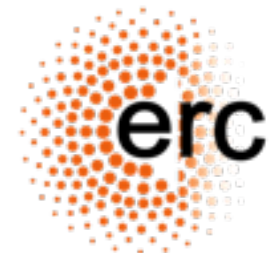


# Adding more Linguistics into Automated Reconstruction

Michael Cysouw



## The Problem:

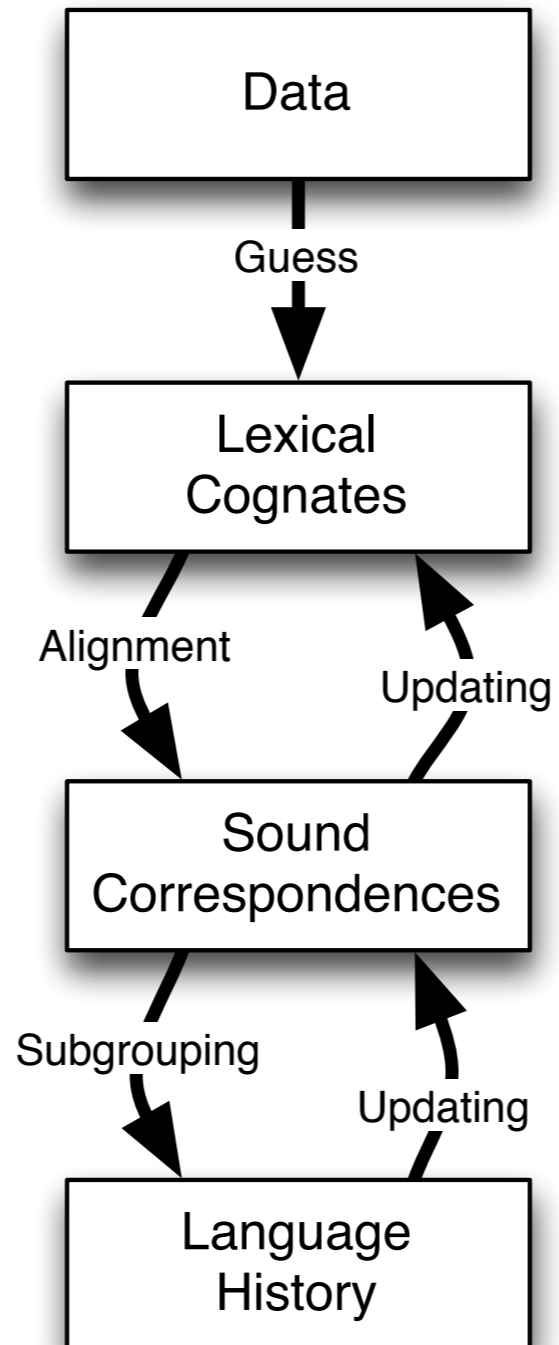
A method that only results in a tree of languages can lead to just three different responses from a linguistic specialist:

- the tree does not correspond to my research, *so the tree is wrong*
- the tree is the same as in my research, *so the method does not tell us anything new*
- I never investigated these languages, *so show me why this should be true*

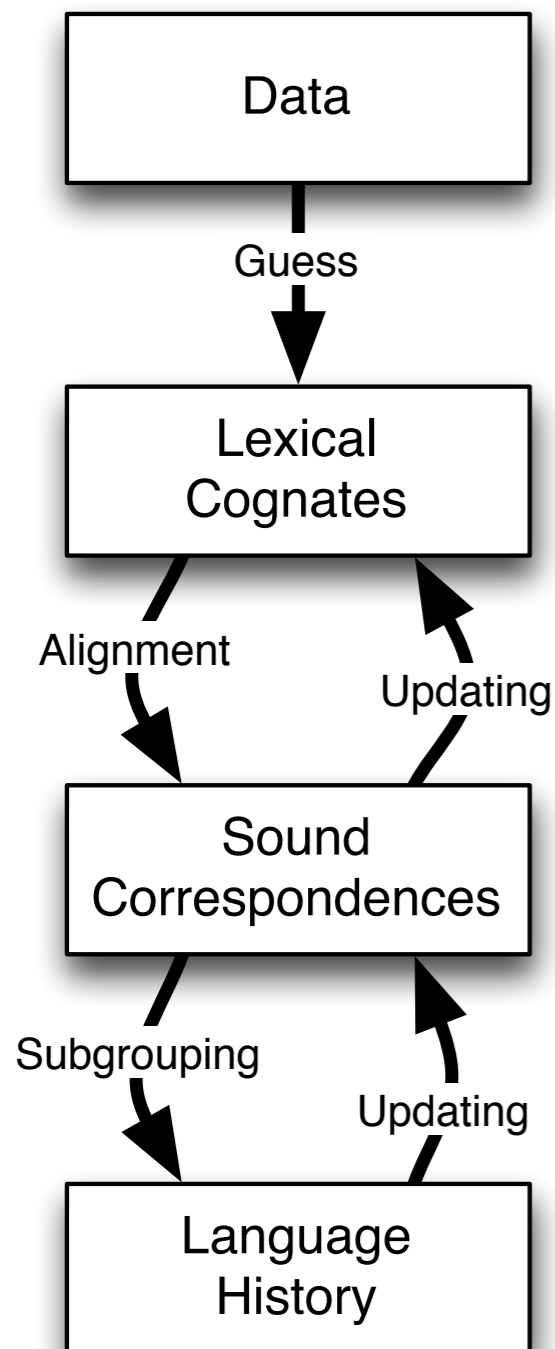
# Crucial steps of the linguistic *Comparative Method*

- **Regularity of Correspondences**  
 (“co-occurrence statistics”)
- **Correspondence Sets**  
 (“multiple alignment”)
- **Reconstruction**  
 (“summary of correspondences”)

