

Using Parallel text to quantify models of linguistic diversity and change

Michael Cysouw
Philipps-University Marburg

Central Approach: Multialignment

- Alignments of words
- Alignments of sounds

***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

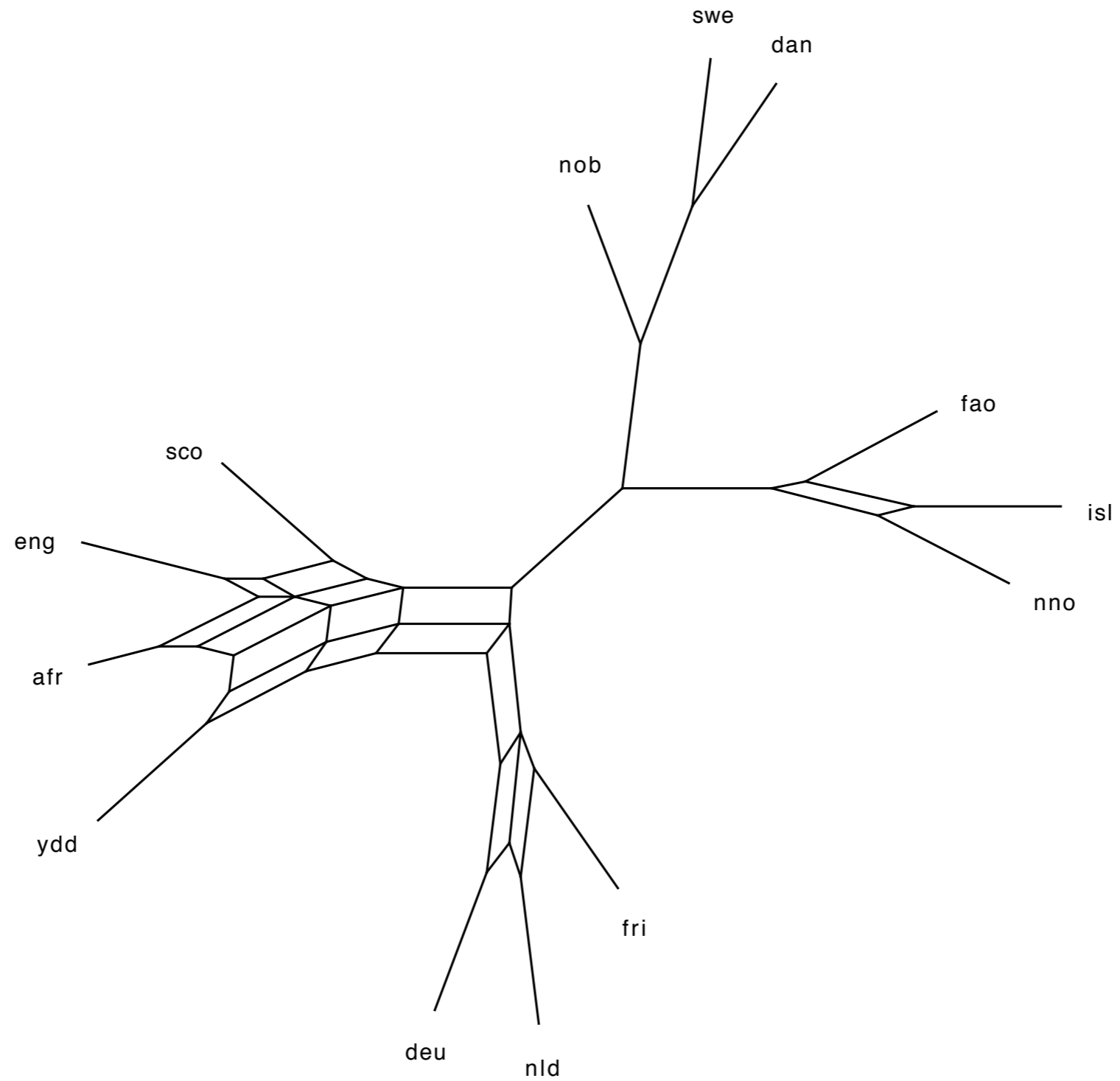
Article 10

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 .”

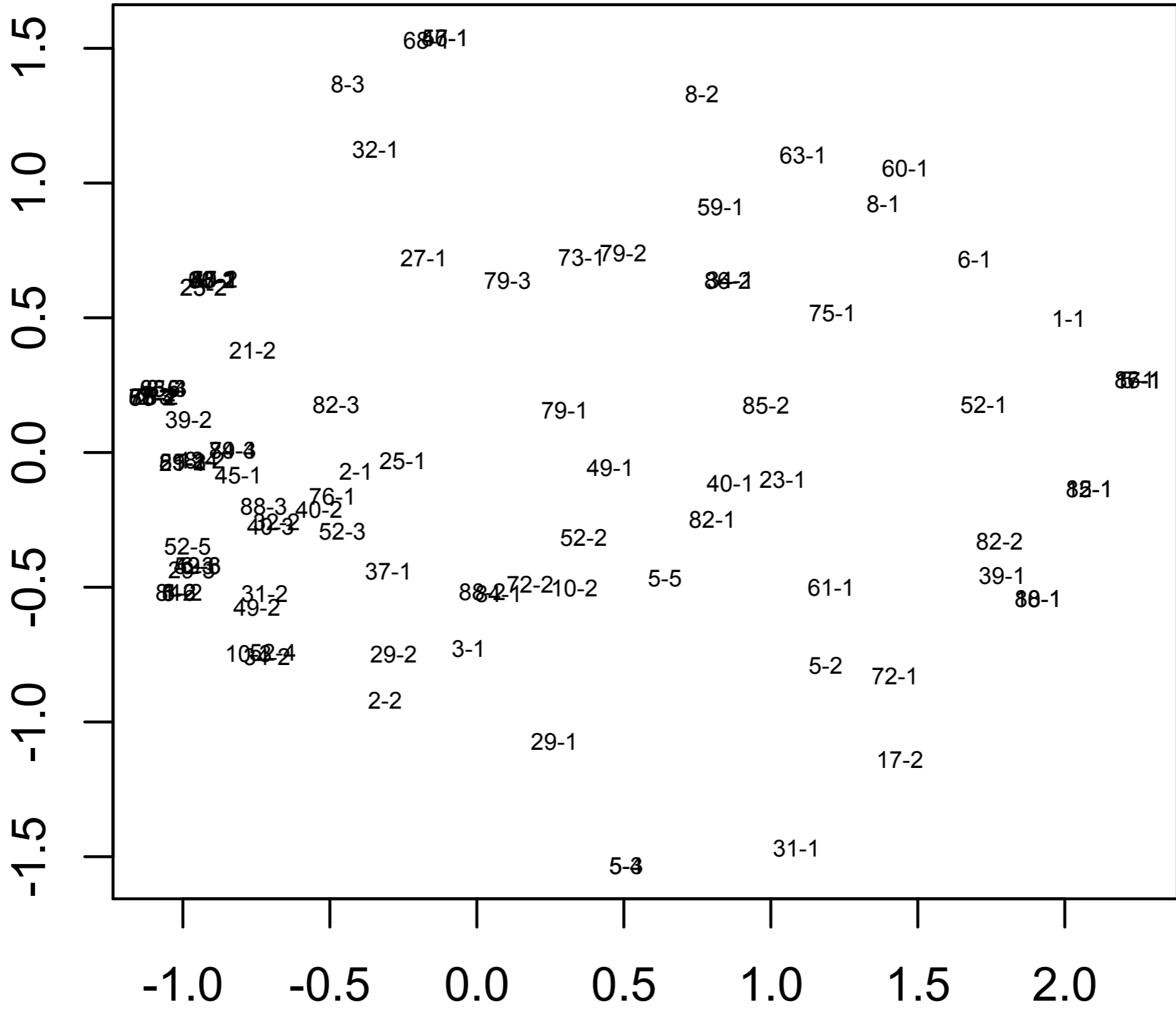
Afrikaans: “elkeen het , in volle gelykheid , die reg tot ’n regverdige en openbare verhoor deur ’n onafhanklike en objektiewe tribunaal , **in die bepaling van sy regte** en verpligtinge en die ondersoek van enige kriminele saak teen hom .”

