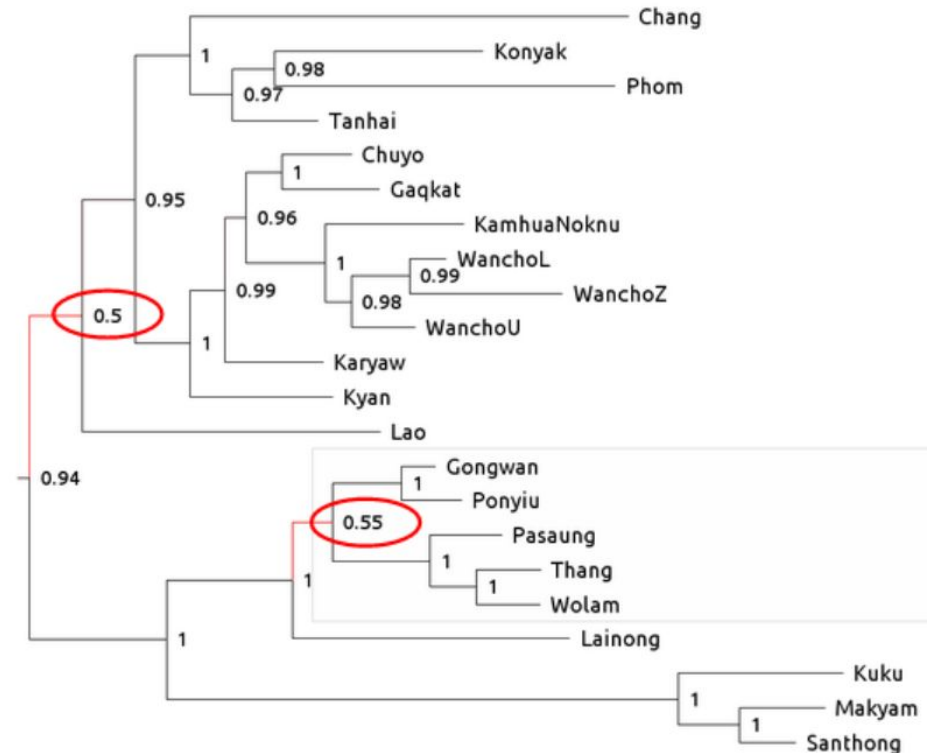- an accounting of **geography** factors

Much more care is required than is often taken.

# Specific conclusions

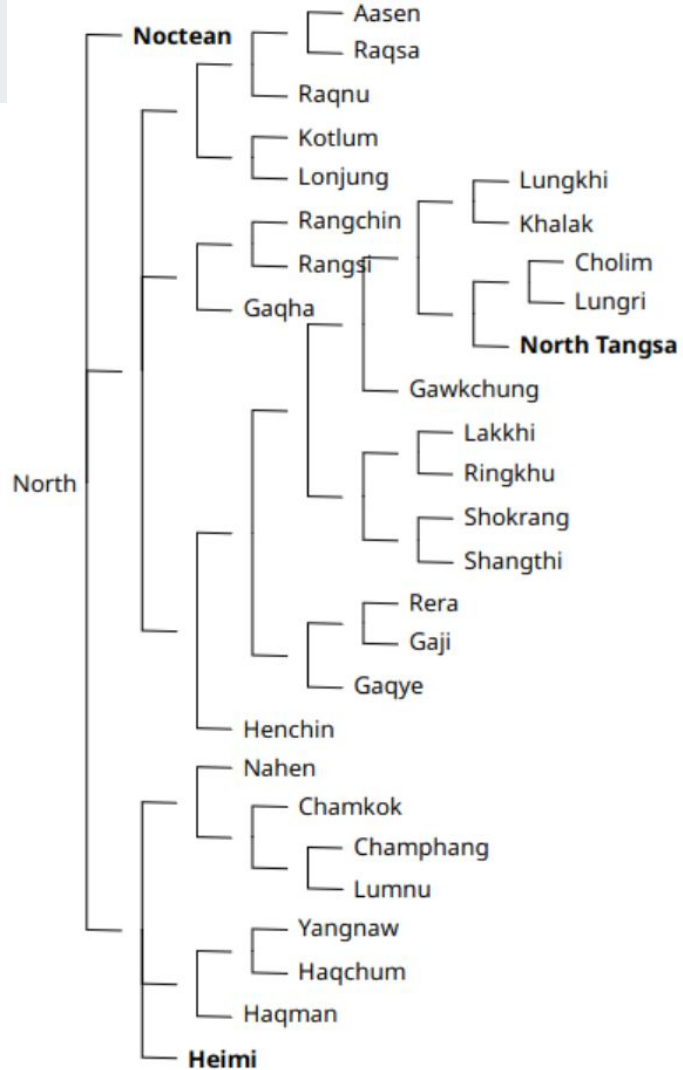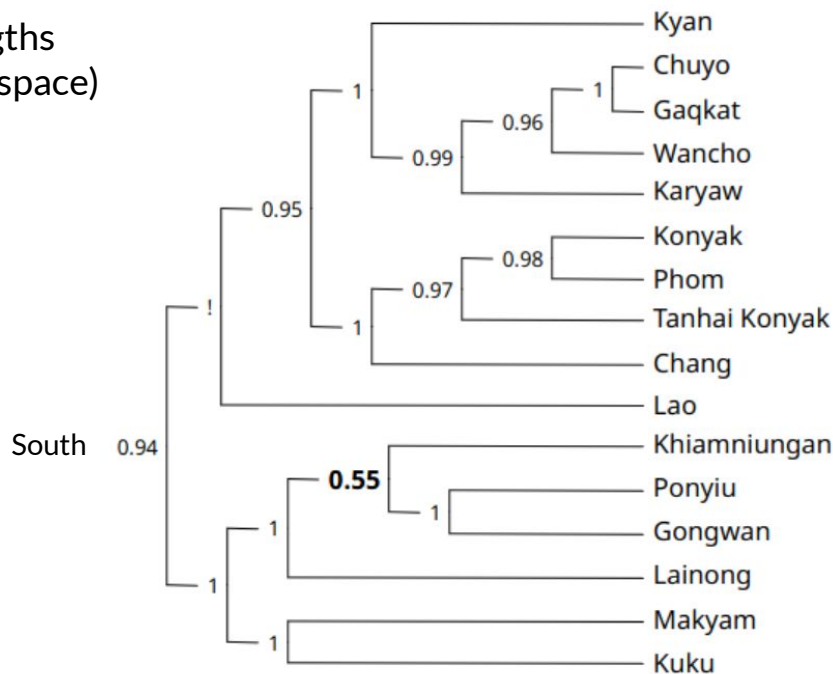"Noctean" as proposed by van Dam & Rahman (2019; 2021) is viable and includes Tutsa, Ollo, Hakhun & Muklom