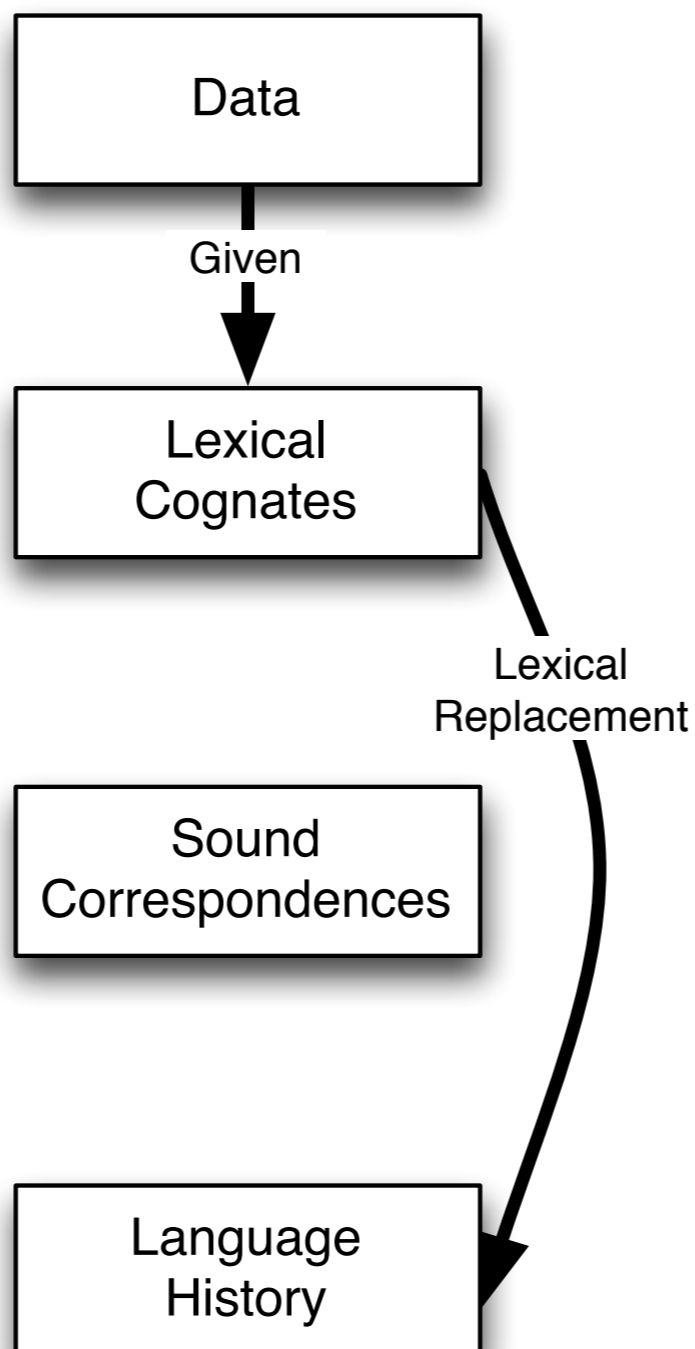
## Comparative Method



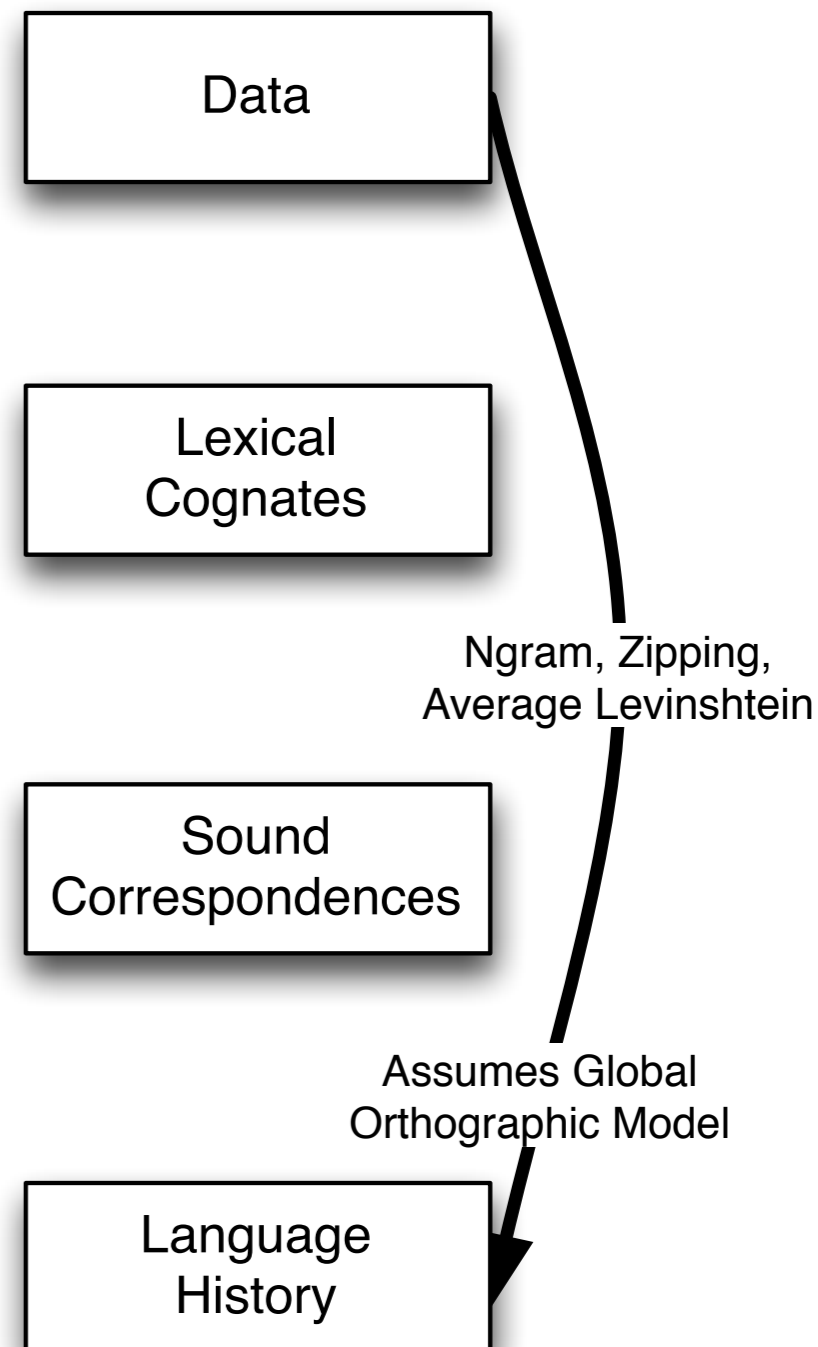
### Comparative Method



### Swadesh Method



### Black Box Method



# Swadesh Method

- **Necessary prerequisite:** known cognacy in parallel wordlist
- **Method:** Reconstruct history from distribution of cognates
- **Distances between language:** Swadesh (1952)
- **Stochastic model with first tree:** Sankoff (1969), Dobson (1969)
- **Modern replacement model:** Gray & Jordan (2000), (Gleason 1959)