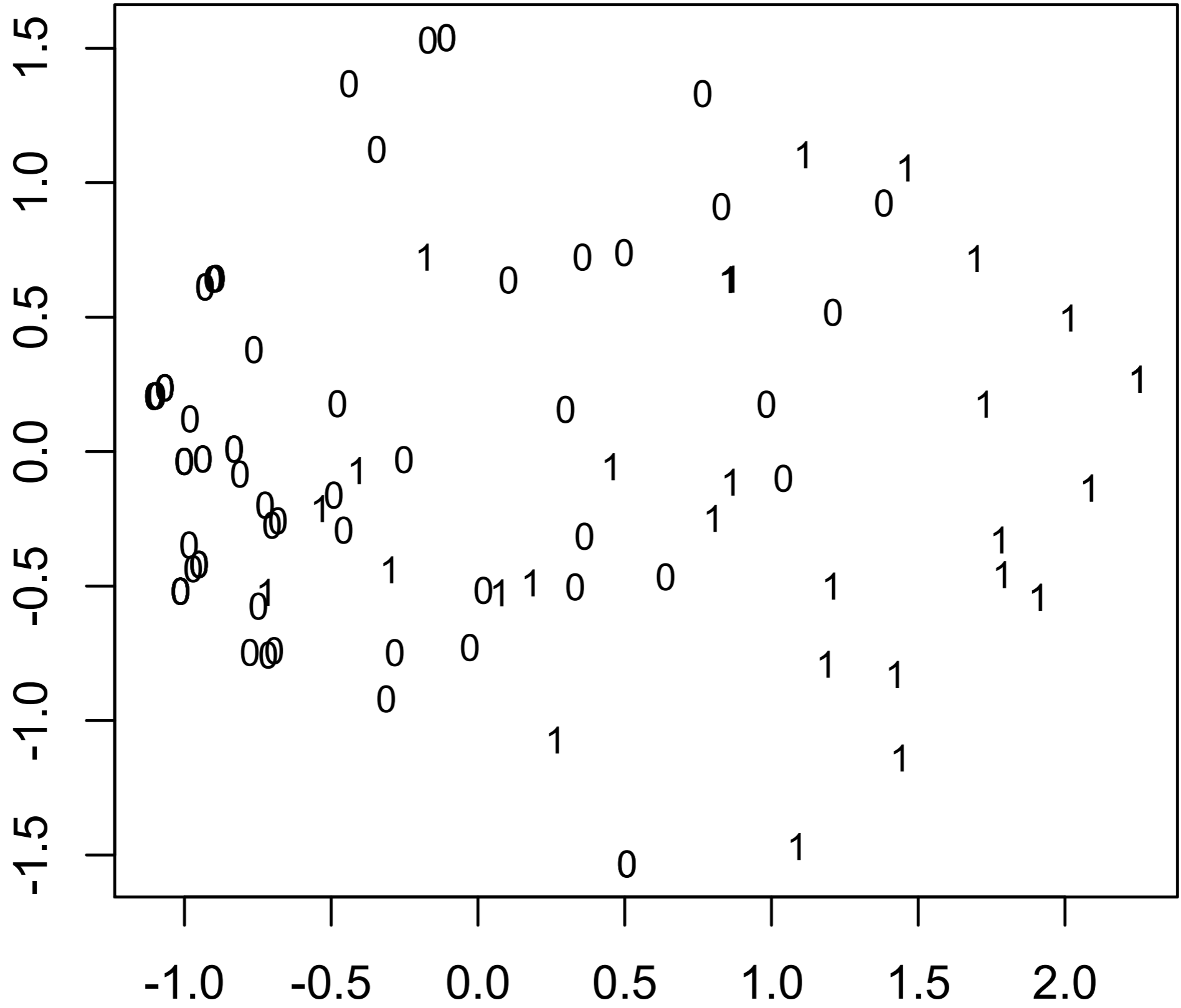
German: “jeder hat **bei der feststellung seiner rechte** und pflichten sowie bei einer gegen ihn erhobenen strafrechtlichen beschuldigung in voller gleichheit anspruch auf ein gerechtes und öffentliches verfahren vor einem unabhängigen und unparteiischen gericht .”

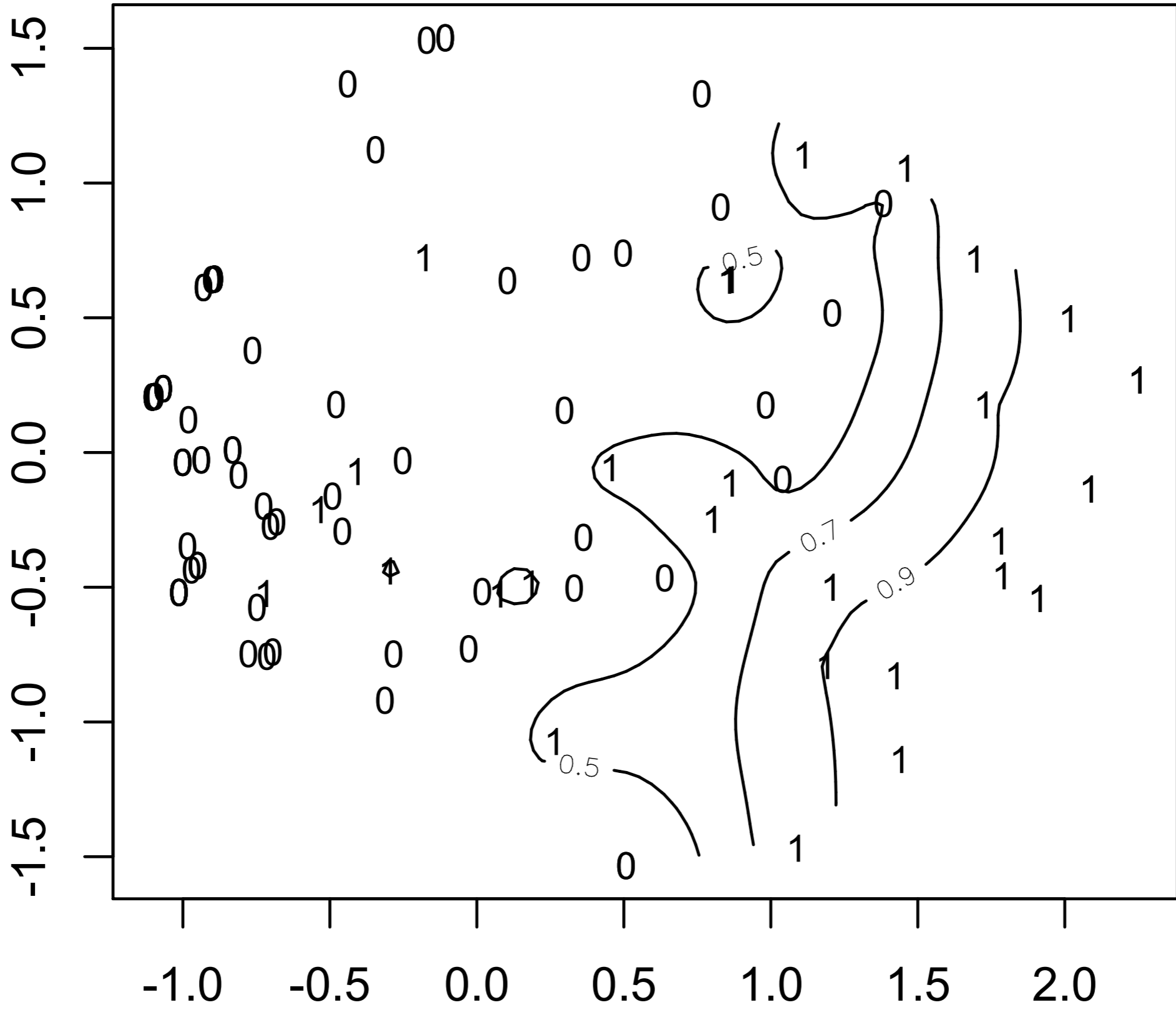
deu	2						
eng	5						
nld	5						
fri	6						
afr	5						
ydd	4						
sco	5						
nno	3						
nob							
swe							
fao	3						
dan							
isl							