Champhangic as including Haqchum is not viable

Khasik is properly placed in the northern branch, not Southwest, but as an earliest branch off Noctean

# The tree(s)
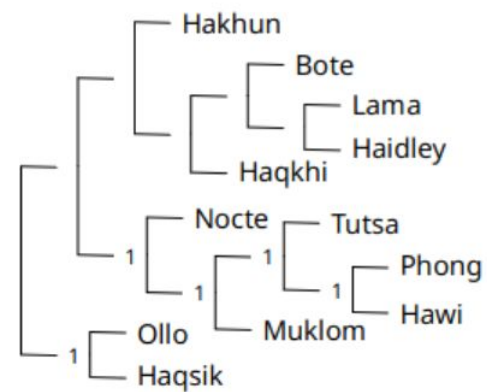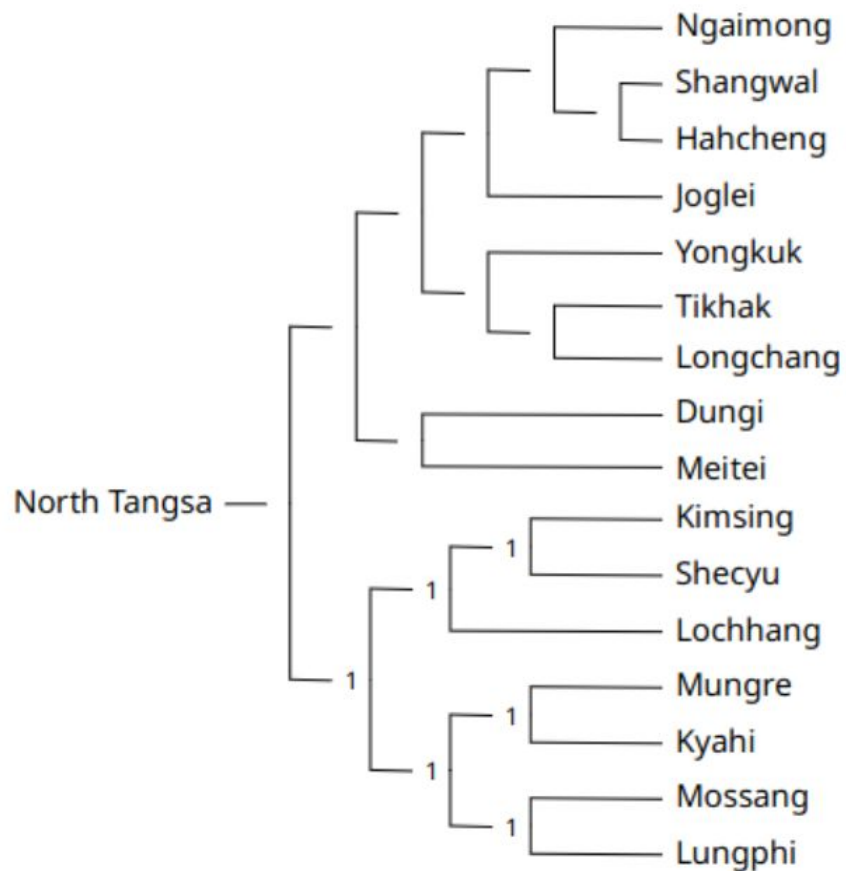
(branch lengths
omitted for space)