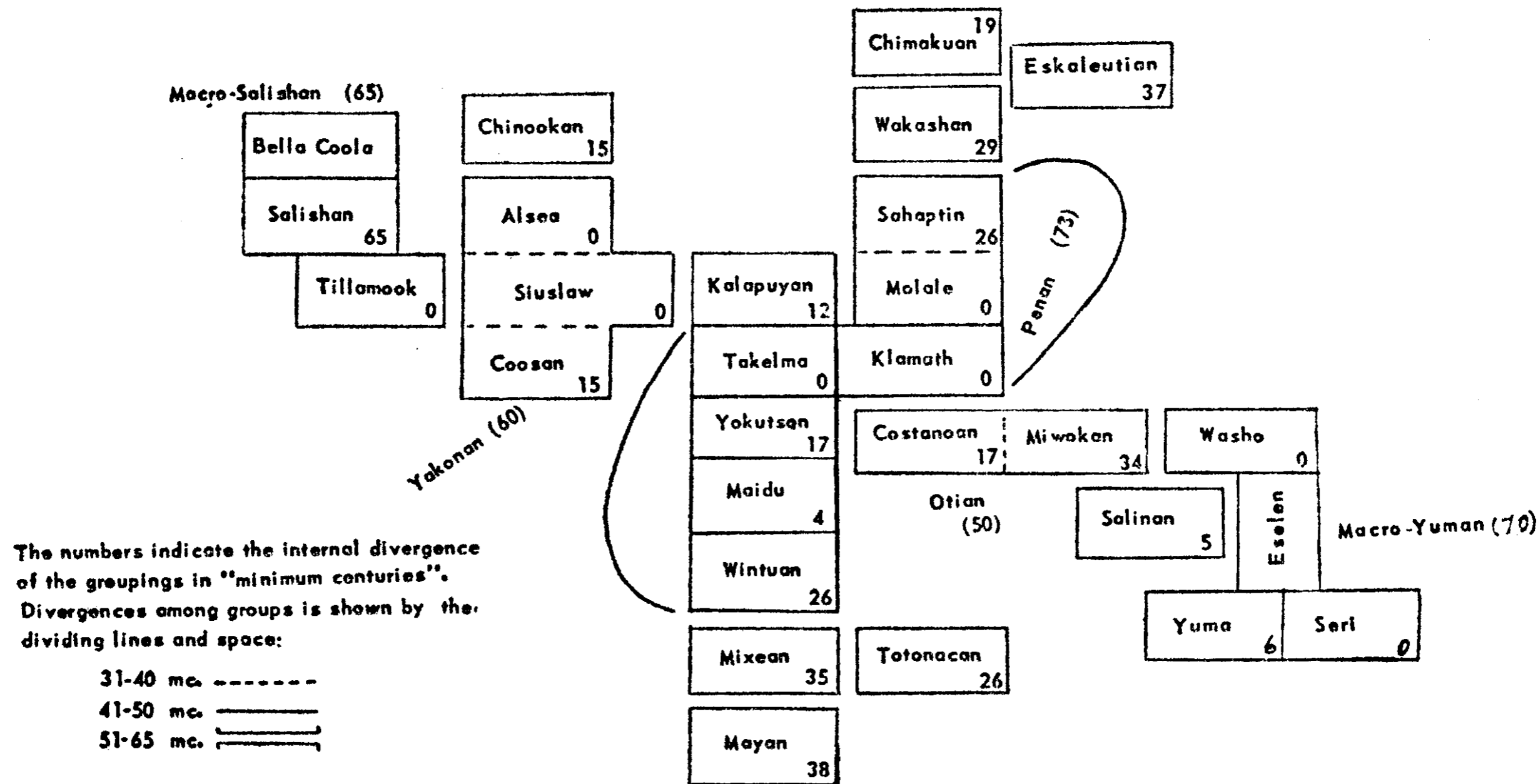


FIG. 2. Internal and external relationships of the Penan micro-phylum.

Swadesh, Morris. 1952. Lexico-Statistic Dating of Prehistoric Ethnic Contacts: With Special Reference to North American Indians and Eskimos. Proceedings of the American Philosophical Society 96(4). 452-463.

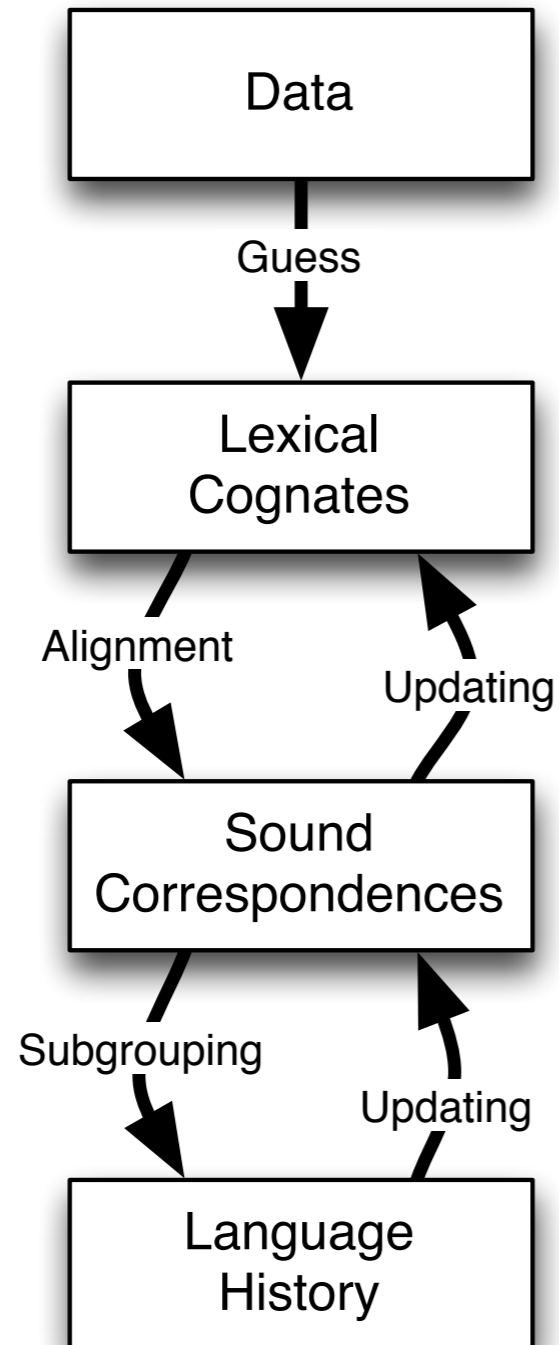
Swadesh, Morris. 1959, "Linguistics As An Instrument of Prehistory. Southwestern Journal of Anthropology 15, no. 1: 20-35.

# Black Box Methods

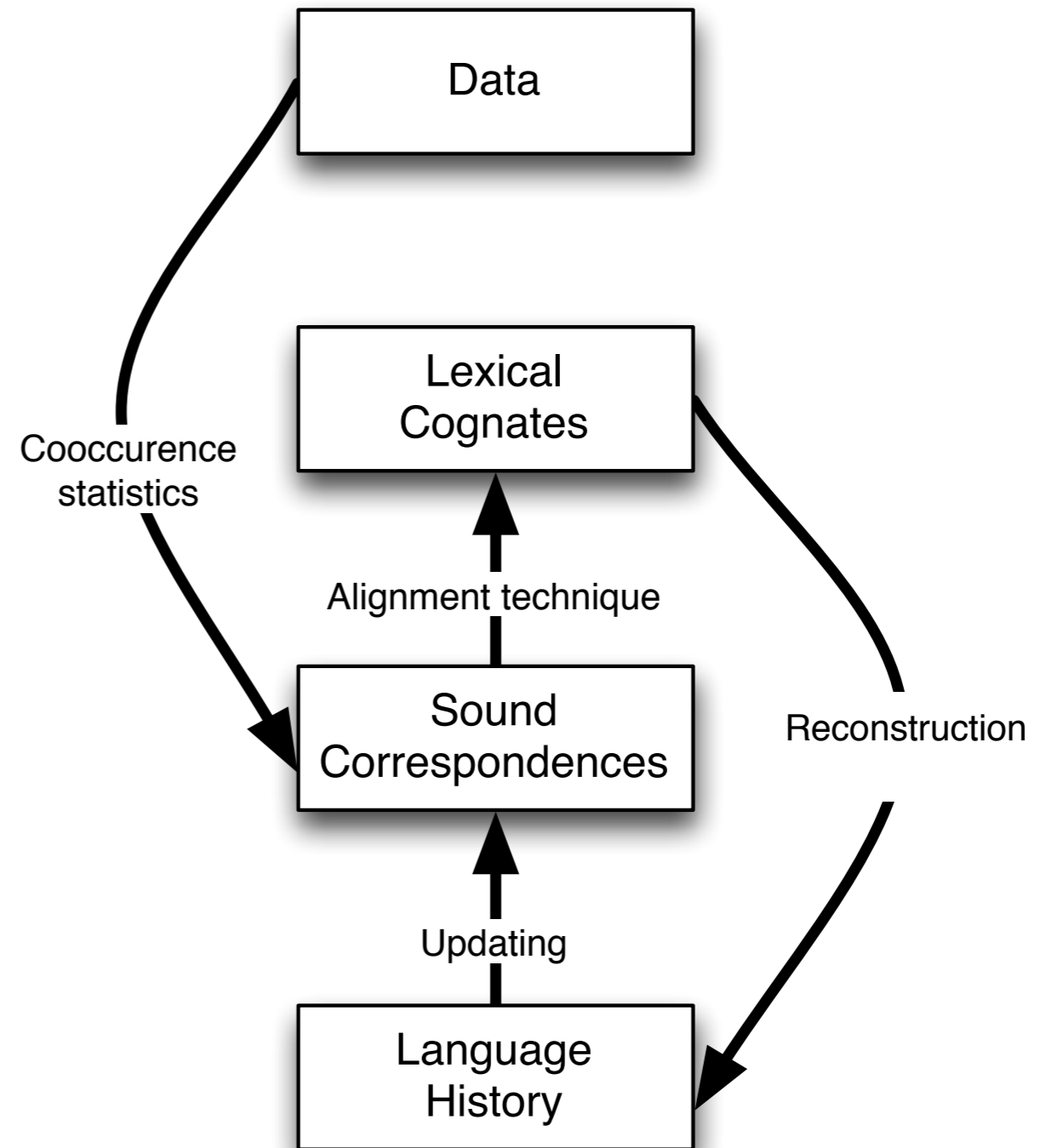
- **Necessary prerequisite:** Global orthographic model (“same alphabet for all data”)
- **Method:** Compare global orthographic similarity
- **Aggregate Levenshtein Distance:** Batagelj et al. (1992), Kessler (1995)
- **N-Gram Similarity:** Huffman (1998)
- **Ziping Distance:** Benedetto et al. (2002)