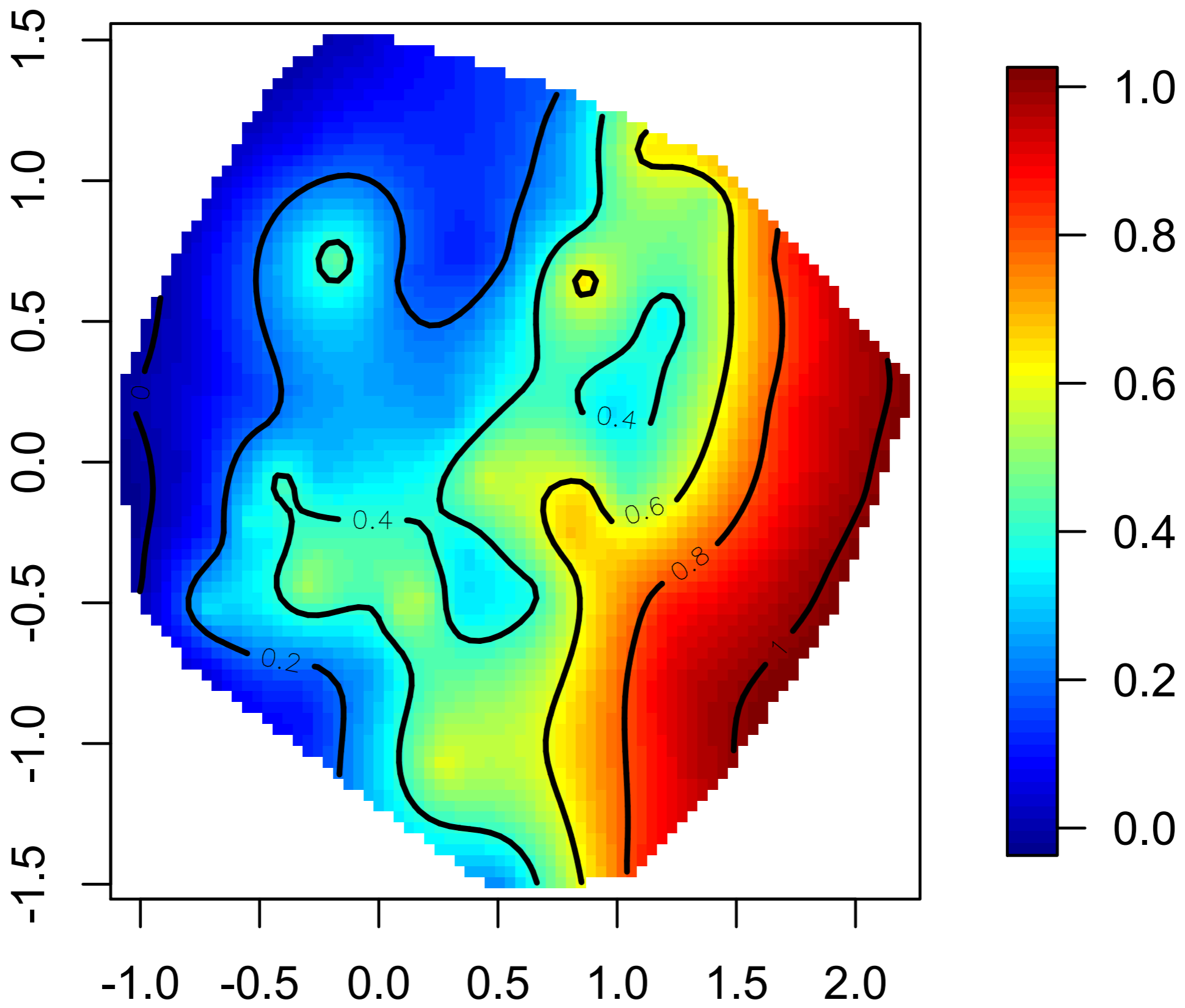


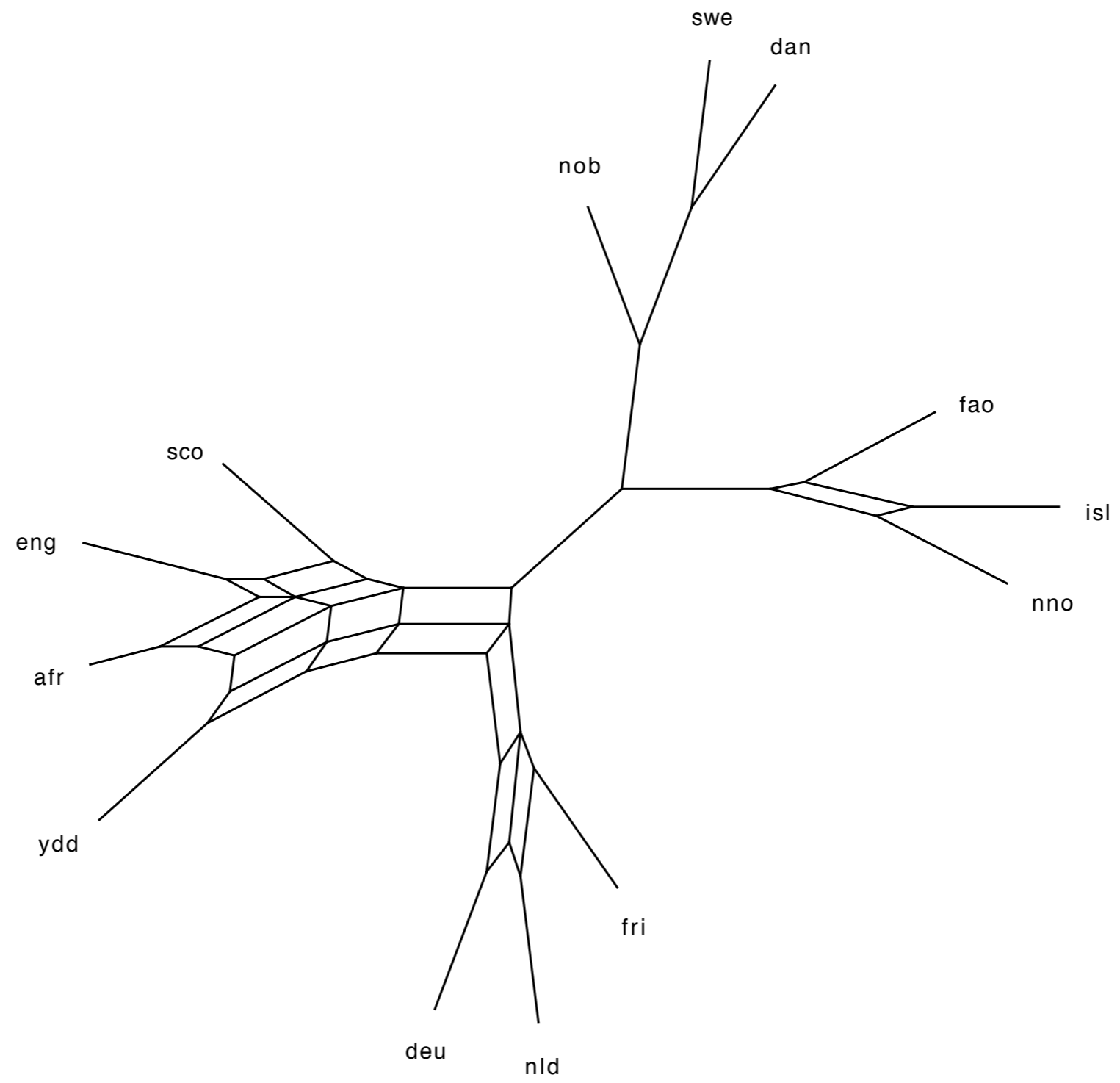
**Consensus network of 4 optimal trees
according to dollo Maximum Parsimony**

