

Quantifying models of linguistic diversity and change

Michael Cysouw
Philipps-University Marburg

- forget about trees,
focus on reconstruction
- treat change in meaning/function
just like change in form

Central Approach: Multialignment

- Alignments of sounds
- Alignments of words

***In linguistics,
an alignment is a central result,
not an intermediate method***

Multialignment of words

- Based on a sentence-by-sentence alignment, induce word-by-word alignment
- Translations can be (and often are!) quite different
- Bi-text alignment is widely researched problem
- Multit-text alignment not so much

Kong Herodes blev skrækslagen , og Jerusalem begyndte at summe af rygter .

Als dies dem König Herodes zu Ohren kam , erschrak er , und mit ihm entsetzte sich auch ganz Jerusalem .

But Herodes the king heard , and was troubled , and all Urishlem with him .

Pegar Heródes heyrði þetta , varð hann skelkaður og öll Jerúsalem með honum .

Als des dr Kenig Herodes ghärt het , isch scha (er) arg vuschrocke un mit nem ganz Jerusalem ,

Kort voor lank het Herodes ook van die geleerdes uit die ooste se storie te hore gekom . Hy was baie omgekrap oor wat hulle oor die nuwe Joodse koning gesê het . So ook die res van Jerusalem .

Der König Herodes war total aufgebracht , als er das hörte , und nicht nur er , alle in Jerusalem waren das .

Kong Herodes blev skrækslagen , og Jerusalem begyndte at summe af rygter .

Als dies dem König Herodes zu Ohren kam , erschrak er , und mit ihm entsetzte sich auch ganz Jerusalem .

But Herodes the king heard , and was troubled , and all Urishlem with him .

Pegar Heródes heyrði þetta , varð hann skelkaður og öll Jerúsalem með honum .

Als des dr Kenig Herodes ghärt het , isch scha (er) arg vuschrocke un mit nem ganz Jerusalem ,

Kort voor lank het Herodes ook van die geleerdes uit die ooste se storie te hore gekom . Hy was baie omgekrap oor wat hulle oor die nuwe Joodse koning gesê het . So ook die res van Jerusalem .

Der König Herodes war total aufgebracht , als er das hörte , und nicht nur er , alle in Jerusalem waren das .

Kong Herodes blev skrækslagen , og Jerusalem begyndte at summe af rygter .

Als dies dem König Herodes zu Ohren kam , erschrak er , und mit ihm entsetzte sich auch ganz Jerusalem .

But Herodes the king heard , and was troubled , and all Urishlem with him .

Pegar Heródes heyrði þetta , varð hann skelkaður og öll Jerúsalem með honum .

Als des dr Kenig Herodes ghärt het , isch scha (er) arg vuschrocke un mit nem ganz Jerusalem ,

Kort voor lank het Herodes ook van die geleerdes uit die ooste se storie te hore gekom . Hy was baie omgekrap oor wat hulle oor die nuwe Joodse koning gesê het . So ook die res van Jerusalem .

Der König Herodes war total aufgebracht , als er das hörte , und nicht nur er , alle in Jerusalem waren das .

Kong Herodes blev skrækslagen , og Jerusalem begyndte at summe af rygter .

Als dies dem König Herodes zu Ohren kam , erschrak er , und mit ihm entsetzte sich auch ganz Jerusalem .

But Herodes the king heard , and was troubled , and all Urishlem with him .

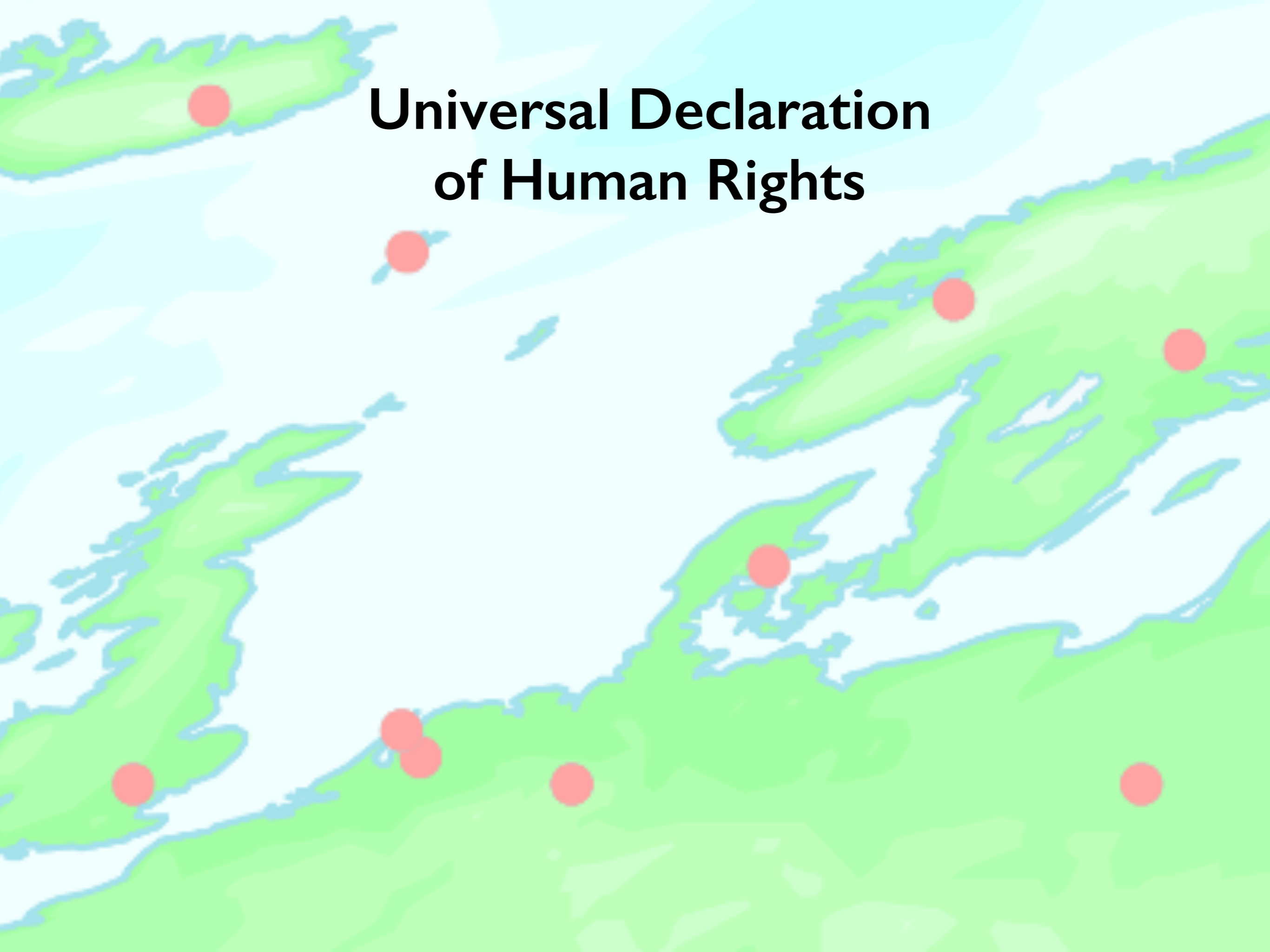
Pegar Heródes heyrði þetta , varð hann skelkaður og öll Jerúsalem með honum .

Als des dr Kenig Herodes ghärt het , isch scha (er) arg vuschrocke un mit nem ganz Jerusalem ,

Kort voor lank het Herodes ook van die geleerdes uit die ooste se storie te hore gekom . Hy was baie omgekrap oor wat hulle oor die nuwe Joodse koning gesê het . So ook die res van Jerusalem .

Der König Herodes war total aufgebracht , als er das hörte , und nicht nur er , alle in Jerusalem waren das .

Universal Declaration of Human Rights





in, im German

in Dutch

in English

yn Frisian

in Afrikaans

in Scots

אין Yiddish

í Icelandic

i Swedisch

i Nynorsk

i Bokmål

í Faroese

i Danish

Article 10

English: “everyone is entitled in full equality to a fair and public hearing by an independent and impartial tribunal , **in the determination of his rights and obligations and of any criminal charge against him .**”

Afrikaans: “elkeen het , in volle gelykheid , die reg tot ’n regverdige en openbare verhoor deur ’n onafhanklike en objektiewe tribunaal , **in die bepaling van sy regte en verpligtinge en die ondersoek van enige kriminele saak teen hom .**”

German: “jeder hat **bei der feststellung seiner rechte und pflichten** sowie bei einer gegen ihn erhobenen strafrechtlichen beschuldigung in voller gleichheit anspruch auf ein gerechtes und öffentliches verfahren vor einem unabhängigen und unparteiischen gericht .”