## Comparative Method



## Sound Change Method



# Regular symbol correspondences

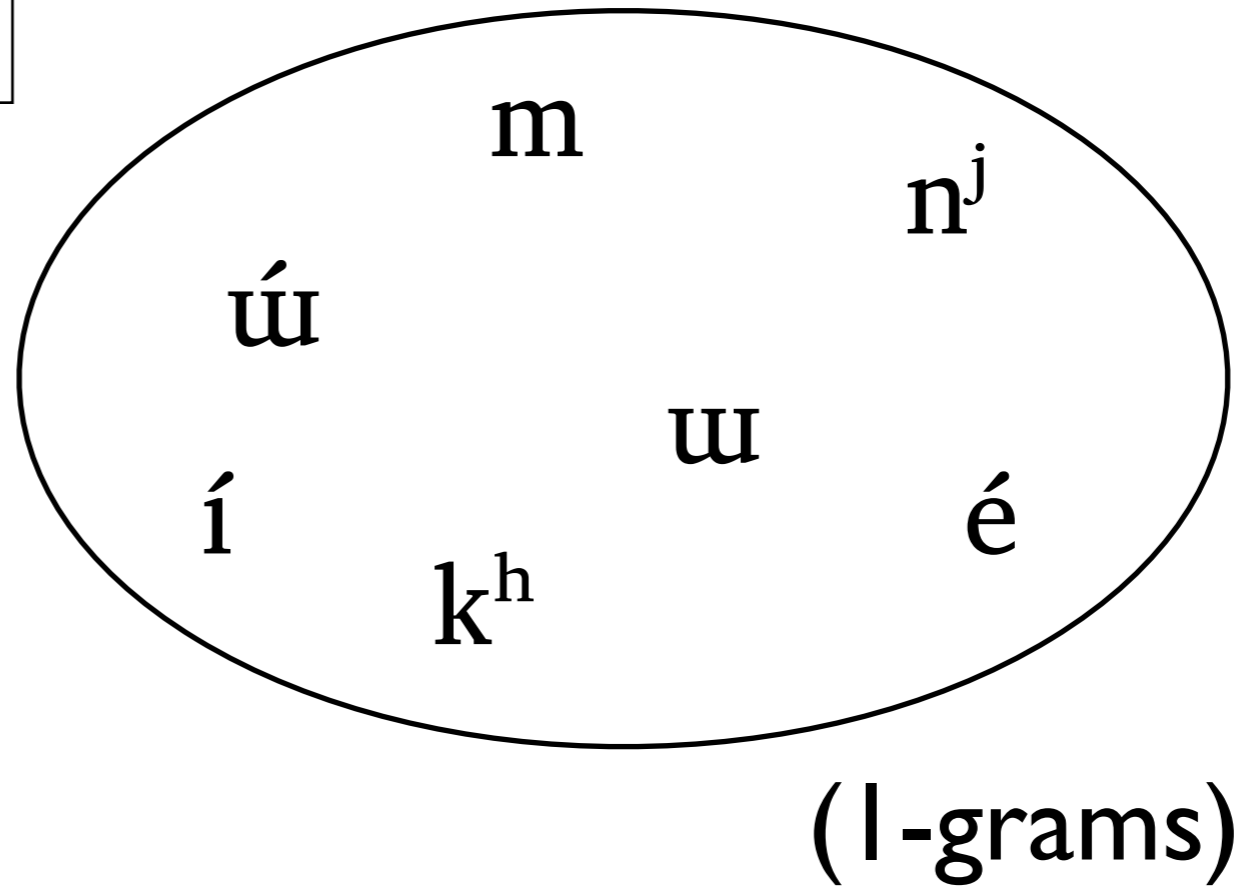
# Parts/Length Ratio

(“informativity of the parts”)

- **DNA**  $10^{-3}$  viz. 4 / 1000s  
(nucleotides vs. length of strings)
- **Protein**  $10^{-1}$  viz. 20 / 100s  
(amino acids vs. length of protein)
- **Words**  $10^{+1}$  viz. 40 / 6  
(sounds vs. length of word)
- **Sentences**  $10^{+2}$  viz. 1000s / 10  
(number of words vs. length of sentence)