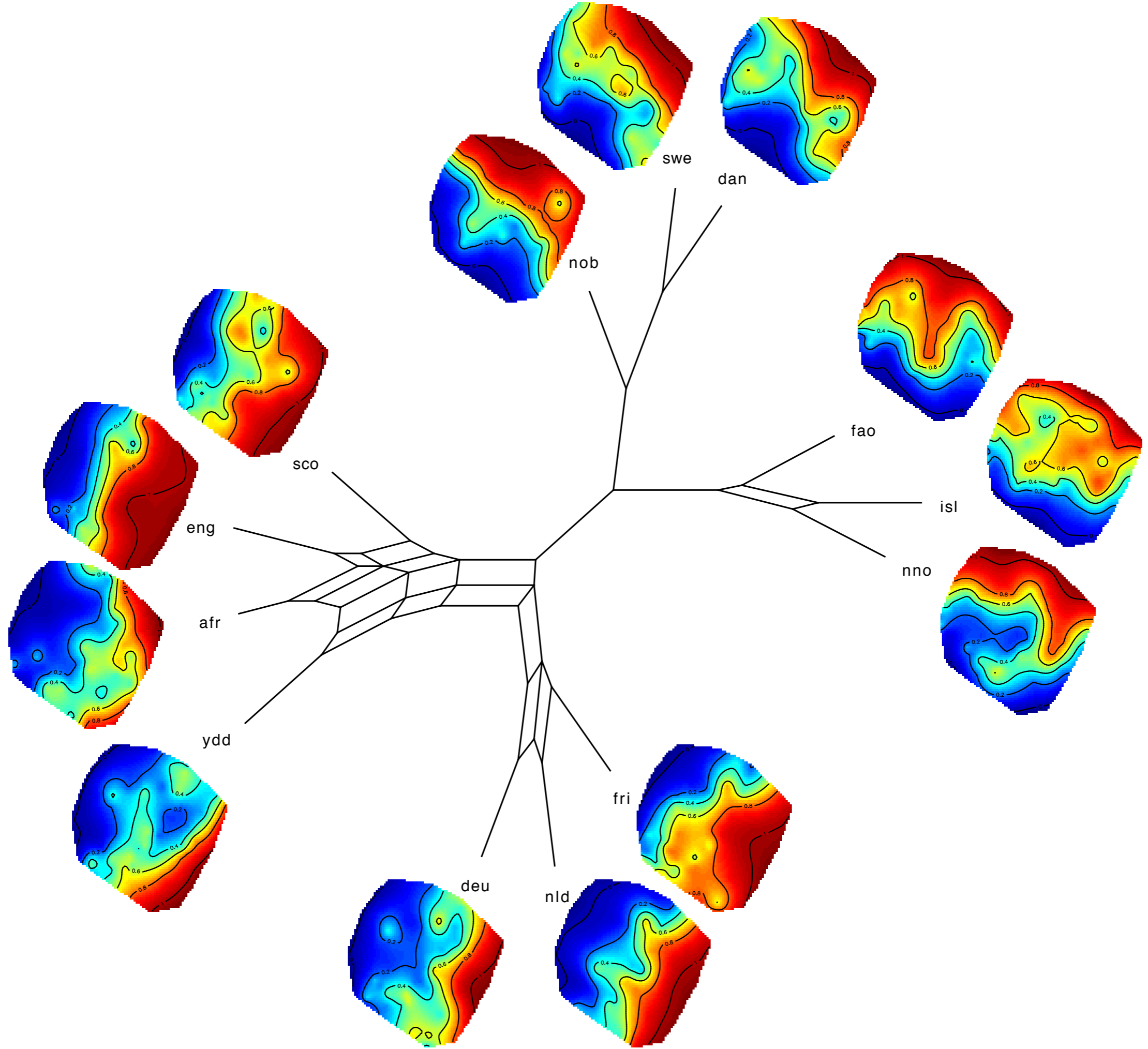












A stylized map of the Pacific region, including parts of Asia, Australia, and the Pacific Islands. The map uses a color palette of light green for land and light blue for water. Several colored dots are scattered across the map: two yellow dots in the northern Pacific, one red dot in the western Pacific, two purple dots in the southern Pacific, and one yellow dot in the southern Pacific near New Zealand. A semi-transparent white rectangular box is centered over the map, containing the title text.

Introducing the *Parallel Text*

Parallel Bible Corpus

- 1169 translations online (soon 1600+)
- 906 different ISO-639/3 codes (soon 1300+)
- In total more than 350 Million tokens
- More than 17 Million different wordforms
- <http://paralleltext.info/data>

Software

- Contact me personally for access

- R-package “qlcMatrix”

<http://cran.r-project.org/web/packages/qlcMatrix/index.html>

<https://github.com/cysouw/qlcMatrix>

- Python library

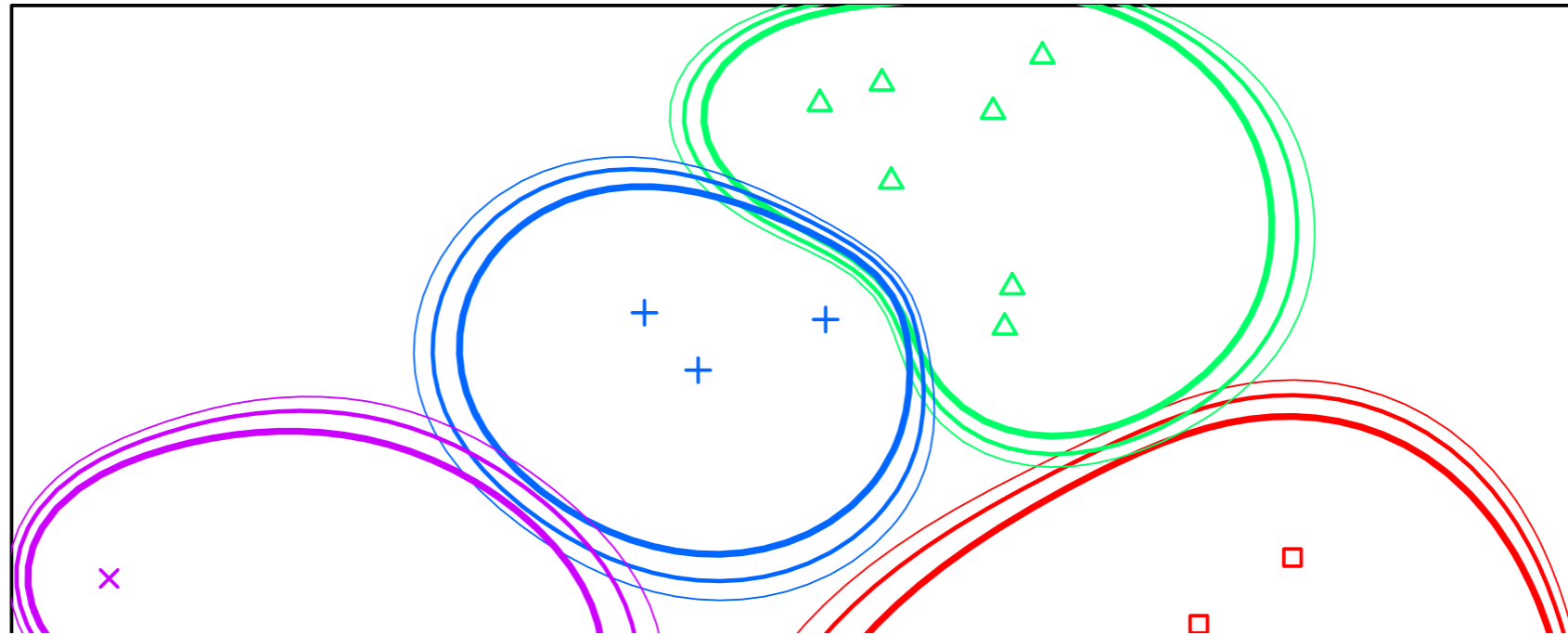
<https://github.com/tmayer/paralleltxtprocessing>

Multiple Alignment

- Small-scale experiment
 - ▶ use *fastalign* for bitext-alignment on all pairs
 - ▶ build multi-text-alignment using graph clustering
- Only for 77 Germanic translations
- New Testament produced almost 100.000 Germanic alignments, which are directly comparable ‘words’