Preamble, sentence 5

English:

“whereas the peoples of the united nations have in the charter reaffirmed their faith in fundamental human rights , in the dignity and worth of the human person and in the equal rights of men and women and have determined to promote social progress and better standards of life in larger freedom ,”

German:

“da die völker der vereinten nationen in der charta ihren glauben an die grundlegenden menschenrechte , an die würde und den wert der menschlichen person und an die gleichberechtigung von mann und frau erneut bekräftigt und beschlossen haben , den sozialen fortschritt und bessere lebensbedingungen in größerer freiheit zu fördern ,”

Preamble, sentence 5

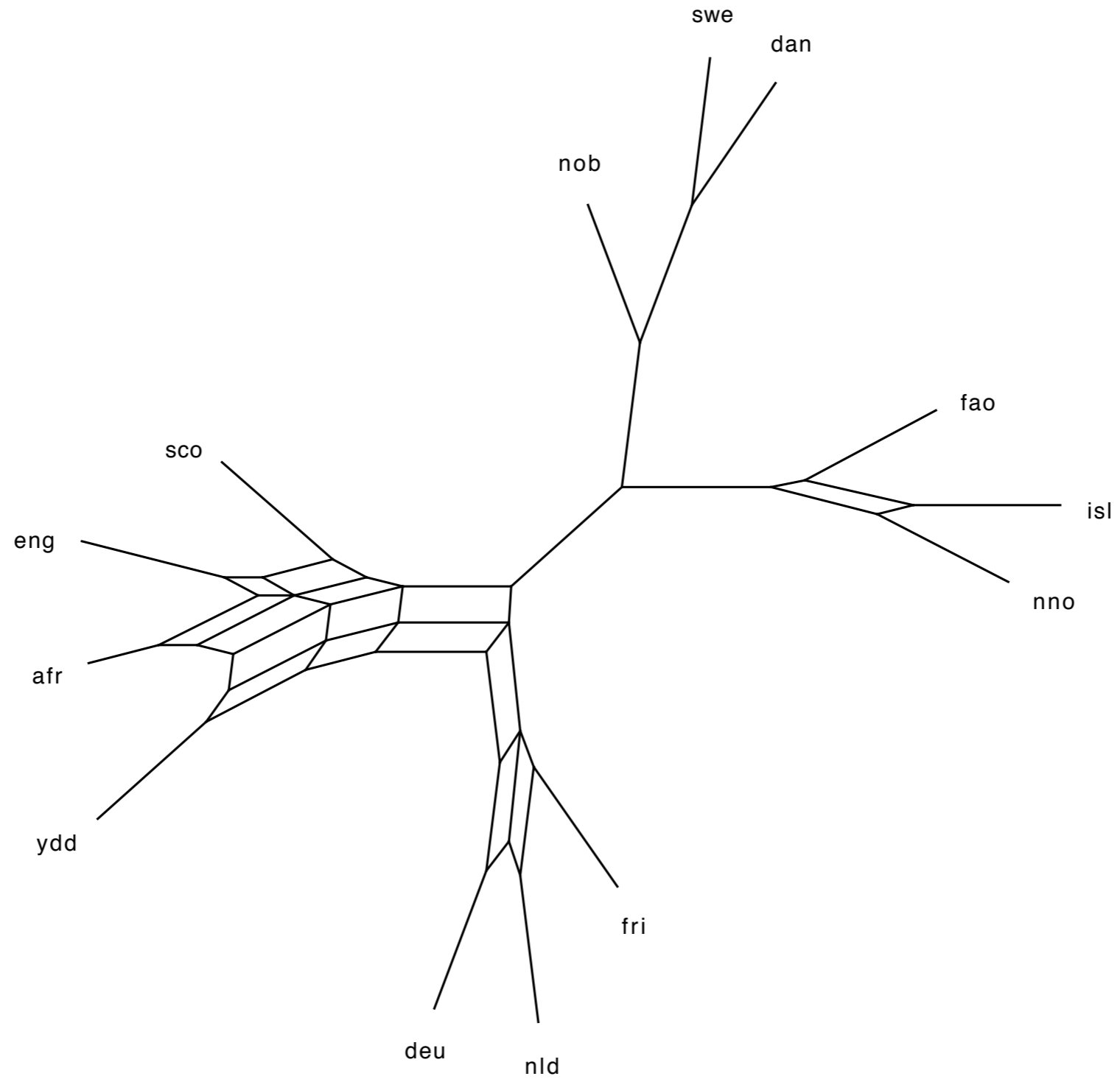
English:

“whereas the peoples of the united nations have *in the charter* reaffirmed their faith *in fundamental human rights* , *in the dignity* and worth of the human person and *in the equal rights* of men and women and have determined to promote social progress and better standards of life *in larger freedom* ,”

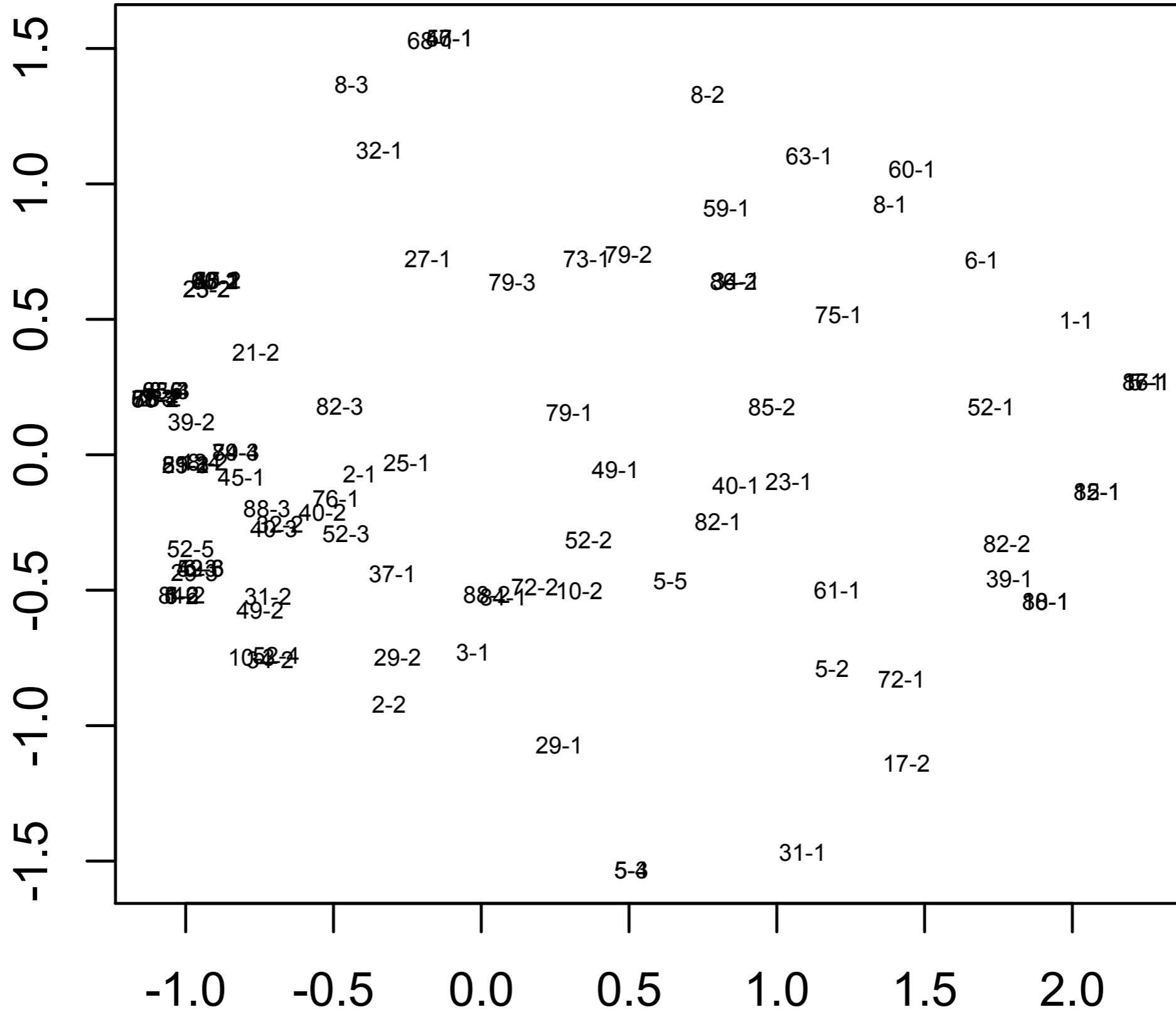
German:

