

Continuous Space Representations of Linguistic Typology and their Application to Phylogenetic Inference

Yugo Murawaki (Kyushu University, Japan)

Background

Computational approaches to phylogenetic inference

- Huge success in the last decade
 - Indo-European
 - Austronesian
 - Bantu, etc
- Previous studies are cognate-based, using either
 1. regular sound changes, or
 2. the rates of birth & death of cognates
- Only applicable to known language families
 - Known because they share cognates!
- **Language isolates** lack cognates to be compared cross-linguistically
 - Ainu
 - Basque
 - Japanese

Linguistic typology as the last hope

- An arbitrary pair of languages can be compared

Japanese:	1	1	2	...	0	4
Korean:	2	2	1	...	0	4
Ainu:	0	1	1	...	0	4

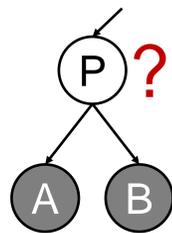
Feature 81A
Order of SOV

- 0: SOV
- 1: SVO
- 2: VSO
- ...

- Generally much more stable than cognates
 - Possibly even on the order of 10,000 years

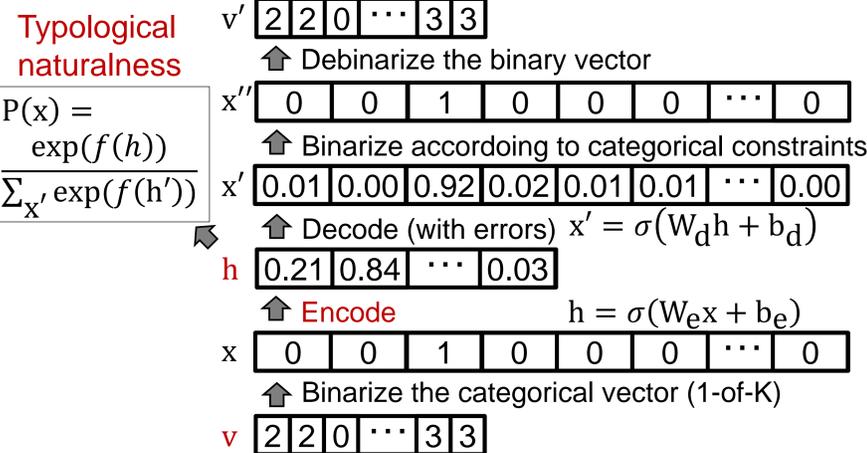
What can be said about the ancestor?

- We know too little about how typological features change over time
- The ancestor would be close to its descendants
- The ancestor must be a **typologically natural** language



- Previous phylogenetic models assume **independence** of features
- Typological studies show they are **interdependent**
 - Object-Verb implies Adjective-Noun [Greenberg, 1963]
- A naïve application of phylogenetic models leads to the reconstruction of unnatural ancestors