trees and wood

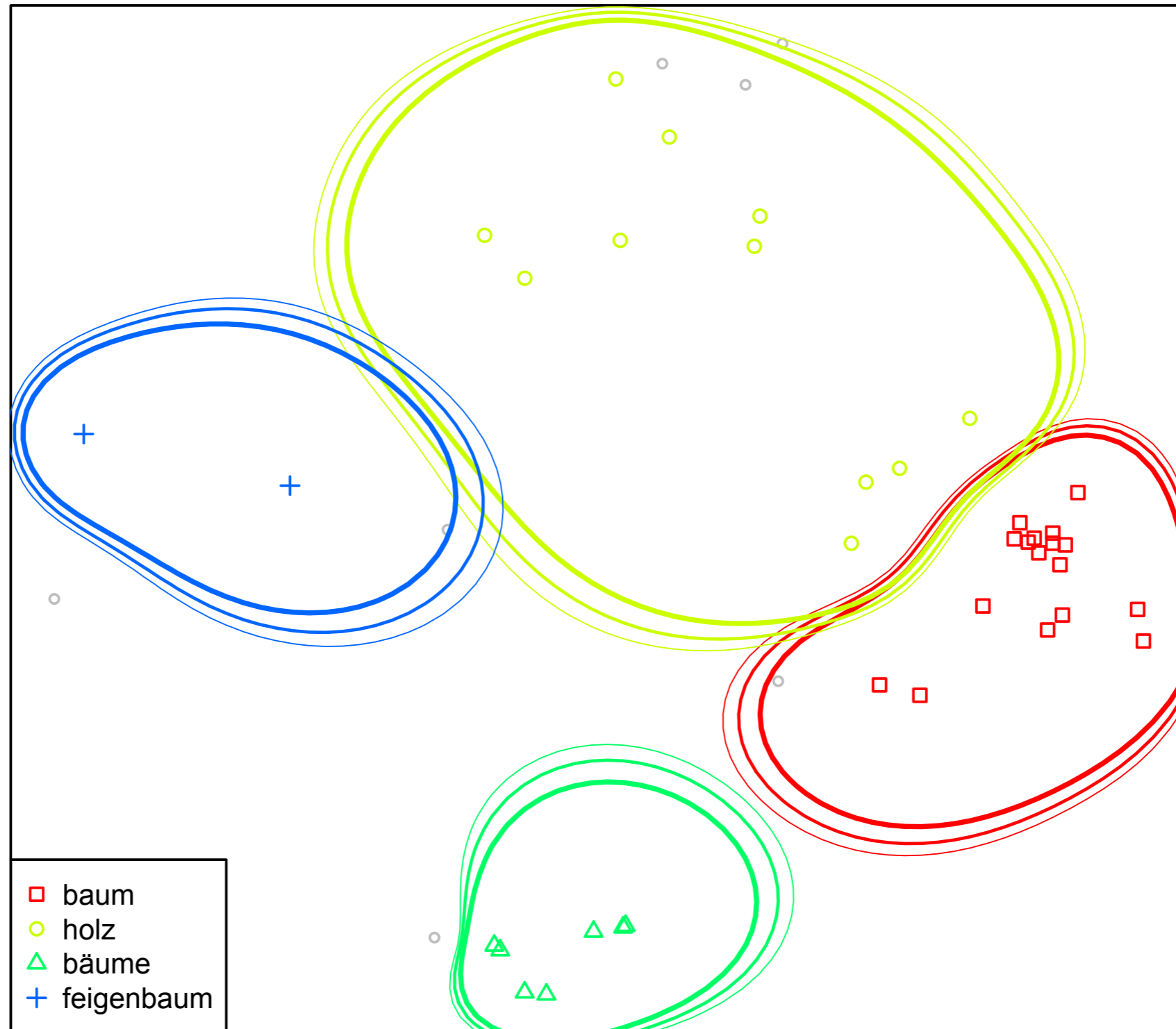
afr-x-bible-1953.txt



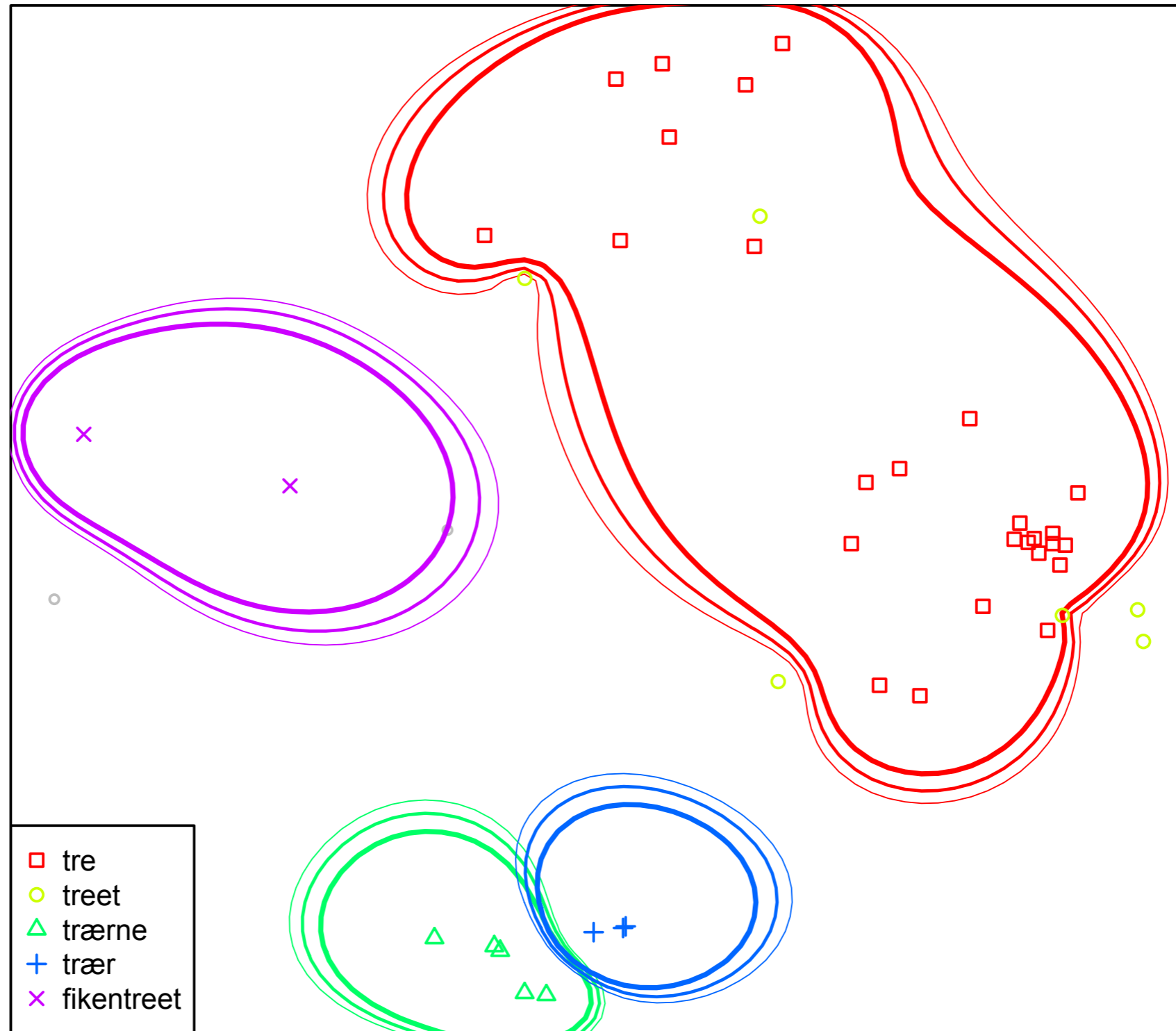
	tree	wood (stuff)	firewood	small forest	large forest
German	<u>Baum</u>	<u>Holz</u>		<u>Wald</u>	
Danish	<u>træ</u>			<u>skov</u>	
French	<u>arbre</u>	<u>bois</u>		<u>forêt</u>	
Spanish	<u>árbol</u>	<u>madera</u>	<u>leña</u>	<u>bosque</u>	<u>selva</u>

Louis Hjelmslev
Prolegomena to a Theory of Language (1963)