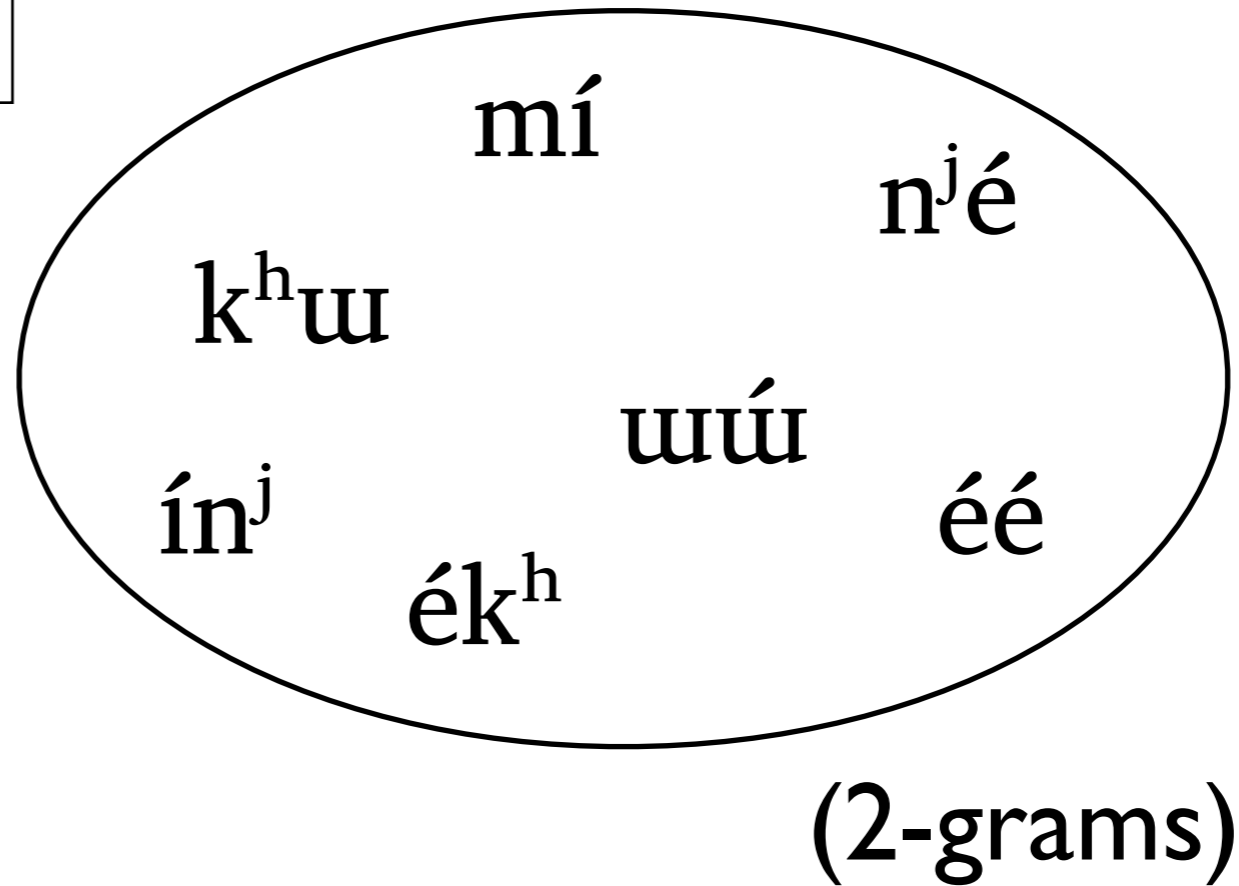
# ‘bag of symbol’ approach

mín<sup>j</sup>éék<sup>h</sup>uú



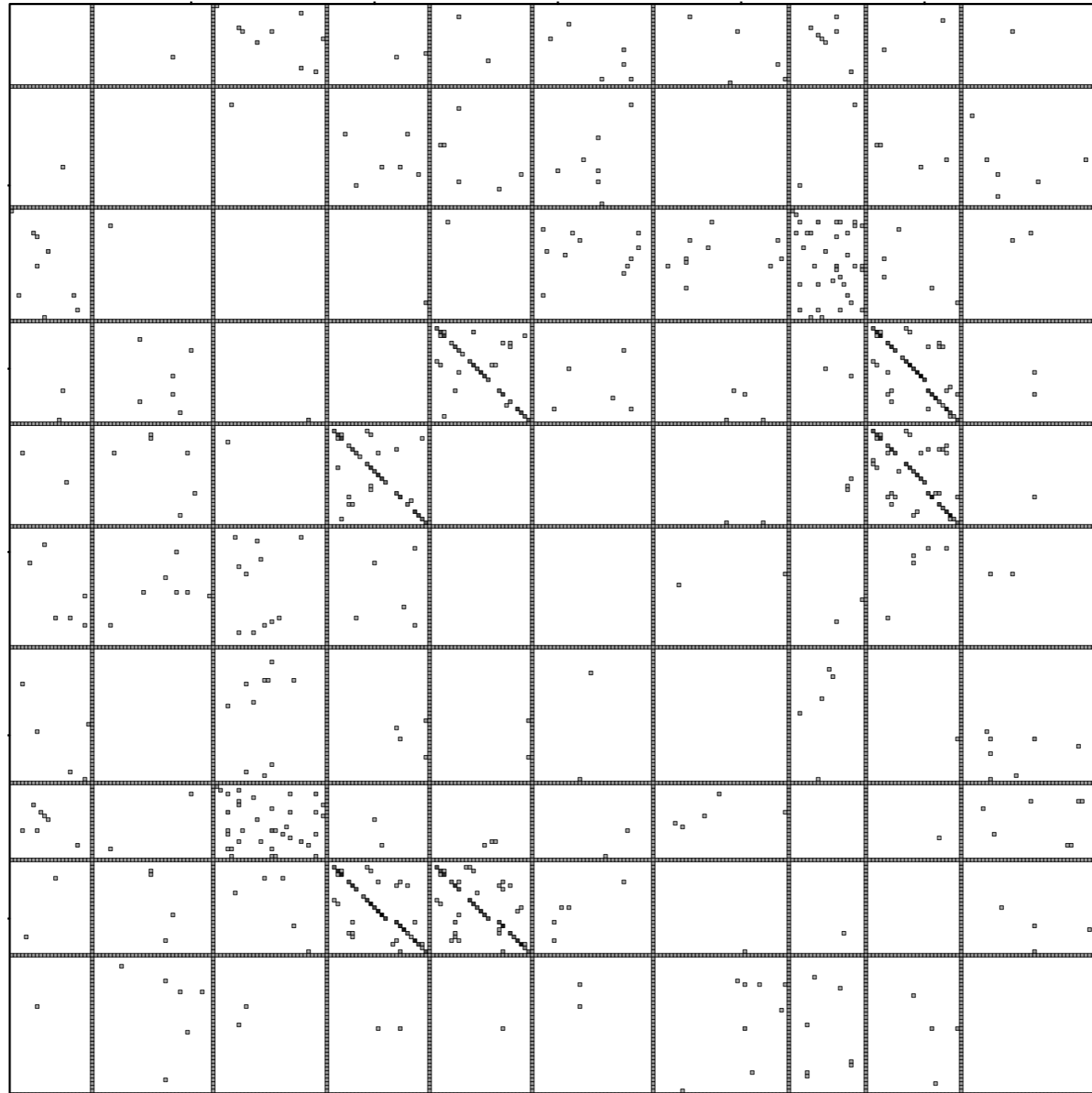
# ‘bag of symbol’ approach

mín<sup>j</sup>éék<sup>h</sup>wú



	<b>Bora</b>	<b>Muinane</b>
down	tʃin <sup>j</sup> e, paári	báari, gíino
bee	íimúʔóexp <sup>h</sup> i, téʔts <sup>h</sup> ipa	níibiri, míibiriʔi
sharp	ts <sup>h</sup> úʔxiβáne	síixéβano
...	...	...

Bora	Muinane	Bora	Muinane	Bora	Muinane
#k	#k <sup>h</sup>	#i	#i	#n	#n
kɨ	k <sup>h</sup> ɯ	#a	#a	#m	#m
se	ts <sup>h</sup> ɨ	di	ti	mɨ	mɯ
xe	xɨ	du	to	nɨ	nɯ
ga	k <sup>w</sup> a	#d	#t	us	ts <sup>h</sup> ɨ
ba	pa	#s	#ts <sup>h</sup>	#t	#t <sup>h</sup>
#b	#p	gi	tʃi	ɨg	ɯk <sup>w</sup>
e#	ɨ#	ni	ni	#ϕ	#p <sup>h</sup>