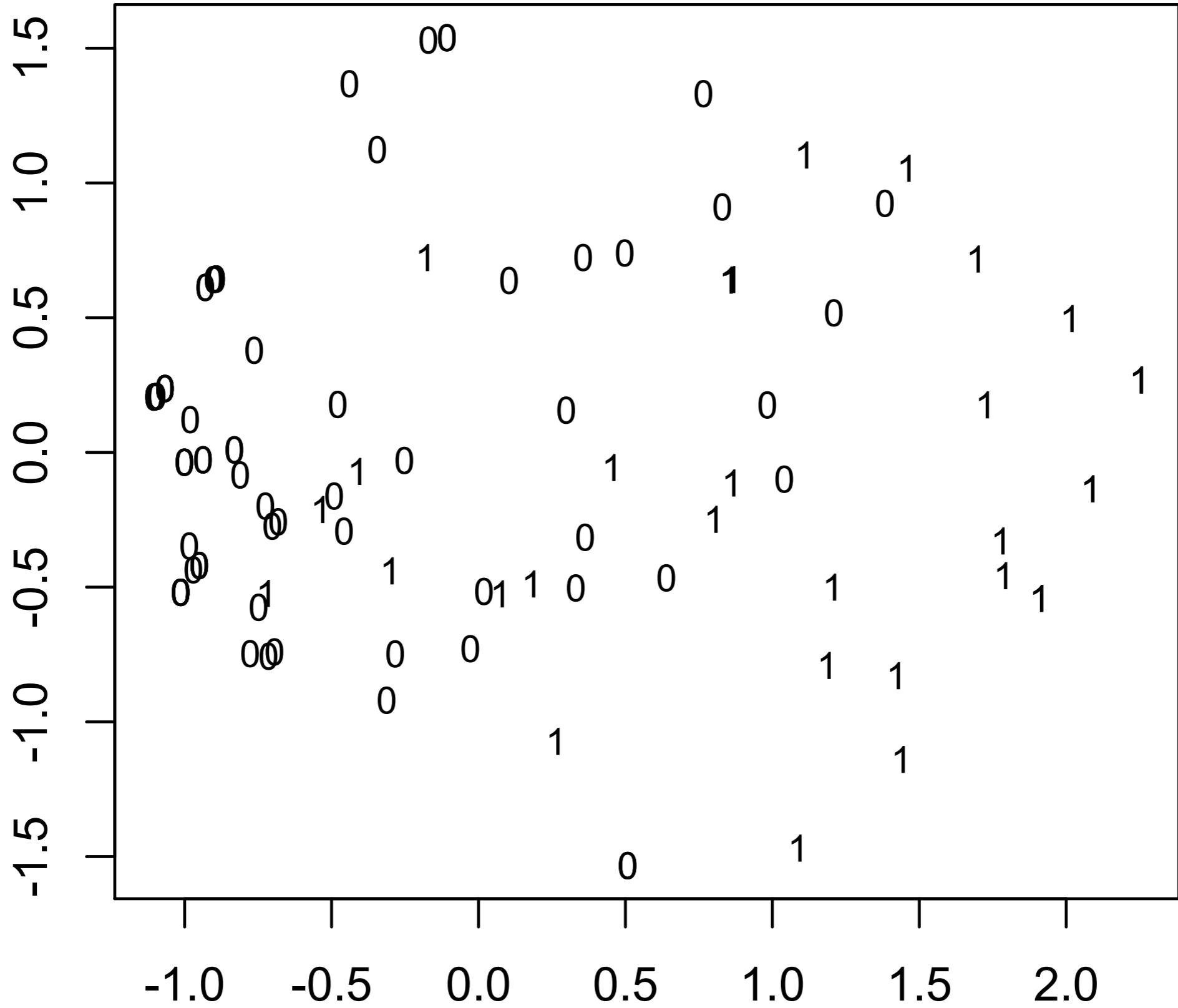
“da die völker der vereinten nationen *in der charta* ihren glauben *an die grundlegenden menschenrechte* , *an die würde* und den wert der menschlichen person und *an die gleichberechtigung* von mann und frau erneut bekräftigt und beschlossen haben , den sozialen fortschritt und bessere lebensbedingungen *in größerer freiheit* zu fördern ,”

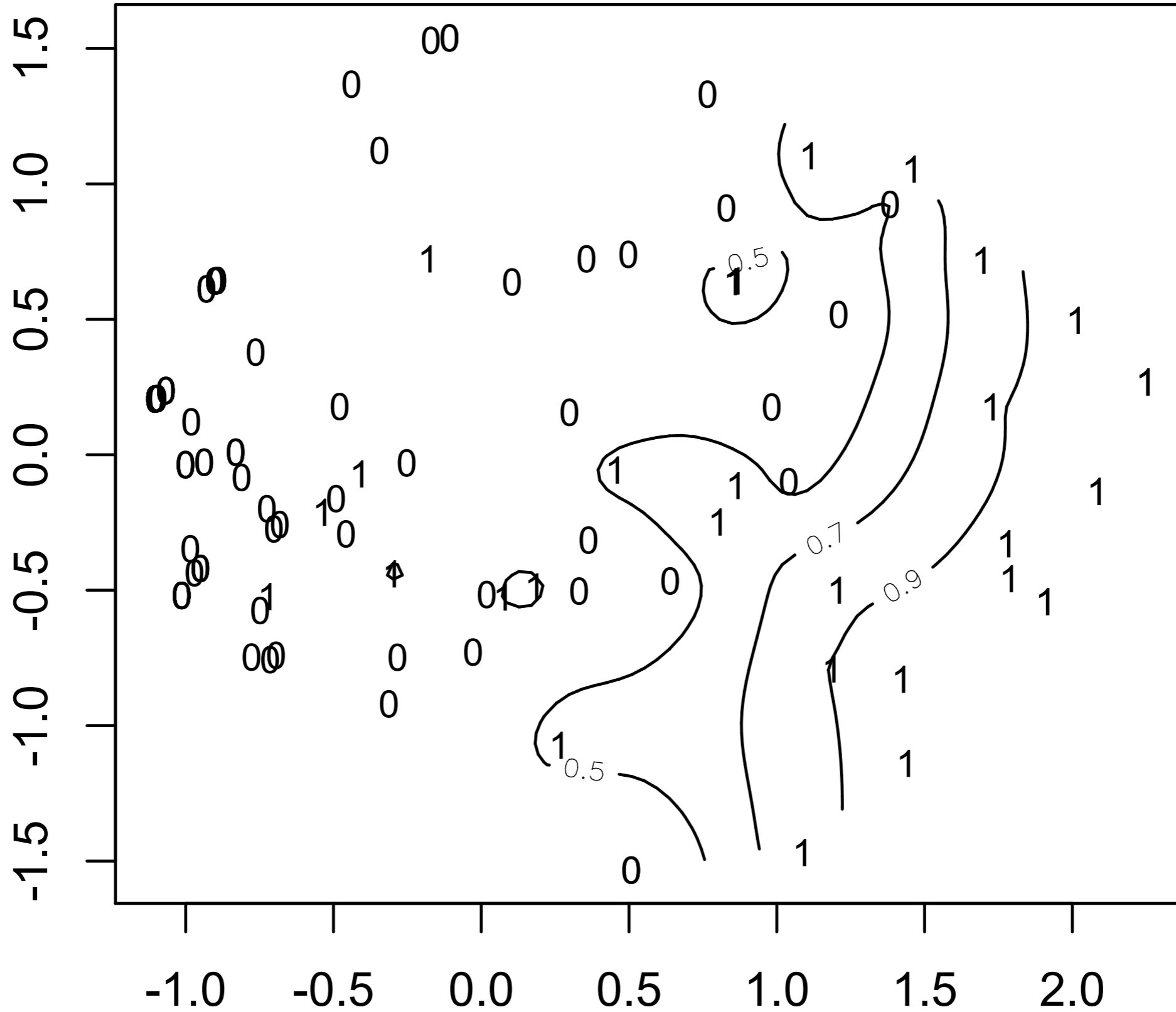
deu	2						
eng	5						
nld	5						
fri	6						
afr	5						
ydd	4						
sco	5						
nno	3						
nob							
swe							
fao	3						
dan							
isl							