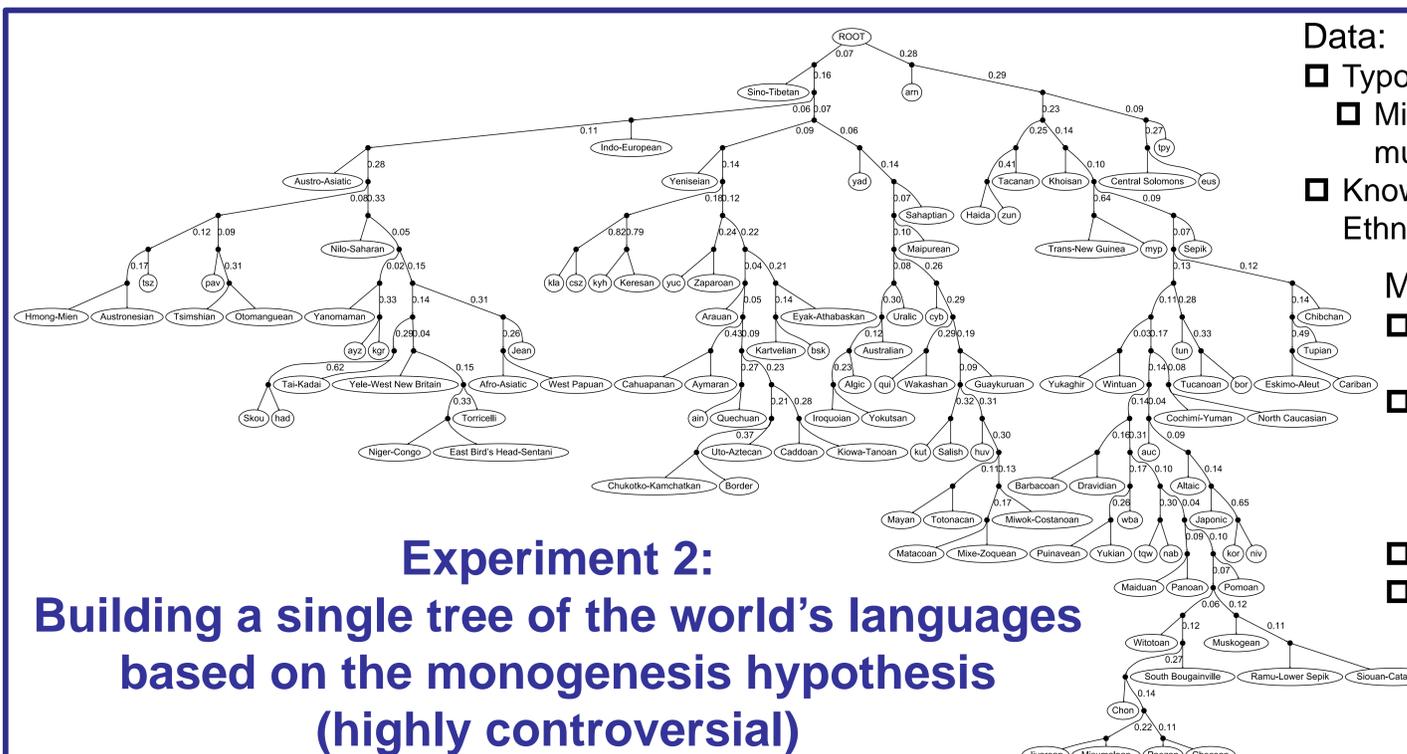
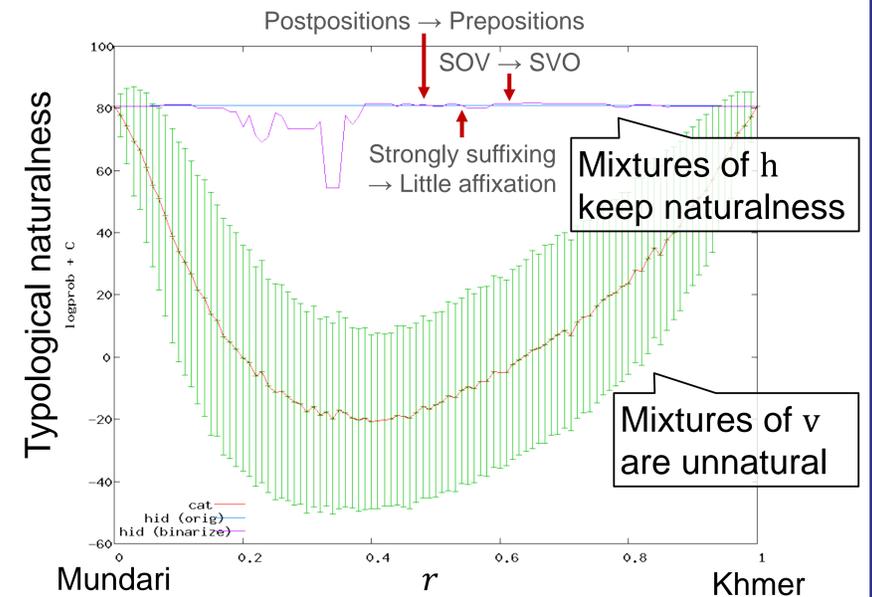
Continuous space representations that capture typological naturalness



- The Autoencoder non-linearly maps the feature vector into the continuous space
- The matrix W_e captures dependencies among features
- The decoder ensures that the original vector is reconstructible
- Energy-based model $P(x)$ assesses naturalness
 - Trained such that observed languages are distinguished from other possible combinations of features

Experiment 1: Mixing languages

- Assumption: the ancestor is close to an intermediate state between two descendants, a and b
- Compare two ways of mixing languages:
 1. Mixtures of h : linear interpolation in the continuous space $h = (1-r)h_a + rh_b$
 2. Mixtures of v : randomly replace the elements of the categorical vector v_a with those of v_b with prob. r



- Data:
- Typology Database: WALS
 - Missing values are imputed by multiple correspondence analysis
 - Known language families: Ethnologue

- Method:
- A simple generative tree model in the continuous space
 - Build a binary tree on top of known language families
 - Components' stability is learned from known trees
 - Run MCMC sampling
 - Summarize samples by a maximum clade credibility tree