deu-x-bible-erben.txt

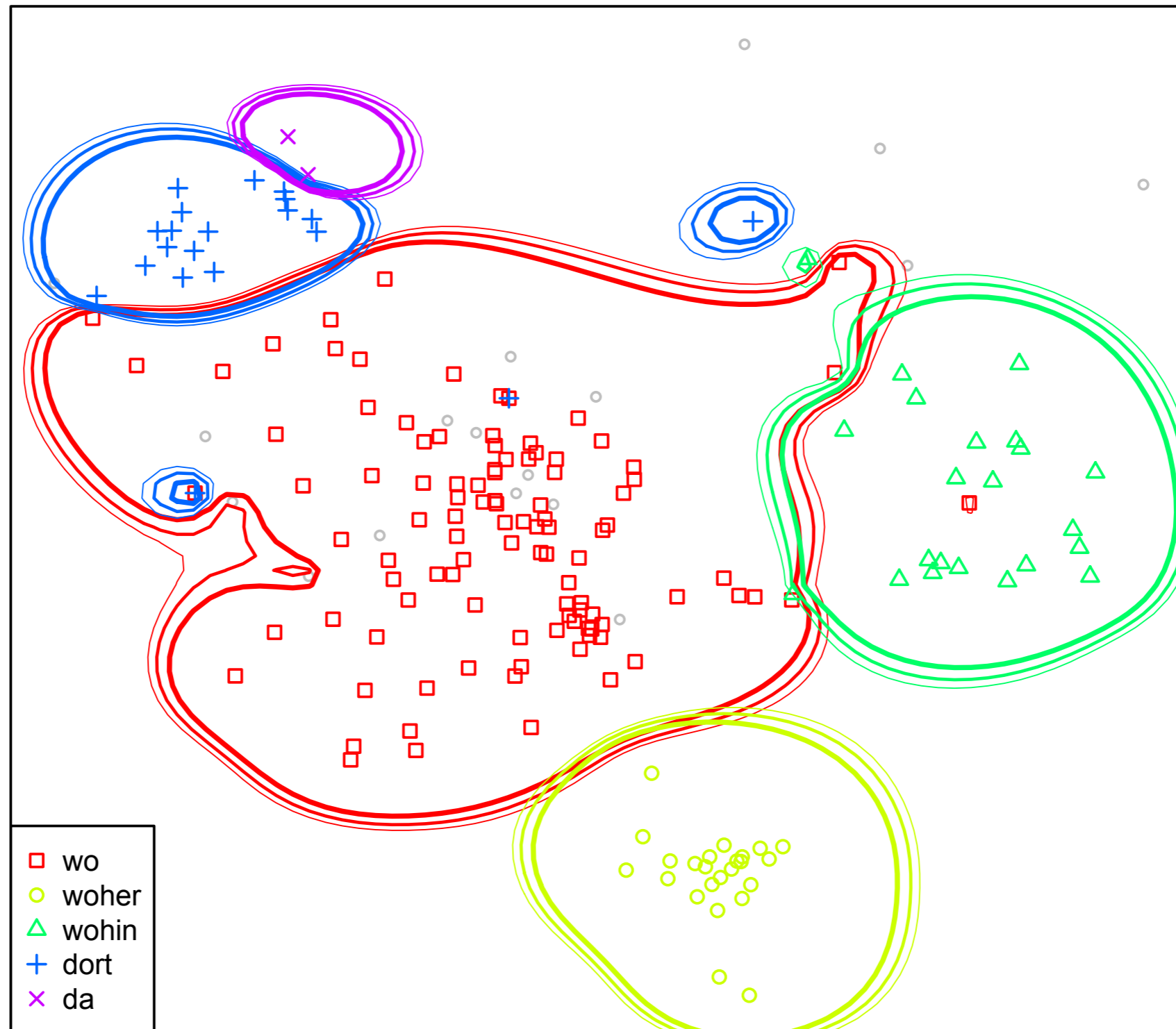


nob-x-bible-2007.txt

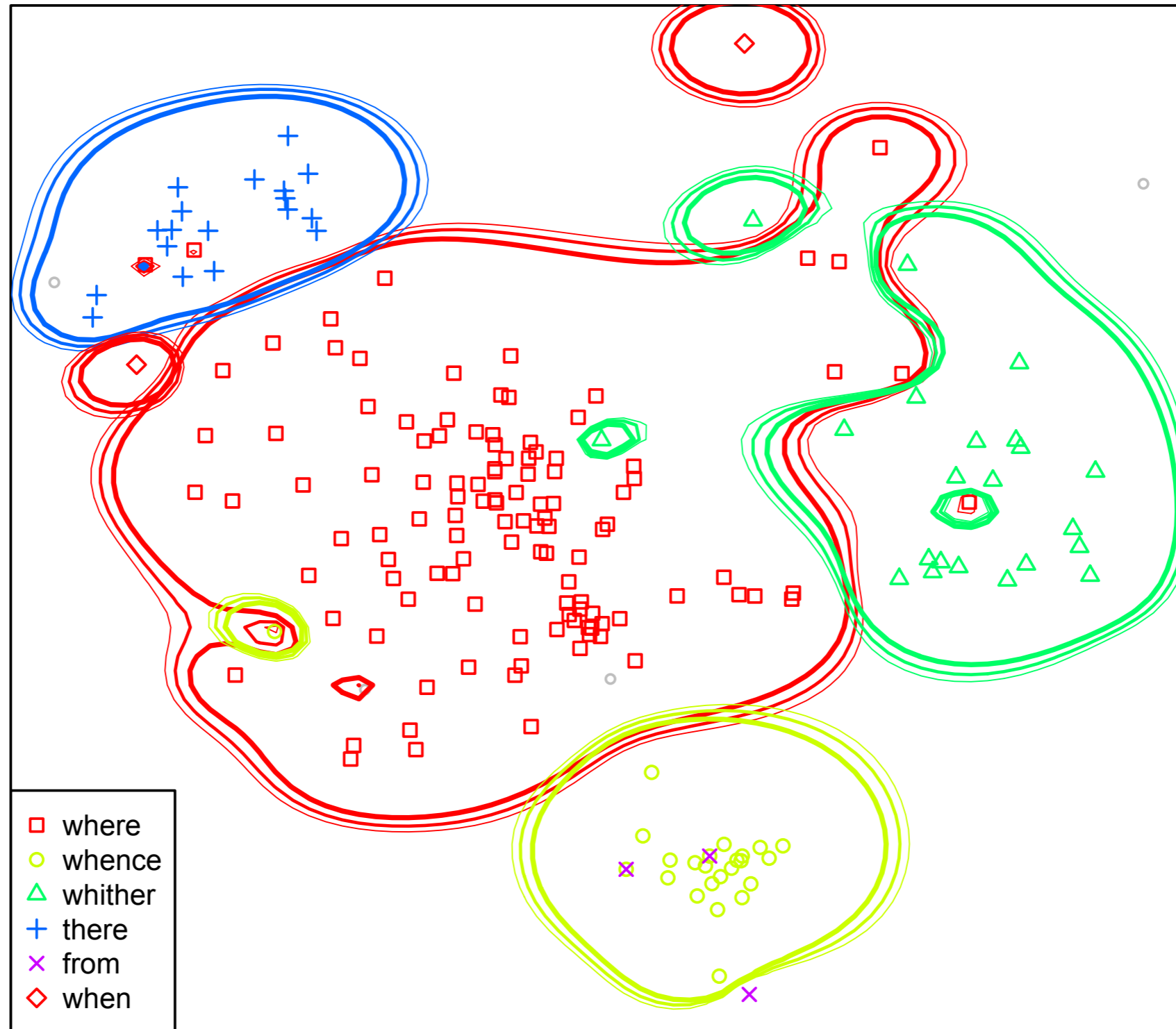


where

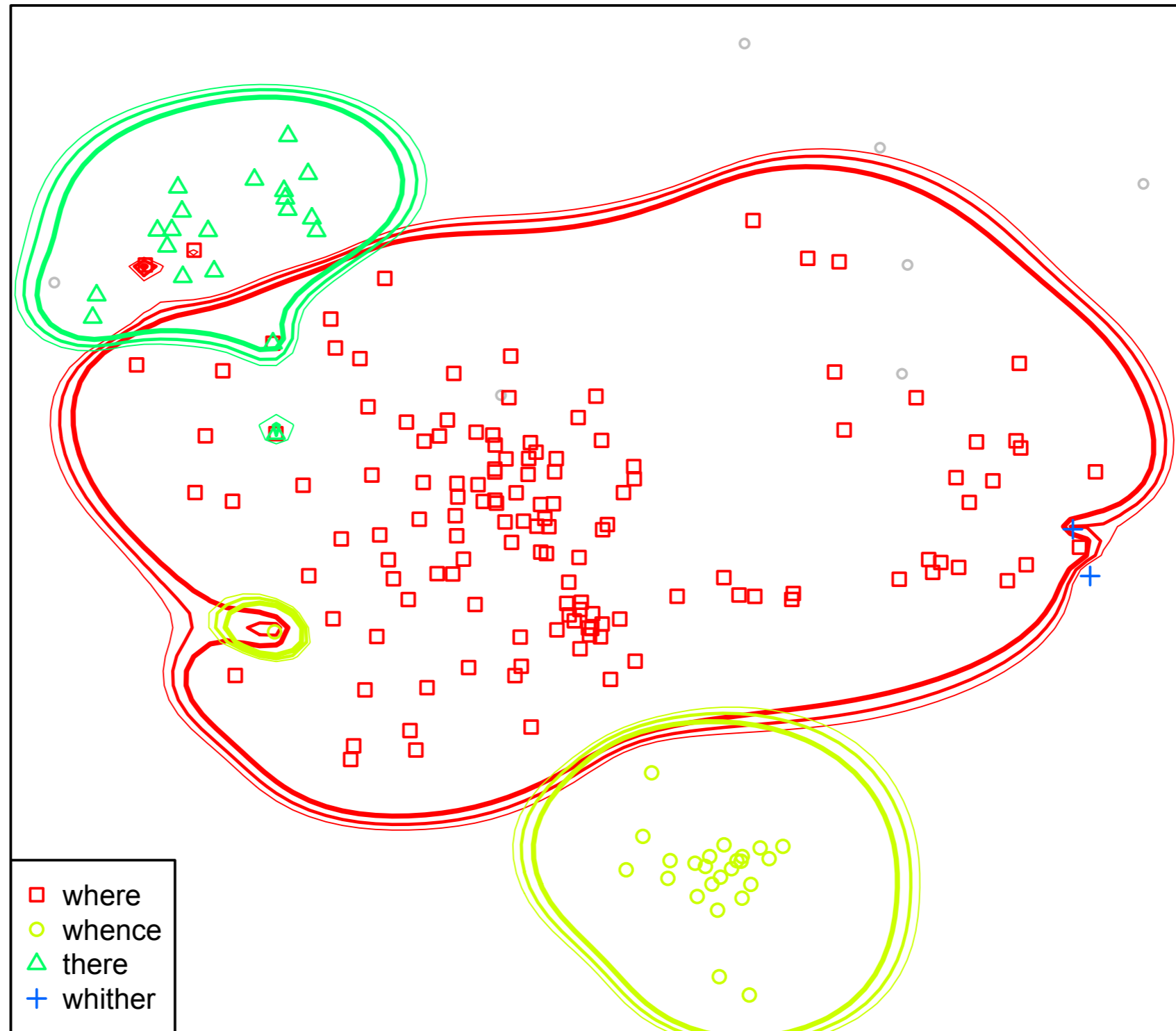
