# helayo: Reconstructing Sanskrit texts from manuscript witnesses

**Charles Li**[1,2]

**1** Centre nationale de la recherche scientifique **2** École des hautes études en sciences sociales

## Summary

For the vast majority of ancient and medieval texts, the original text itself is no longer extant in a material form. Instead, what we have are manuscripts which are copies of copies of copies, made over the course of hundreds or thousands of years, which accumulate errors and other changes each time they are copied by hand. In order to reconstruct the original text from these imperfect copies, scholars create a stemma — analogous to an evolutionary tree — in order to determine the relationships between manuscripts and to trace those textual changes over time.

## Statement of need

Due to the similarities in the methods used in the fields of textual reconstruction and evolutionary biology, textual scholars have begun to employ software created for biologists in order to analyze texts. Specifically, textual scholars are now using sequence alignment algorithms and phylogenetic tree-building packages to help with the reconstruction of ancient texts (Maas, 2013; Phillips-Rodriguez, 2007; Salemans, 2000). However, as bioinformatics becomes increasingly sophisticated, their models and algorithms have become more specific, and less applicable to non-biological sequences.

`helayo` has been designed from the ground up to perform multiple sequence alignment specifically for Sanskrit texts. Since Sanskrit has been written in over a dozen different scripts, each with their own orthographic peculiarities depending on their time and place, `helayo` performs a crucial pre-processing step in which the texts are normalized so that they can be compared meaningfully (Li, 2017). `helayo` can also tokenize texts either as individual characters or as *akṣaras*, since the Brahmic scripts used to write Sanskrit are abugidas, in which consonant and vowel pairs are written as a single unit.

In addition, a web-based `matrix editor` can be used to edit an alignment. It is also capable of automatically reconstructing a text based on an alignment and a phylogenetic tree, using the Fitch algorithm (Fitch, 1971). A full tutorial, with example files, is available at https://chchch.github.io/sanskrit-alignment/docs.

## Implementation

`helayo` is written in Haskell, and implements the Center Star multiple sequence alignment algorithm (Gusfield, 1997, pp. 347–350) with an affine gap penalty model (Li, 2021). It can be run in three different tokenization modes (character, akṣara, or whitespace-delimited word) and outputs a TEI XML file which can then be edited using the `matrix editor`.

37 The `matrix editor` is written in Javascript and can be used either online or offline. It
38 loads both TEI XML alignments produced by `helayo` as well as phylogenetic trees in NeXML
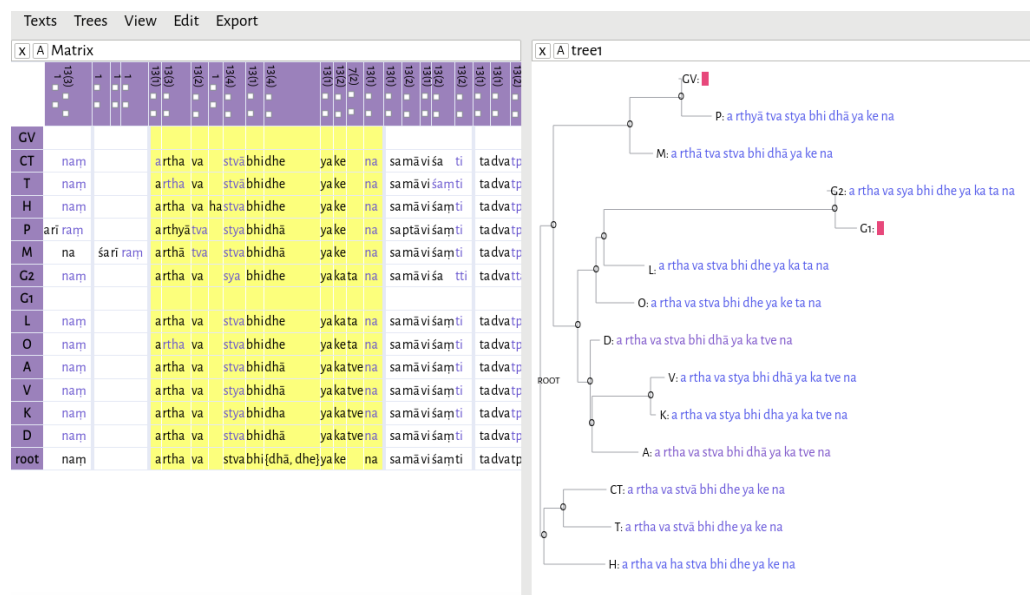39 format, which can be used together to reconstruct a text.



**Figure 1:** The matrix editor.

# References

41 Fitch, W. M. (1971). Defining the course of evolution: Minimum change for a specific tree
42 topology. *Systematic Zoology*, *20*, 406–316.

43 Gusfield, D. (1997). *Algorithms on strings, trees, and sequences*. Cambridge University Press.

44 Li, C. (2017). Critical diplomatic editing: Applying text-critical principles as algorithms. In
45 P. Boot, A. Cappellotto, W. Dillen, F. Fischer, A. Kelly, A. Mertgens, A.-M. Sichani, E.
46 Spadini, & D. van Hulle (Eds.), *Advances in digital scholarly editing. Papers presented at*
47 *the DiXiT conferences in The Hague, Cologne, and Antwerp* (pp. 305–310). Sidestone
48 Press.

49 Li, C. (2021). Align-affine: Sequence alignment with an affine gap penalty model. In *GitHub*
50 *repository*. GitHub. https://github.com/chchch/align-affine

51 Maas, P. A. (2013). On what to do with a stemma – towards a critical edition of the
52 Carakasaṃhitā Vimānasthāna 8.29. In D. Wujastyk, A. Cerulli, & K. Preisendanz (Eds.),
53 *Medical texts and manuscripts in indian cultural history*. Manohar.

54 Phillips-Rodriguez, W. J. (2007). *Electronic techniques of textual analysis and edition for*
55 *ancient texts: An exploration of the phylogeny of the Dyūtaparvan* [PhD thesis]. University
56 of Cambridge.

57 Salemans, B. J. P. (2000). *Building stemmas with the computer in a cladistic, Neo-*
58 *Lachmannian, way: The case of fourteen text versions of Lanseloet van Denemerken*
59 [PhD thesis]. Katholieke Universiteit Nijmegen.