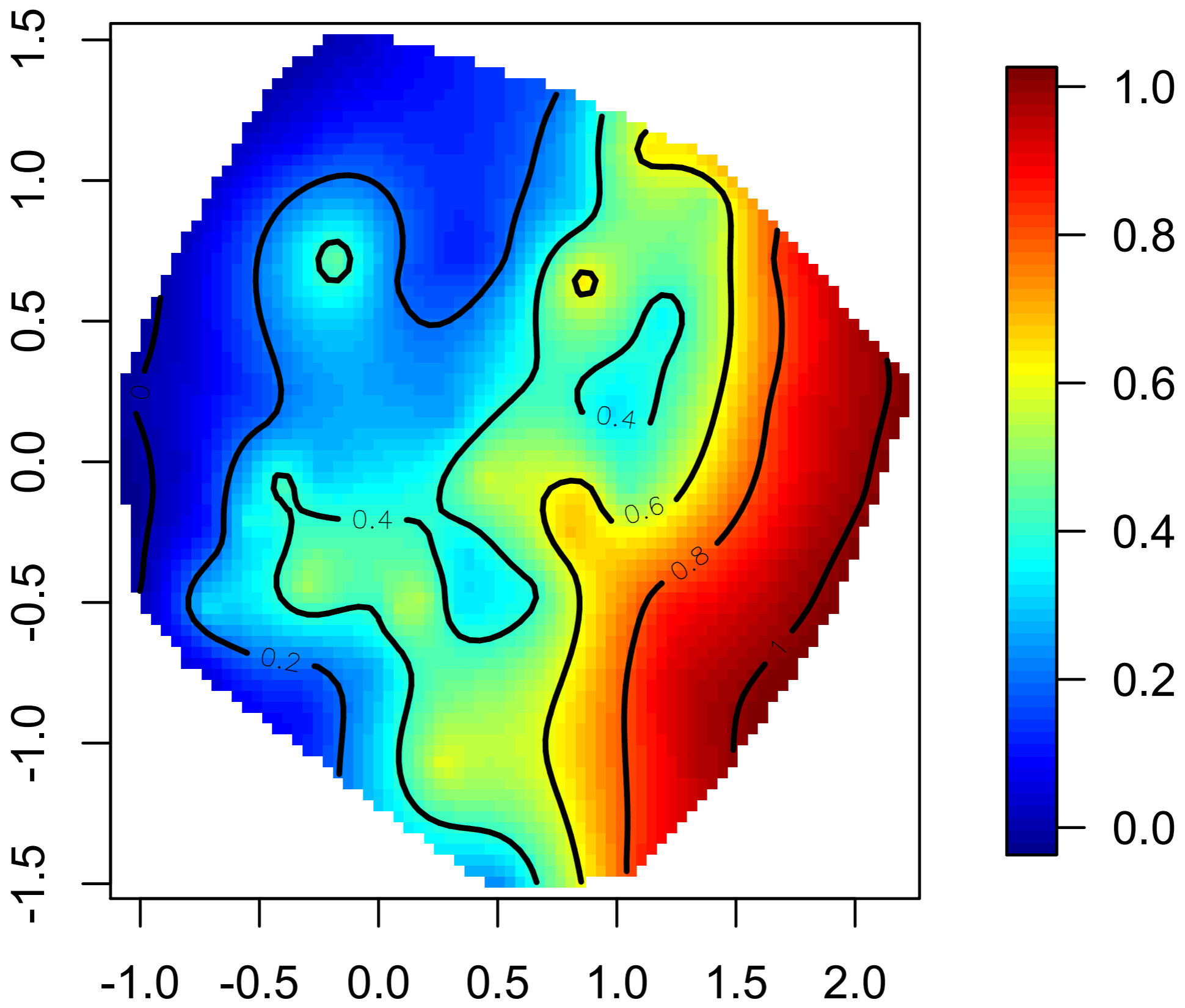


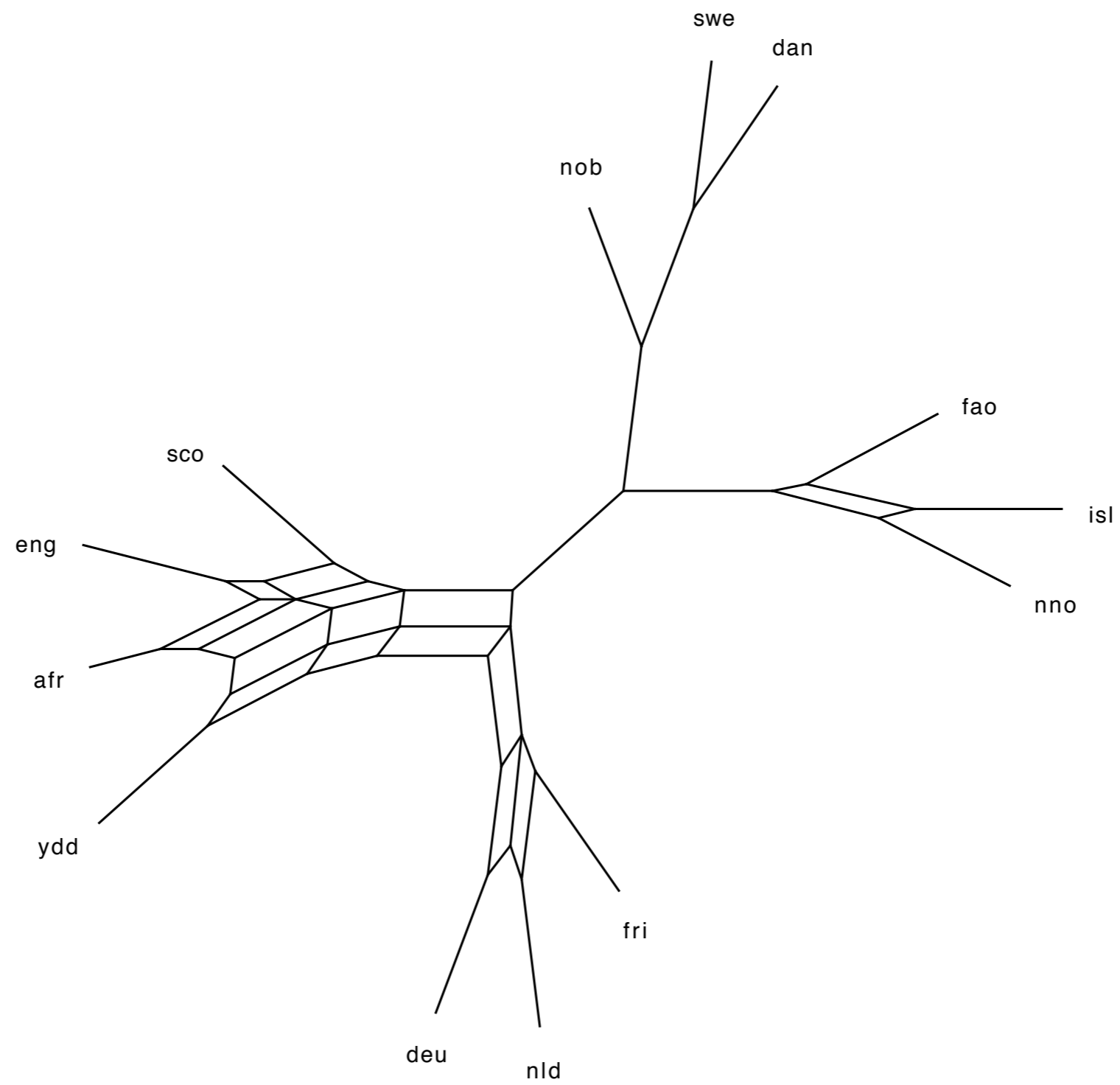
**Consensus network of 4 optimal trees
according to dollo Maximum Parsimony**