deu-x-bible-pattloch.txt



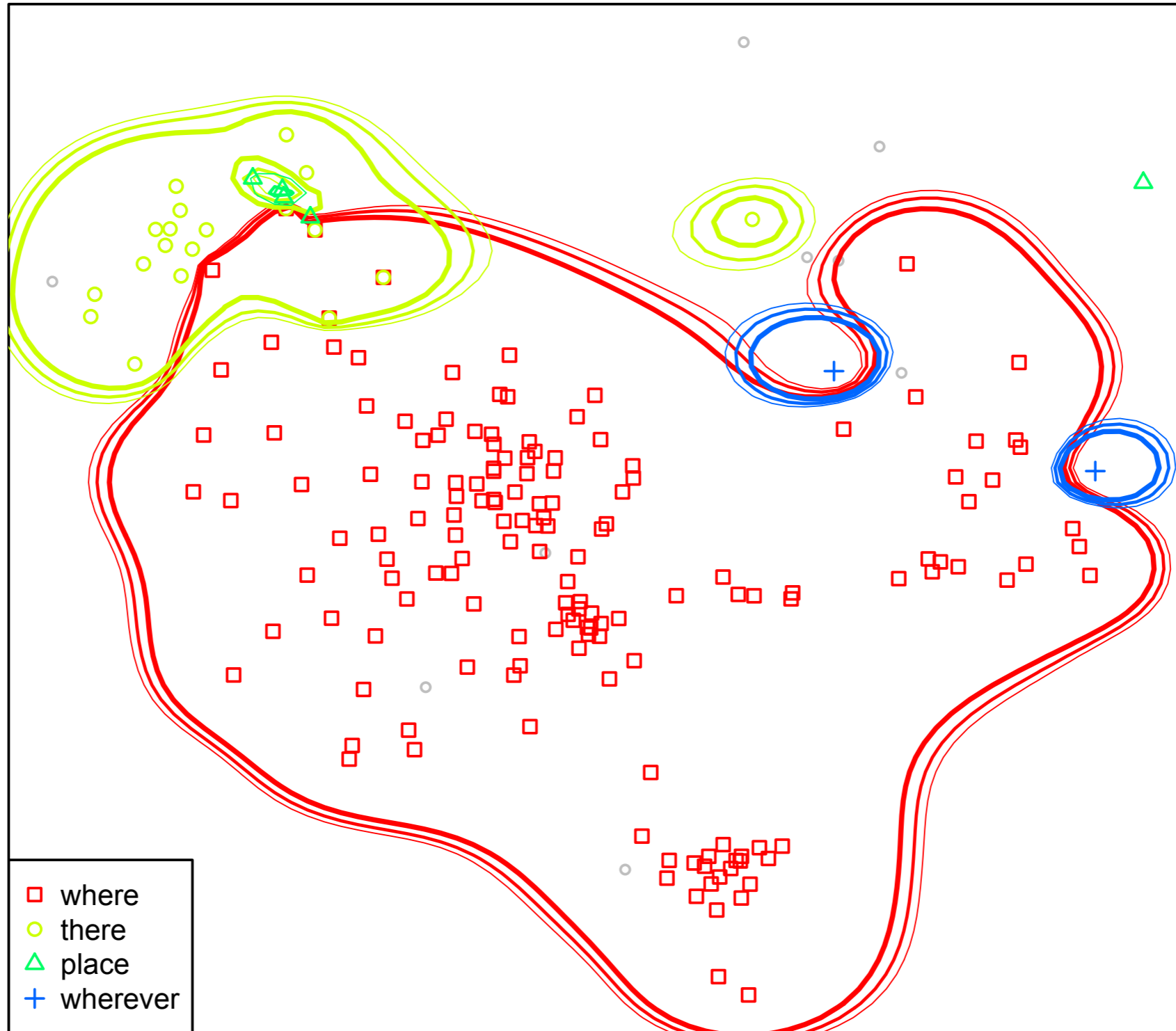
eng-x-bible-kingjames.txt



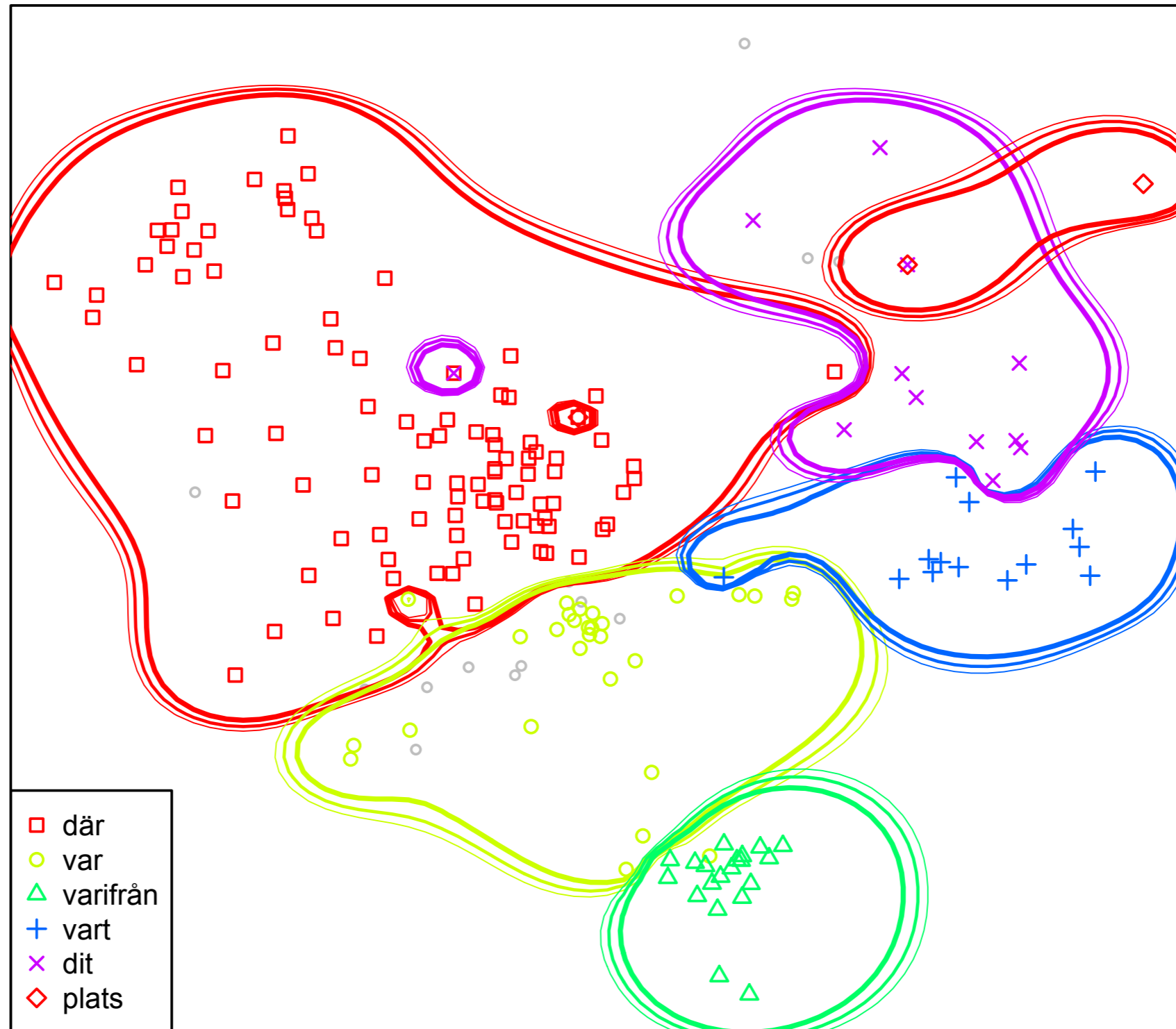
eng-x-bible-darby.txt



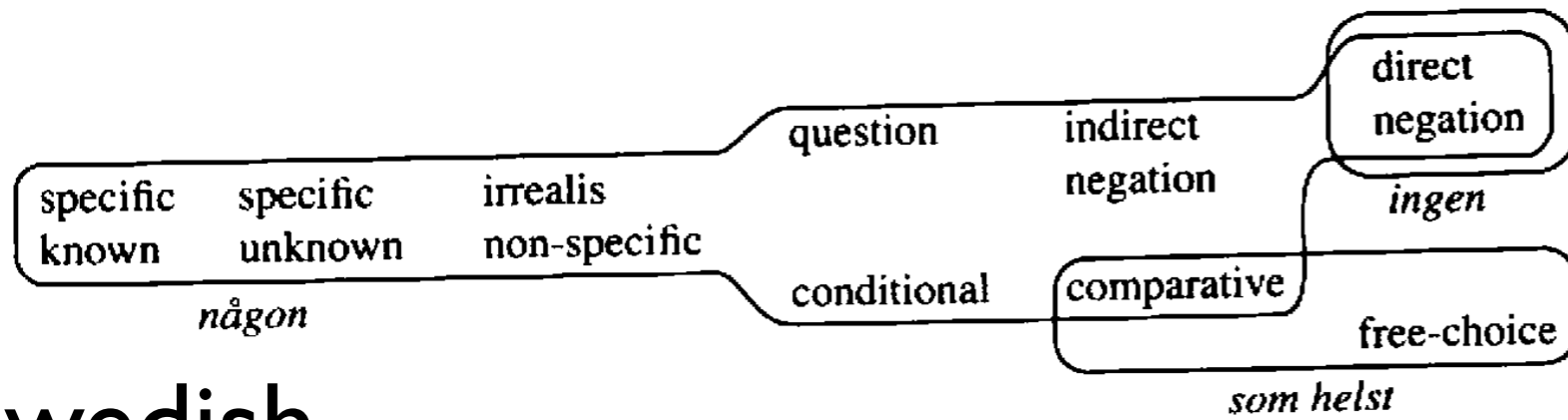