North Tangsa

Ngaimong
Shangwal
Hahcheng
Joglei
Yongkuk
Tikhak
Longchang
Dungi
Meitei
Kimsing — 1
Shecyu
Lochhang
Mungre — 1
Kyahi
Mossang — 1
Lungphi

Hakhun
Bote
Lama
Haidley
Haqkhi
Nocte — 1
Tutsa — 1
Phong — 1
Hawi
Muklom
Ollo — 1
Haqsik

Figure 3: Noctean

Nahen — 1
Champang — 1
Lumnu
Chamkok
Kaisan — 1
Sheiknyo
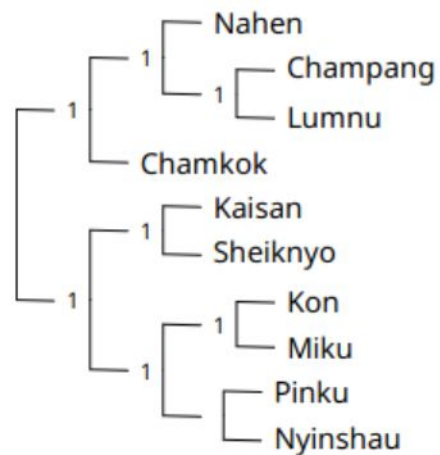Kon — 1
Miku
Pinku — 1
Nyinshau

Figure 4: Heimi

# Next steps

- Additional filtering by coverage

- Additional data collection for low-coverage varieties

- Collection of data for entirely missing varieties, including other 'liturgical' language doculects

*tʰaimi hai* (Wolam Khiamniungan)

*ketʒu əʒuŋ* (Muishaung Tangsa)

*ɲem pʰoi* (Kamhua Noknu Wancho)

**thank you** (West Michigan American English)

contact: **kellenparker@gmail.com**

Special thanks:
**Natalia Chousou-Polydouri** (UZH), **Dipjyoti Goswami** (La Trobe University),
**Keen Thaam** (Nagaland University), **Anui Sainu** (Myanmar), **Linda Konnerth** (Bern),
**Nathan Statezni** (SIL), **Stephen Morey** (La Trobe University)

# References

- DeLancey, Scott. 2015. The origins of postverbal negation in Kuki-Chin. North East Indian Linguistics 7, 203-212. Canberra, Australian National University: Asia-Pacific Linguistics Open Access.
- Konnerth, L.A., 2014. A grammar of Karbi. University of Oregon. // — 2023, personal communication
- List, Johann-Mattis and Forkel, Robert (2021): LingPy. A Python library for historical linguistics. Version 2.6.9. URL: https://lingpy.org, DOI: https://zenodo.org/badge/latestdoi/5137/lingpy/lingpy. With contributions by Greenhill, Simon, Tresoldi, Tiago, Christoph Rzymski, Gereon Kaiping, Steven Moran, Peter Bouda, Johannes Dellert, Taraka Rama, Frank Nagel. Leizpig: Max Planck Institute for Evolutionary Anthropology.
- Morey, S., 2019. Pangwa Tangsa agreement markers and verbal operators. Himalayan Linguistics, 18(1).
- Orlandi, G., 2021. Once again on the history and validity of the Sino-Tibetan bifurcate model. Journal of Language Relationship, 19(3-4), pp.263-292.
- Ronquist, F., M. Teslenko, P. van der Mark, D.L. Ayres, A. Darling, S. Höhna, B. Larget, L. Liu, M.A. Suchard, and J.P. Huelsenbeck. 2012. MRBAYES 3.2: Efficient Bayesian phylogenetic inference and model selection across a large model space. Syst. Biol. 61:539-542.
- Sagart, L., Jacques, G., Lai, Y., Ryder, R.J., Thouzeau, V., Greenhill, S.J. and List, J.M., 2019. Dated language phylogenies shed light on the ancestry of Sino-Tibetan. Proceedings of the National Academy of Sciences, 116(21), pp.10317-10322.
- Zhang, M., Yan, S., Pan, W. and Jin, L., 2019. Phylogenetic evidence for Sino-Tibetan origin in northern China in the Late Neolithic. Nature, 569(7754), pp.112-115.