Pairwise symbol co-occurrences between 10 languages



# Correspondence Sets (“multiple Alignment”)

# tongue

American English	t <sup>h</sup>	ʌ	ŋ	-	-
Canadian English	t <sup>h</sup>	ʌ	ŋ	-	-
Central German (Cologne)	ts	ʊ	ŋ	-	-
Central German (Honigberg)	ts	aɪ	-	-	-
Central German (Luxembourg)	ts	ɔ	ŋ	-	-
Central German (Murrhardt)	ts	ʊ	ŋ	-	-
Danish	t <sup>h</sup>	ʊ	ŋ	-	ə
Dutch (Antwerp)	t	ʌ	ŋ	-	-
Belgian Dutch	t	ʊ	ŋ	-	-
Dutch (Limburg)	t	ʊ	ŋ	-	-
Dutch (Ostend)	t	ʊ	ŋ	-	ə
Dutch	t <sup>h</sup>	ʌ	ŋ	-	-
New Zealand English (Auckland)	t <sup>h</sup>	e	ŋ	-	-
English (Buckie)	t <sup>h</sup>	ʌ	ŋ	-	-
Indian English (Delhi)	t	e	ŋ	-	-
Nigerian English (Igbo)	t <sup>h</sup>	ʌ	ŋ	g	-
South African English (Johannisburg)	t <sup>h</sup>	ɜ	ŋ	-	-
English (Lindisfarne)	t	ɔ	ŋ	-	-
English (Liverpool)	t̪ <sup>h</sup>	ʊ	ŋ	g	-
English (London)	t <sup>h</sup>	e	ŋ	-	-
English (North Carolina)	t <sup>h</sup>	ɜ:	ŋ	-	-
Australian English (Perth)	t <sup>h</sup>	e	ŋ	-	-
English (Singapore)	t	ɑ	ŋ	-	-
English	t <sup>h</sup>	e	ŋ	-	-
English (Tyrone)	t <sup>h</sup>	ɔ	ŋ	-	-
Faroese	t <sup>h</sup>	ʊ	ŋ	k	a
German	ts	ʊ	ŋ	-	ə
High German (North Alsace)	ts	ʊ	ŋ	-	-
High German (Biel)	ts	ʊ	ŋ	-	ə
High German (Bodensee)	ts	ʊ	ŋ	-	ə
High German (Graubuenden)	ts	ʊ	ŋ	g	e
High German (Herrlisheim)	ts	ʊ	ŋ	-	-
High German (Ortisei)	ts	ʊ	ŋ	g	ɛ
High German (Tuebingen)	ts	u	ŋ	g	-
High German (Walser)	ts	ʊ	ŋ	g	ə
Icelandic	t <sup>h</sup>	ʊ	ŋ	k	a
Low German (Achterhoek)	t <sup>h</sup>	ʊ	ŋ	-	ə
Low German (Bargstedt)	t <sup>h</sup>	ʊ	ŋ	-	-
Norwegian (Stavanger)	t <sup>h</sup>	ʊ	ŋ	-	ə
Scottish	t <sup>h</sup>	ʌ	ŋ	-	-
Swedish (Skane)	t <sup>h</sup>	øʏ	ŋ	j	e
Swedish (Stockholm)	t <sup>h</sup>	ʊ	ŋ	-	ə
West Frisian (Grou)	t	ɔ	ŋ	-	ə
Yiddish (New York)	ts	u	ŋ	g	-

(LingPy library, maintained by Mattis List)

# foot

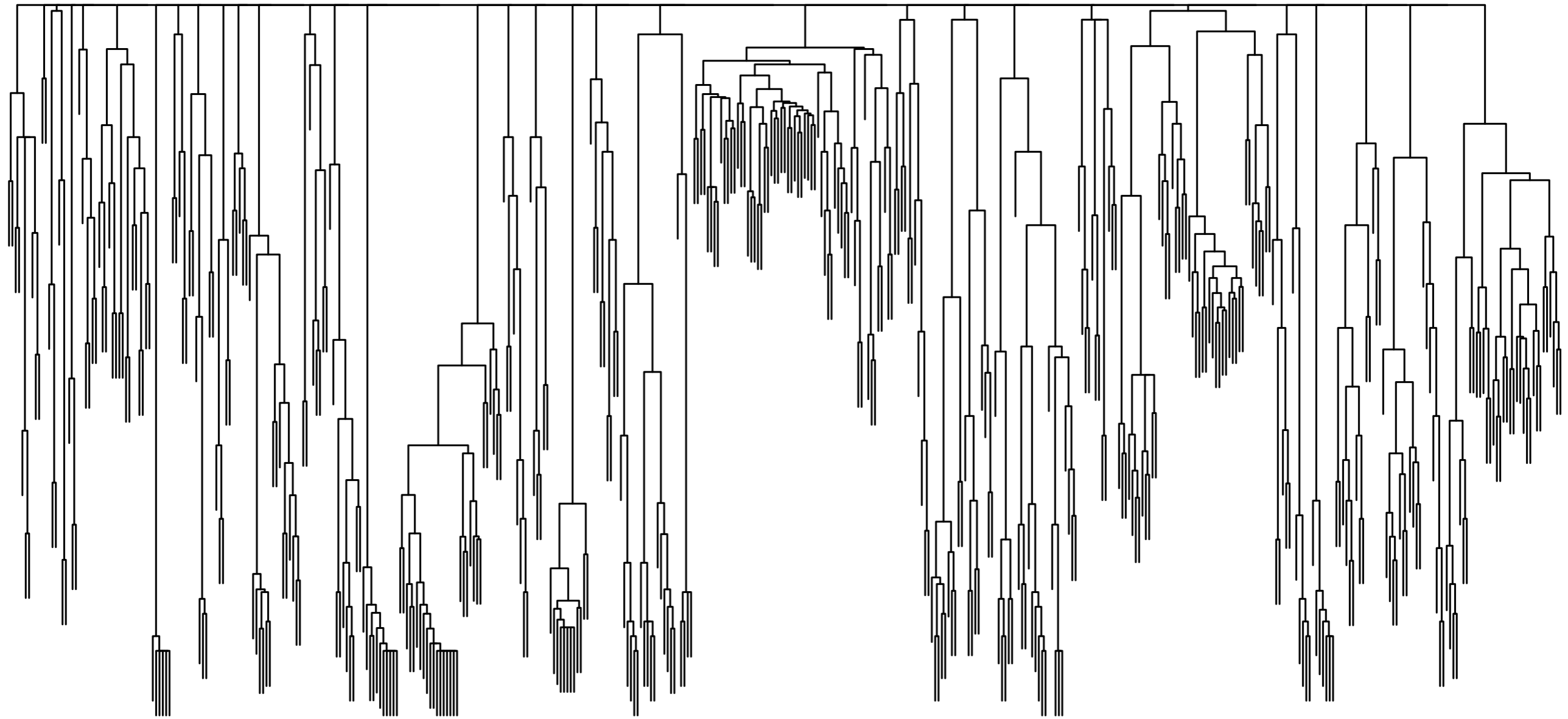
American English	f	ʊ	t
Canadian English	f	ʊ	t
Central German (Cologne)	f	ʊ	s
Central German (Hönigberg)	f	eu	s
Central German (Luxembourg)	f	œʊ	s
Central German (Murrhardt)	f	uə	s
Danish	f	ʊ	lʰ
Dutch (Antwerp)	f	u	t
Belgian Dutch	v	u	t
Dutch (Limburg)	f	oɪ	t
Dutch (Ostend)	v	uə	t
Dutch	f	u	t
New Zealand English (Auckland)	f	ʊ	t
English (Buckie)	f	æ	t <sup>h</sup>
Indian English (Delhi)	f	ʊ	t
Nigerian English (Igbo)	f	u	t
South African English (Johannesburg)	f	ʊ	t <sup>h</sup>
English (Lindisfarne)	f	ɔ	t
English (Liverpool)	f	ʊ	θ
English (London)	f	ʰ	ʔ
English (North Carolina)	f	øɪ	t
Australian English (Perth)	f	ʊ	t
English (Singapore)	f	u	t
English	f	ʊ	t <sup>h</sup>
English (Tyrone)	f	ɔ	t
Faroese	f	ðu	t
German	f	uɪ	s
High German (North Alsace)	ϕ	eɪ	s
High German (Biel)	f	oə	s
High German (Bodensee)	f	ʊɛ	s
High German (Graubünden)	f	ʊɔ	s
High German (Herrlisheim)	ϕ	eɪ	s
High German (Ortisei)	f	ʊə	s
High German (Tübingen)	f	uə	s
High German (Walser)	ϕ	oɛ	s
Icelandic	f	ɔʊ	t
Low German (Achterhoek)	β	oɪ	t
Low German (Bargstedt)	f	ɔɪ	t
Norwegian (Stavanger)	f	uɪ	t
Scottish	f	ʰ	t <sup>h</sup>
Swedish (Skane)	f	ʰu	t
Swedish (Stockholm)	f	o	tɪ
West Frisian (Grou)	f	oə	t
Yiddish (New York)	f	u	s