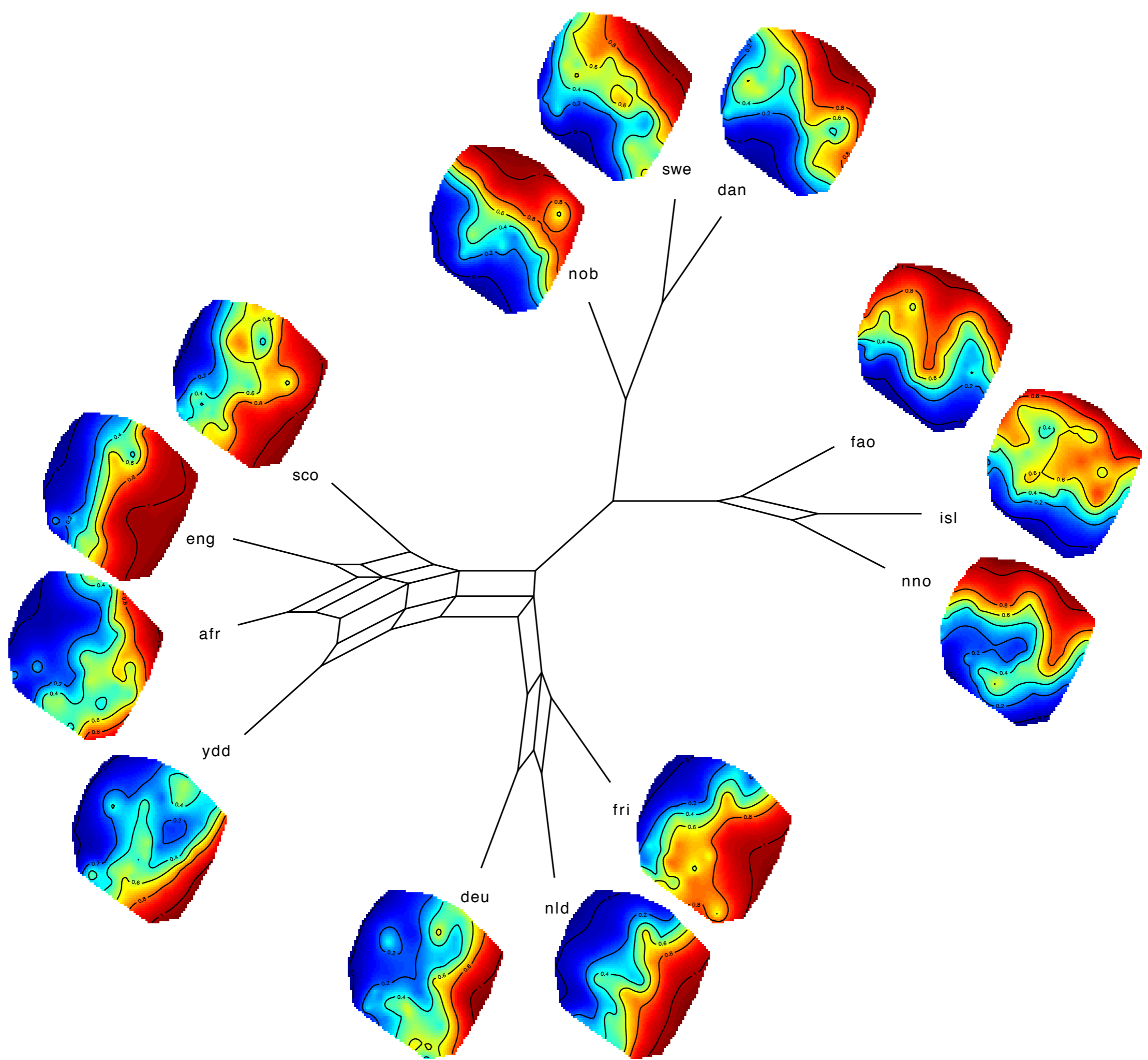












Multialignment of Sounds

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter	178 'vɪntə ^h ʁ	ʁ = kontinuant ə bereit velarisier
* 2	fliegen	56 'flaɪ→ə ^h θə	'fliegen die', Segmentierung unklar - kein geminiertes [t]
3	Blätter	23 'blɛ:də ^h ʁ	
4	Luft	103 lu ^h ft ^h	ʁ = kontinuant
5	hört	89 hɪ ^h t ^h	ʁ = kontinuant
6	gleich	78 k ₊ lɛɪ ^h k ₊	folgt P
7	schneien	130 ʃnɛɪ ^h n	
8	Wetter	174 /	statt demoa 'vɪtə ^h ʁɔŋə
9	tu	151 dɛ→v	

Ort der Mundart/Kreis Astfeld/ Günthersheim	Aufnahme-Nr. I/62	Transkribent Angelika Braun	Listentyp A
---	----------------------	--------------------------------	----------------

Phonetischer Atlas von Deutschland

- Wenker-sentences recorded in the 1960s (with additions in the 1970s)
- Selected words from the recordings were transcribed on paper in the 1980s
- A joint project between Marburg and Groningen digitized the data in the 2000s
- In total 29530 words distributed over 183 locations and 186 cognate sets

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter	178 vɪntɐ	k = kontinuant b = kontinuant
2	fliegen	56 flaɪ̯ə	'fliegen die', Sequenzierung unklar - kein genügendes [tʰ]
3	Blätter	23 blɛtɐ	
4	Luft	103 lʊft	= kontinuant
5	hört	89 hœrt	= kontinuant
6	gleich	118 glɛɪ̯ç	
7	schneien	130 ʃnɛi̯ən	
8	Wetter	17 vɛtɐ	statt dem 'vɛtəʀɔgə'
9	tu	151 tu	

Ort der Mundart/Kreis Astfeld/ Gandersheim	Aufnahme-Nr. I/62	Transkribent Angelika Trauer	Listentyp A
Planrechteck X 29	Aufnahmedatum 20.11.1965	Transkribiert von 14.6.1985 bis 24.7.1985	

Besprochen von 24.07.1985
25.07.1985 UStv

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter	'vɪntə ^h	ɸ = kontinuant ə bereit velarisier
* 2	fliegen	'flaɪ→ə ^h	'fliegen die', Sequenzie- rung unklar - kein geminiertes [t]
3	Blätter	'blɛ:də ^h	
4	Luft	lu ^h ft ^h	ɸ = kontinuant
5	hört	hɪ ^h t ^h	ɸ = kontinuant
6	gleich	klɛ ₊ ɪk _~	folgt P
7	schneien	'ʃnɛ→ɪən	
8	Wetter	/	statt demoa 'vɪtə ^h ʀɔŋə
9	tu	dɛ→u	

Ort der Mundart/Kreis	Aufnahme-Nr.	Transkribent	Listentyp
-----------------------	--------------	--------------	-----------

Phonetischer Atlas von Deutschland

X 27	20.11.1965	bis 24.7.1985
------	------------	---------------

- Digitised in X-SAMPA, converted back to match original transcriptions, minor corrections for consistency of encoding

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
----------	-----------	---------------	-------------

- The data is transcribed in high phonetic detail (3786 different phonetic segments)

- We make the complete data available

- ▶ electronically, separated by phonetic segments
- ▶ as close as possible to the original source
- ▶ including all idiosyncrasies
- ▶ github.com/cysouw/PAD

9	tu	151	de-u
---	----	-----	------

LOCATION	WORD
Aachen	a:ph
Adorf	ɑ:b ^h ə
Ahrbergen	o→ɔphə
Albersloh	ɑ:p ^h ə
Allna	ɑϕh
Altenberg	ʌfɛ
Altentrüdin	af
Altlandsberg	ɑ'fə'
Altwarp	o:ph
Astfeld	ɒ':p ^h ə
Atzendorf	afɛ
Ballhausen	ʌ'fə
Bardenfleth	ɔ:p̄ϕ
Barssel	ɒ:p ^h ə
Bempflingen	af:
Bennin	ɔp ^h
Billingsbach	af
Bockelwitz	ʌvə
Bonn	ɑ:p'
Borstendorf	ɣf:
Breddin	ɒ:ph
Brelingen	ɑfβə
Bremscheid	ɒ':p ^h ə
...	...

A	FF	E
a:	ph	-
ɑ:	b ^h	ə
o→ɔ	ph	ə
ɑ:	p ^h	ə
ɑ	ϕh	-
ʌ	f	ɛ
a	f	-
ɑ'	f	ə'
o:	ph	-
ɒ':	p ^h	ə
a	f	ɛ
ʌ'	f	ə
ɔ:	p̄ϕ	-
ɒ:	p ^h	ə
a	f:	-
ɔ	p ^h	-
a	f	-
ʌ	v	ə
ɑ:	p'	-
ɣ	f:	-
ɒ:	ph	-
ɑ	f̄β	ə
ɒ':	p ^h	ə
...

● **Workflow:**

- ▶ Tokenisation of segments (github.com/cysouw/qlcData)
- ▶ Automatic alignment using **LingPy** (github.com/lingpy)
- ▶ Manual correction using **Alignment Editor** (github.com/digitallinguist/msa-editor)
- ▶ Separation of cognates (e.g. *Samstag* vs. *Sonnabend*)
- ▶ Annotation of columns (e.g. many-to-one alignments, metathesis)
- ▶ Merging of complex columns and removing boundaries