# Reconstruction

(“summary of correspondences”)

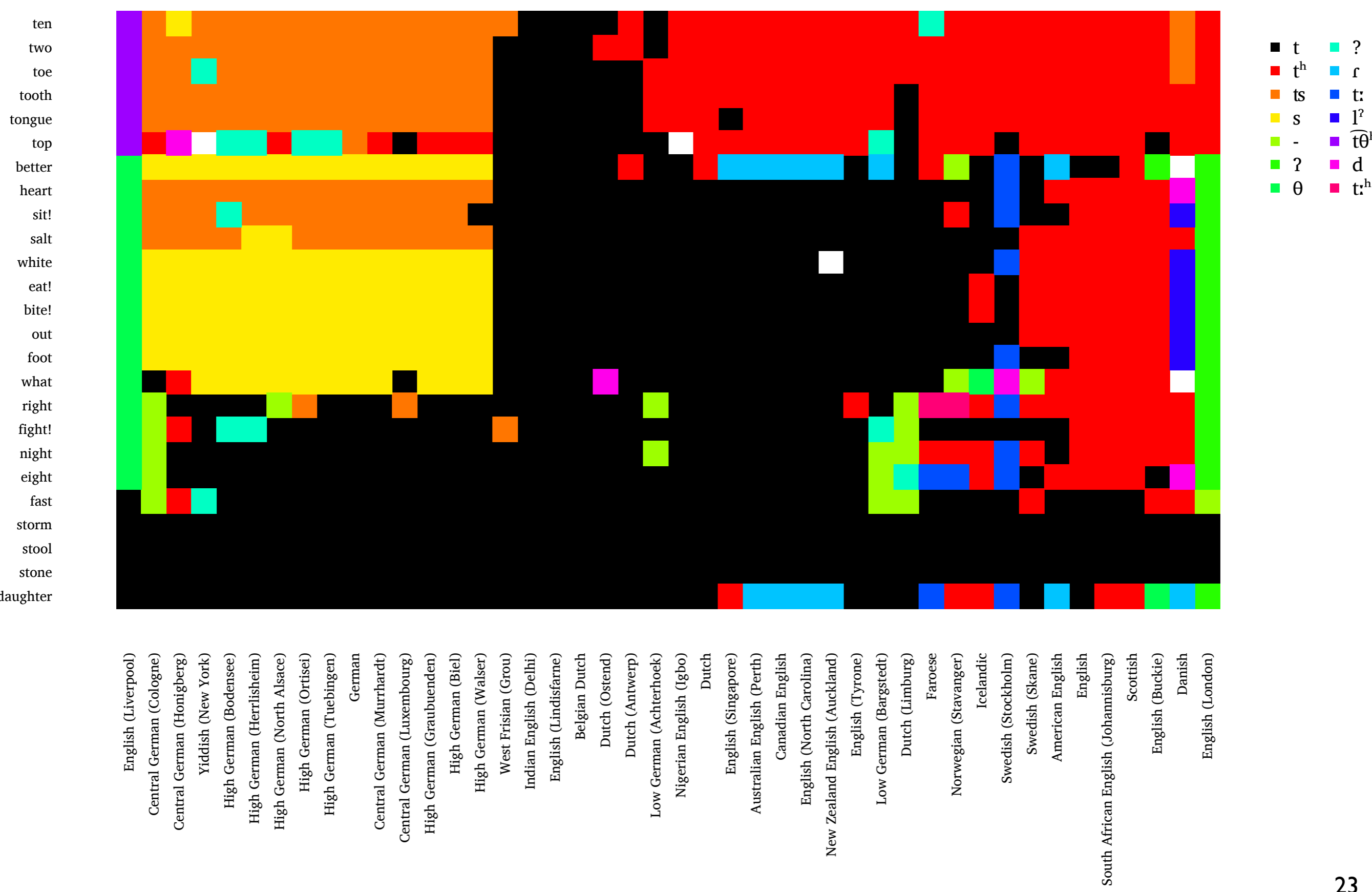
# Reconstruction

- What should be the (computational) representation of reconstructions ?
- Linguists mostly think of them as actual sounds
  - ▶ that is a highly problematic point-estimate of the actual variation
- Proposal: a proto-sound is a cluster of correspondence sets
  - ▶ clusters of “columns in an alignment”



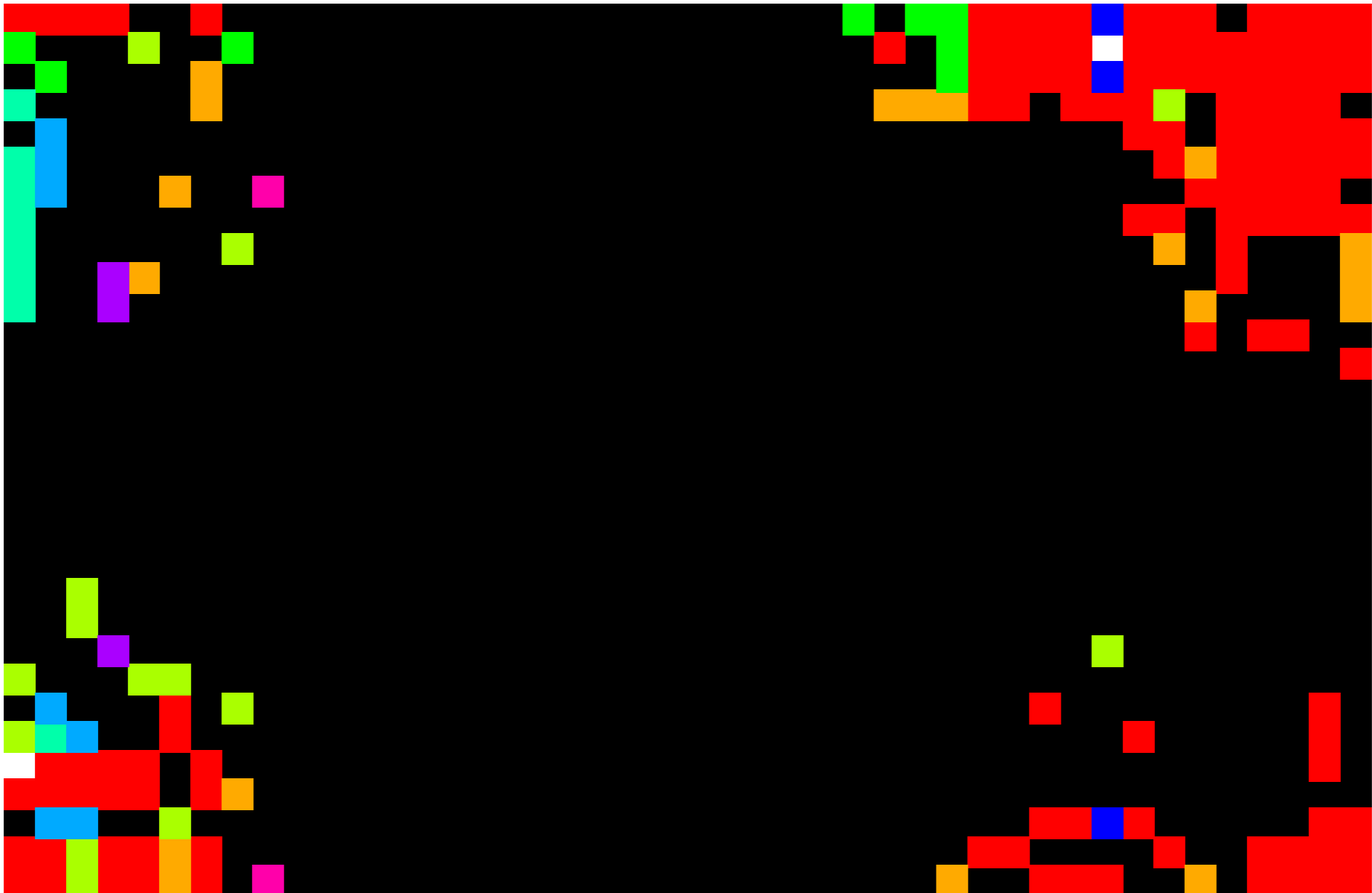
**greedy clustering (“complete clustering”)  
of all correspondences (“columns in multi-alignments”)  
in a Germanic data set**