MSA Editor

Choose Files 3 files Augenblick_1013.msa View Edit Reload Save

COLUMNID	1	2	3	4	5	6	7	8	9
STANDARD	Au	g	e	n	b	l	i	ck	(e)
Adorf	æ→u	ʁ	-	ŋ	b	l	ɛ	k ^h	-
Ahrbergen	æ→õ	ʁ	ə	m	b	l	ɪ	k'	-
Albersloh	æ→u	-	-	m	β	l	ɪ	k	-
Allna	ɔɪ	-	-	-	p	l	æ ^c	x	-
Altenberg	ʁʁ	ʁ	ã	-	b	l	ɪ	k	-
Altentrüdin	æ→u	ʁ	ə	-	p	l	ɪ	g	-
Altlandsberg	ɑ'→u	ʁ	-	ŋ	b	l	ɪ	k̄x	-
Altwarf	õu	-	-	ŋ	b	l	ɪ	k	-
Astfeld	ʊɪ	ʁ	ə	m	b	l	ɪ	k ^h	ə
Ballhausen	ɑ→u	ʁ	-	ŋ	p	l	ɪ	k	-
Bardenfleth	oɪ	g	-	ŋ	b	l	ɛ	k̄x ₊	-
Barssel	oɪ	g	-	ŋ	p	l	ɪ	k ₊	-
Bempflingen	æ→u	g	ə	-	b	l	ɪ	c ^h	-
Bennin	oɪ	-	-	ŋ	b	l	ɪ	x	-
Billingsbach	ɑɪ	x	ə	-	p	l	ɪ	k̄x	-
Bockelwitz	ɑ→u	ʁ	-	ŋ	b	l	ɪ	k	-
Borstendorf	ɔɪ	ʁ̄x	-	ŋ	p	l	ɛ	k	-

github.com/digitallinguist/msa-editor

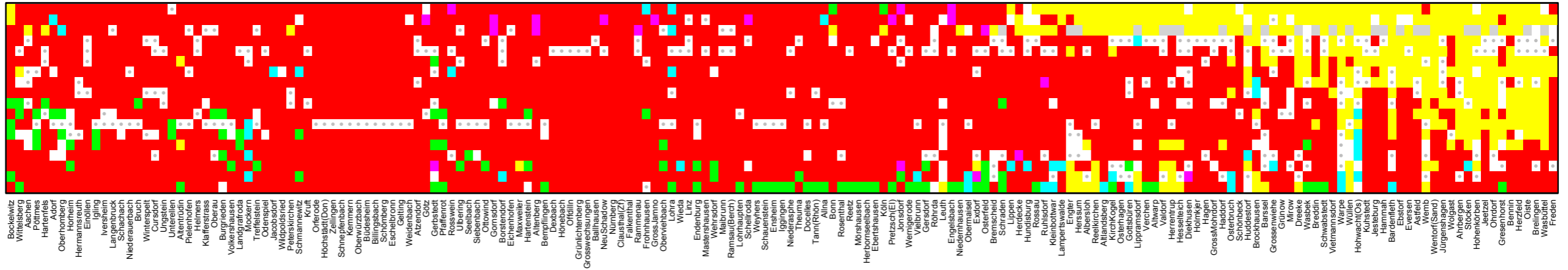
Correspondence Sets

- more than 700 columns of aligned segments (“correspondences”)
- Comparative-historical linguistics uses clusters of correspondences (“correspondence sets”)
- Automatic clustering of columns is a good start, but needs correction
- Visualisations in R
github.com/cysouw/qlcVisualize

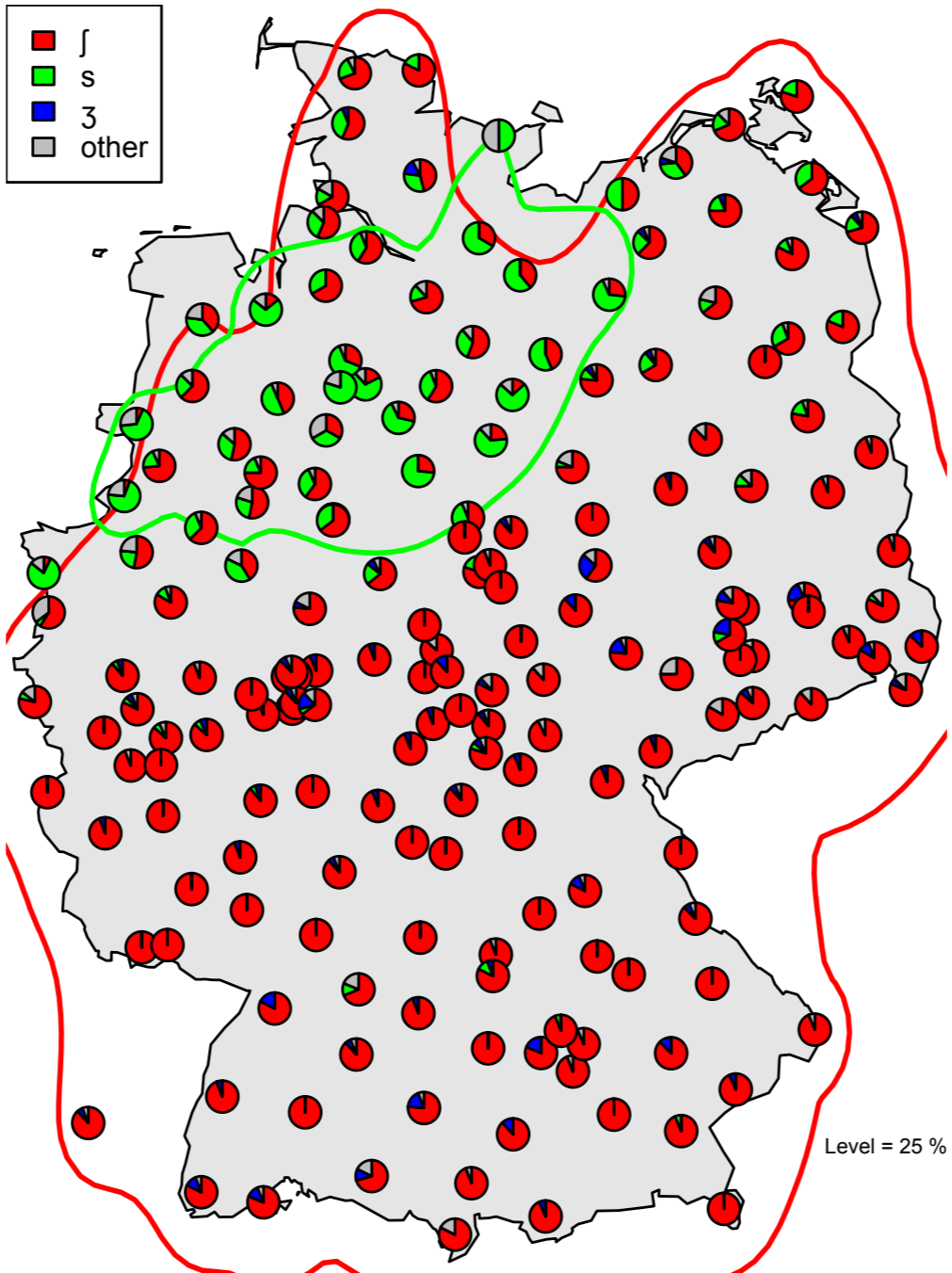
limage (“level-image”)
from R-package *qlcVisualize*
(Cysouw 2015)



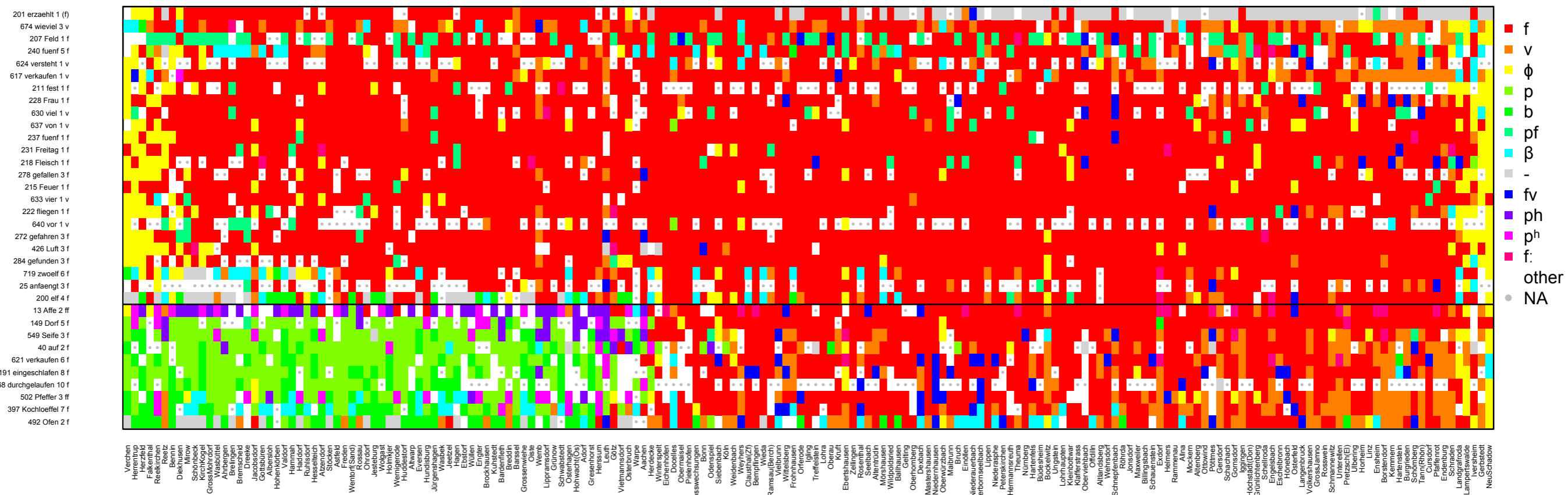
700 Würst 5 s
 172 Durst 4 s
 142 Donnerstag 8 s
 312 gestohlen 3 s
 626 versteht 4 s
 85 bestellt 3 s
 581 Stueckchen 1 s
 319 gestorben 3 s
 534 schwarz 1 sch
 188 eingeschlafen 5 sch
 522 Schnee 1 sch
 525 schneien 1 sch
 516 schlechte 1 sch
 538 Schwester 1 sch
 530 schoene 1 sch
 221 Tisch 3 sch
 587 Tisch 3 sch
 156 Dreschen 7 sch



Correspondences "sch"

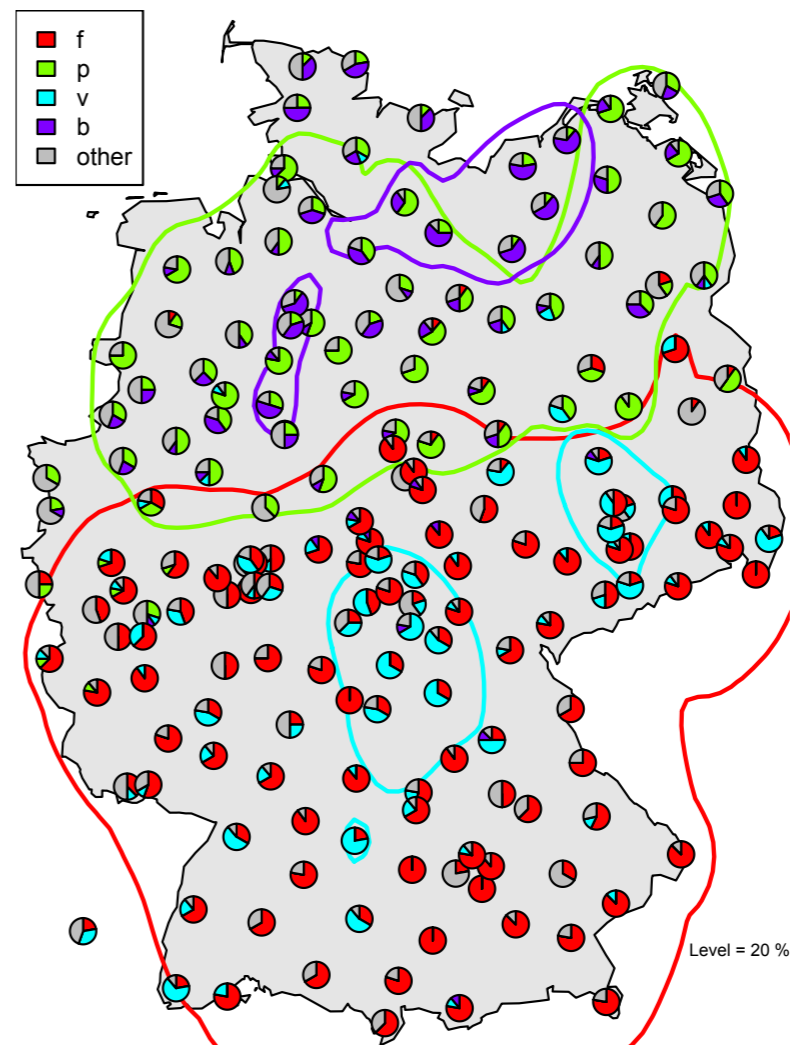
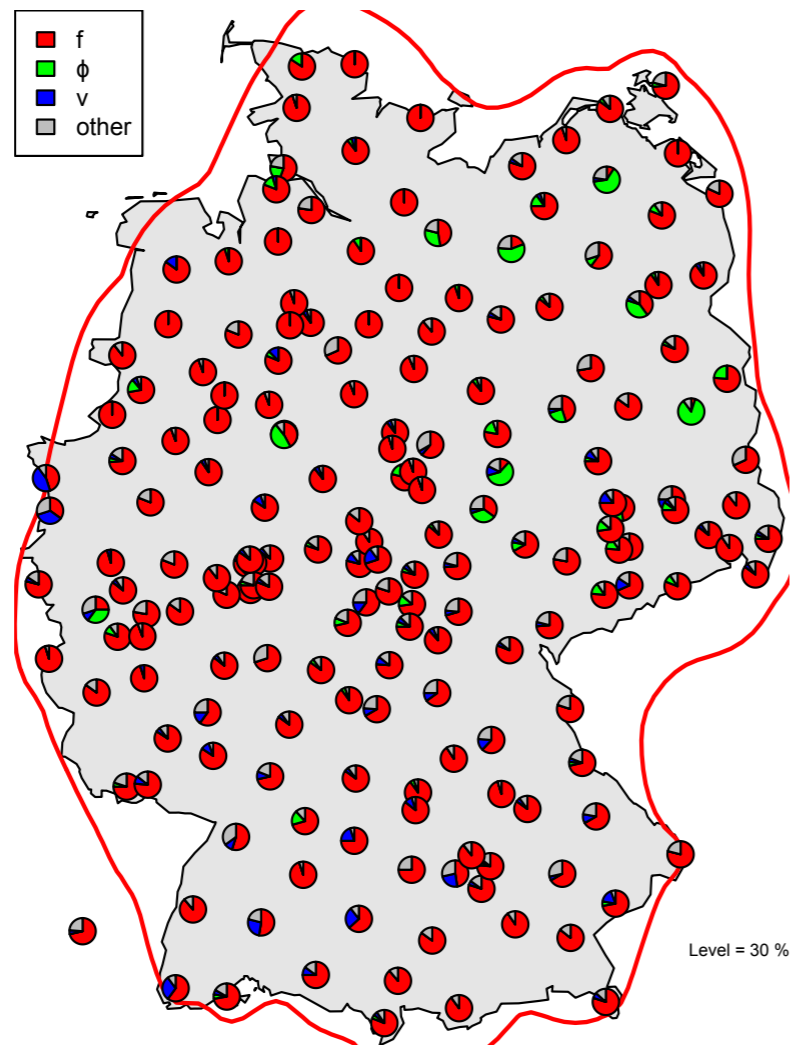


■ j
■ s
■ 3
■ 6
■ -
■ j:
 other
 ● NA



Correspondences "f"

Correspondences "f/p"