# Correspondences with high frequency of [t]



# Correspondences with high frequency of [n]

seven  
open  
oven  
rain  
stone  
one  
green  
bone  
hand  
wind  
hound  
moon  
hundred  
nail  
knee  
name  
snow  
night  
north  
needle  
nine  
new  
honey  
naked  
thorn  
horn  
in  
thunder  
corn  
ten  
nine

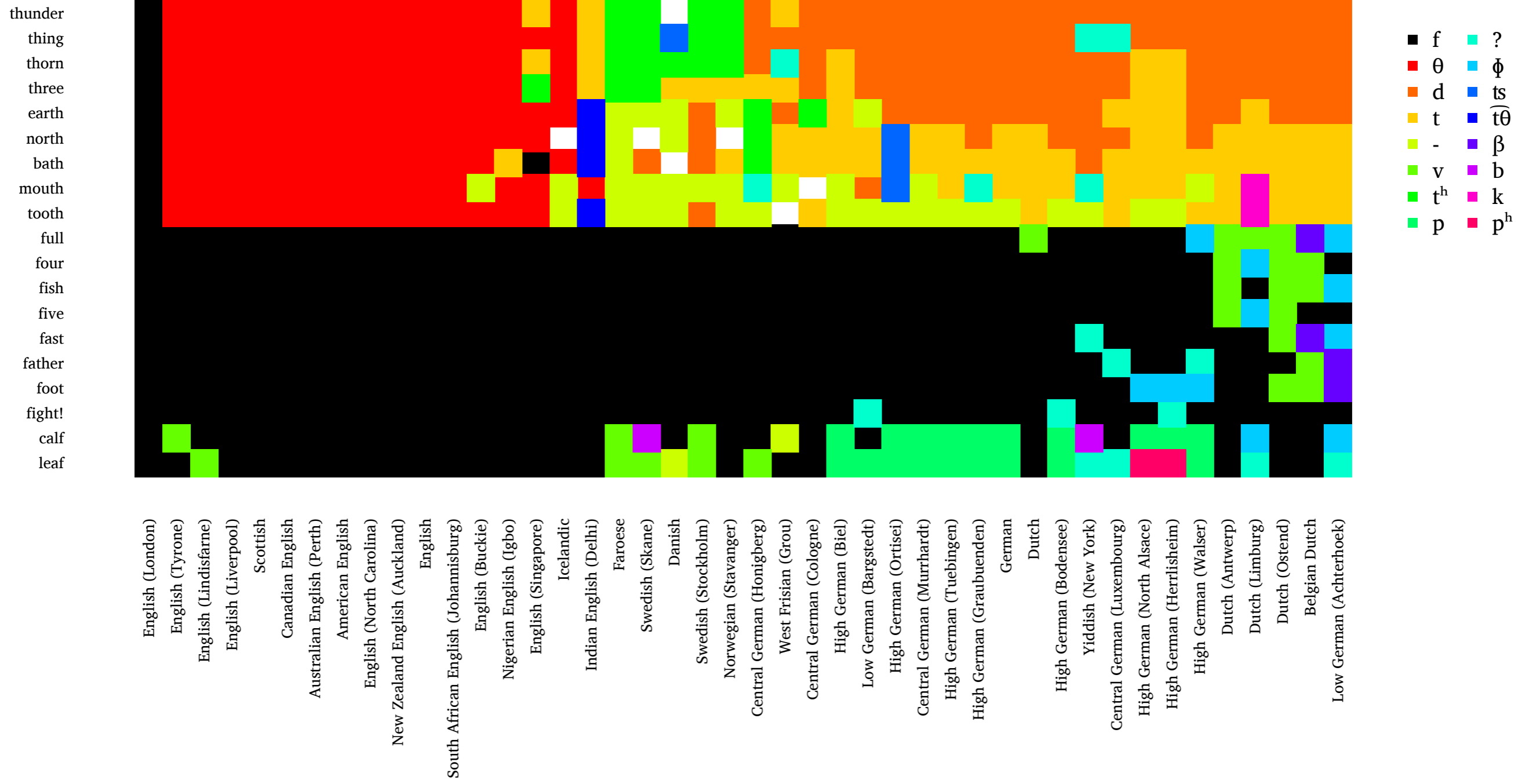


- n
- -
- ŋ
- ?
- m
- nʔ
- ŋ̊
- ŋ̊
- n:
- ŋ̊

Danish  
Icelandic  
Faroese  
Norwegian (Stavanger)  
Swedish (Skane)  
Central German (Luxembourg)  
Swedish (Stockholm)  
West Frisian (Grou)  
Central German (Honigberg)  
English (Lindisfarne)  
English (Liverpool)  
Nigerian English (Igbo)  
Yiddish (New York)  
English (North Carolina)  
South African English (Johannisburg)  
New Zealand English (Auckland)  
English (Buckie)  
Canadian English  
English  
American English  
English (Singapore)  
Scottish  
Australian English (Perth)  
English (Tyrone)  
Belgian Dutch  
English (London)  
Indian English (Delhi)  
High German (Ortisei)  
Low German (Achterhoek)  
German  
Low German (Bargstedt)  
High German (Herrlisheim)  
High German (North Alsace)  
Dutch (Antwerp)  
Dutch  
Dutch (Ostend)  
High German (Graubuenden)  
High German (Biel)  
Central German (Cologne)  
Central German (Murrhardt)  
High German (Bodensee)  
High German (Tuebingen)  
High German (Walser)  
Dutch (Limburg)



# Correspondences with high frequency of [f]



# Prospects

- Add more linguistics into reconstruction
- Retaining variation in the reconstruction process is important for deep phylogenies
- Visual UI for non-computational linguists seems highly important