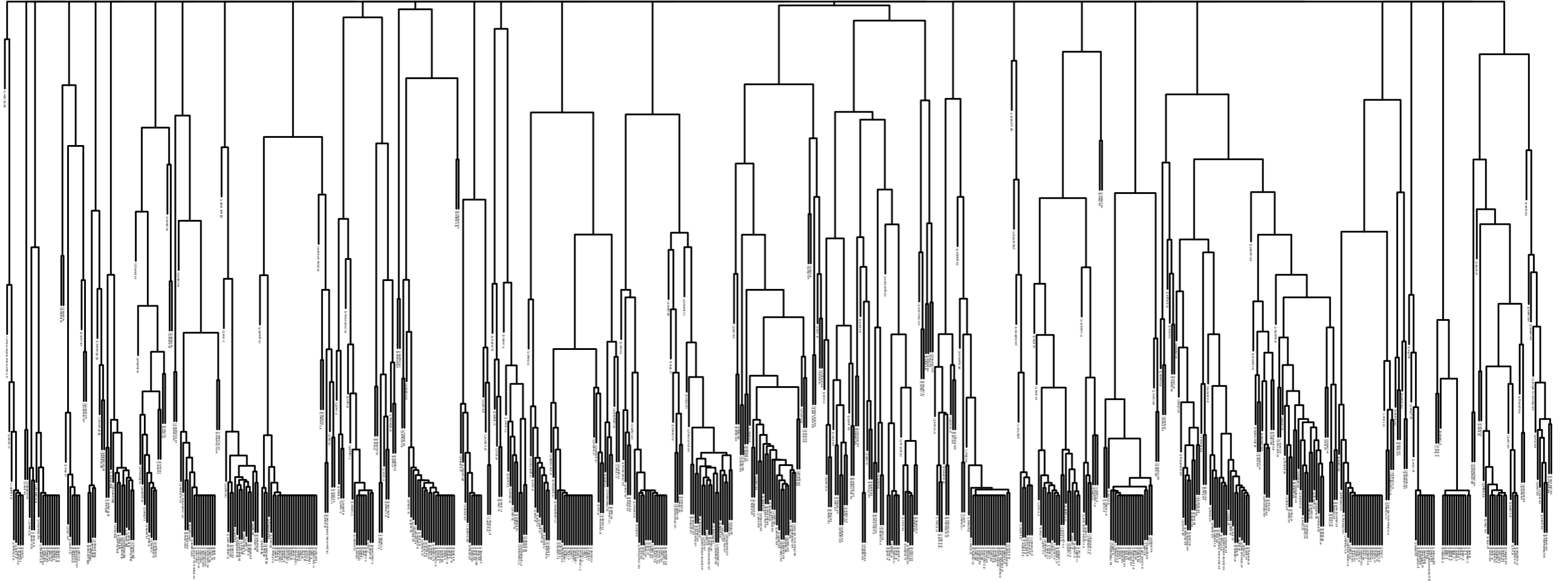
Bulgarian dialect data

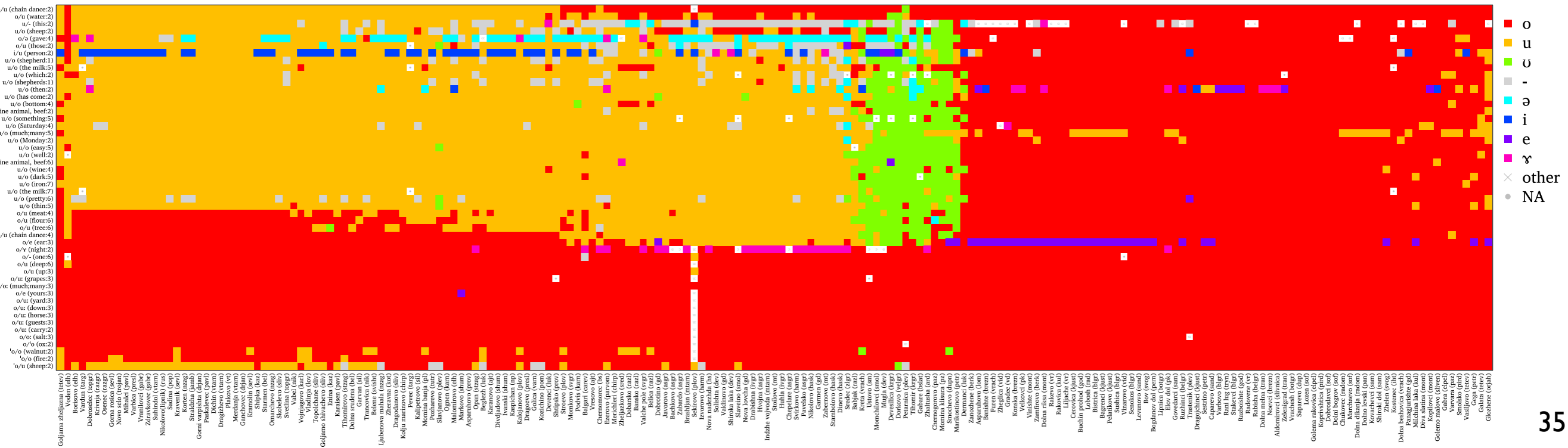
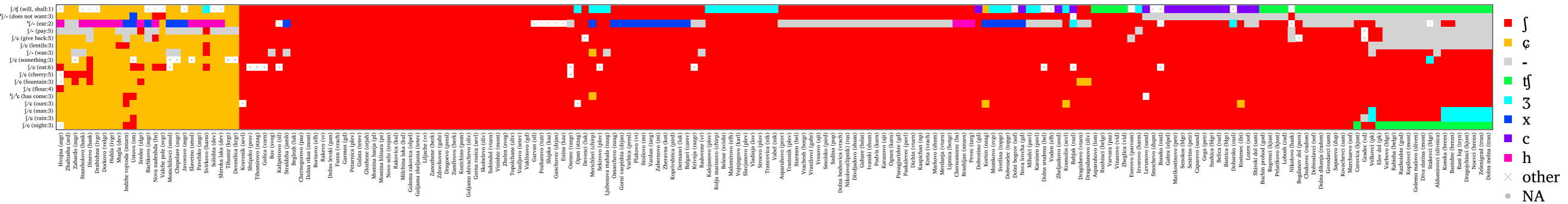
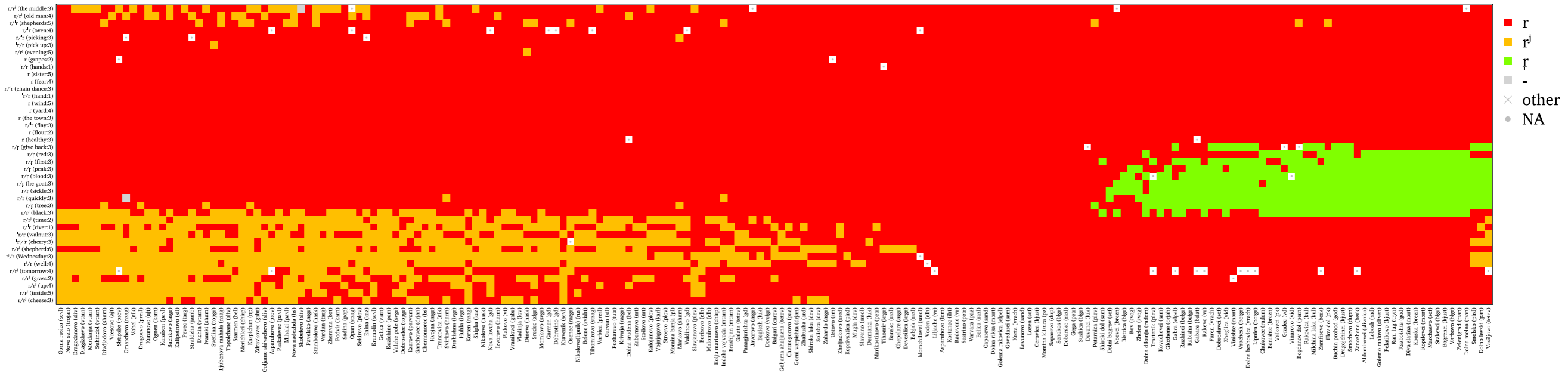
- Data from Jelena Prokić
- Taken from the *benchmark database of phonetic alignments* (BDPA)

List, Johann-Mattis and Jelena Prokić. (2014). A benchmark database of phonetic alignments in historical linguistics and dialectology. In: Proceedings of the International Conference on Language Resources and Evaluation (LREC), May 2014, Reykjavik. 288-294.

Bulgarian dialect data

- Data from 197 villages for 152 words, resulting in 807 aligned columns
- There are many possible methods to cluster these into correspondence sets
- And they give different results!
(though mostly smaller or larger sets)
- Here: a simple greedy clustering into 27 sets



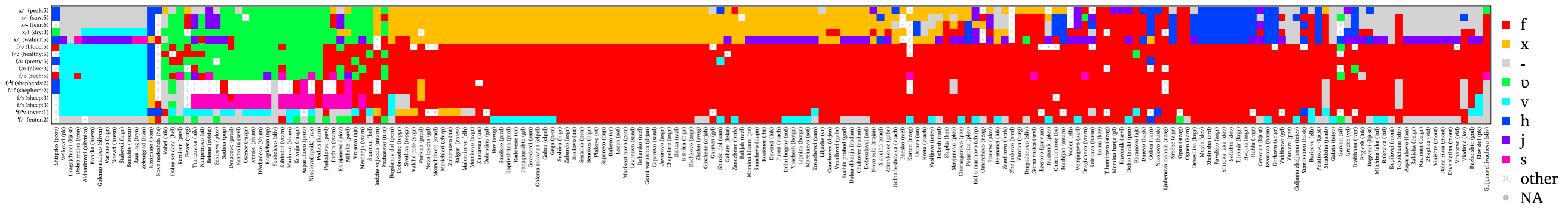


Proto-Bulgarian

	Labial	Dental	Palatal	Velar
Plosive	p b	t d		k g
Affricate		\overline{ts}	$\overline{tʃ}$	
Fricative	f v	s z	ʃ ʒ	
Nasal	m	n	n ^j	
Lateral		l		
Trill		r		
Approximant			j	

Old Church Slavonic

	Labial	Dental	Palatal	Velar
Plosive	p b	t d		k g
Affricate		\overline{ts} \overline{dz}	$\overline{tʃ}$	
Fricative		s z	ʃ ʒ	x
Nasal	m	n	n ^j	
Lateral		l	l ^j	
Trill		r	r ^j	
Approximant	v		j	



Reconstructed /f/ or /x/ or both ?

Conclusion

- Value the actual decisions that you make in the analysis, and share them
- Multialignments and reconstructions are central decisions, that can be profitably discussed and collectively improved