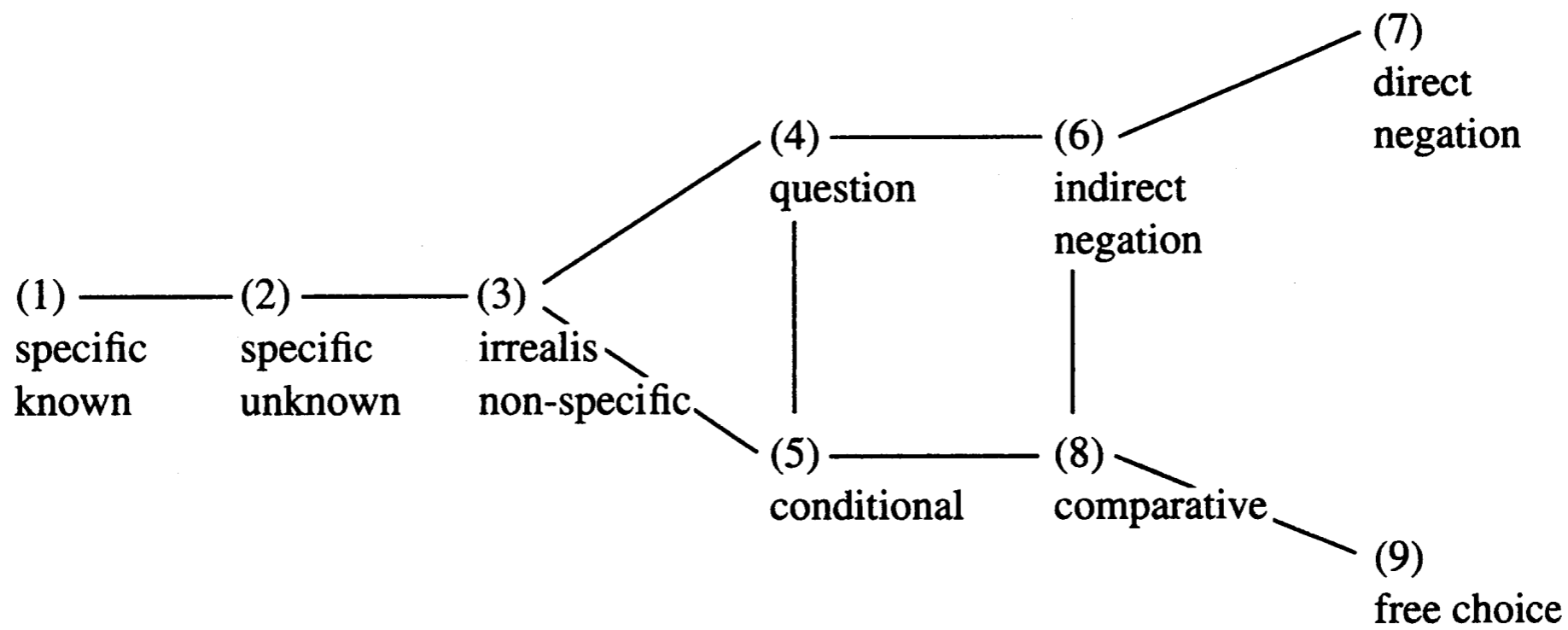
eng-x-bible-treeoflife.txt



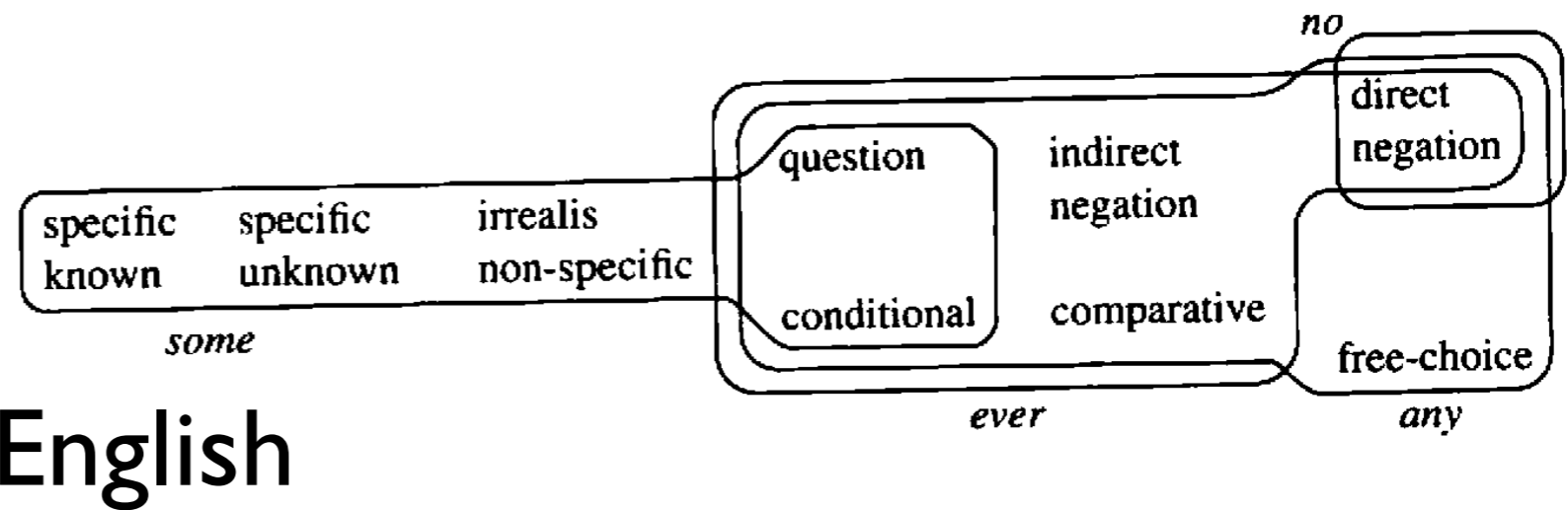
swe-x-bible-folk1998.txt



**Indefinite person
(someone, anyone)**

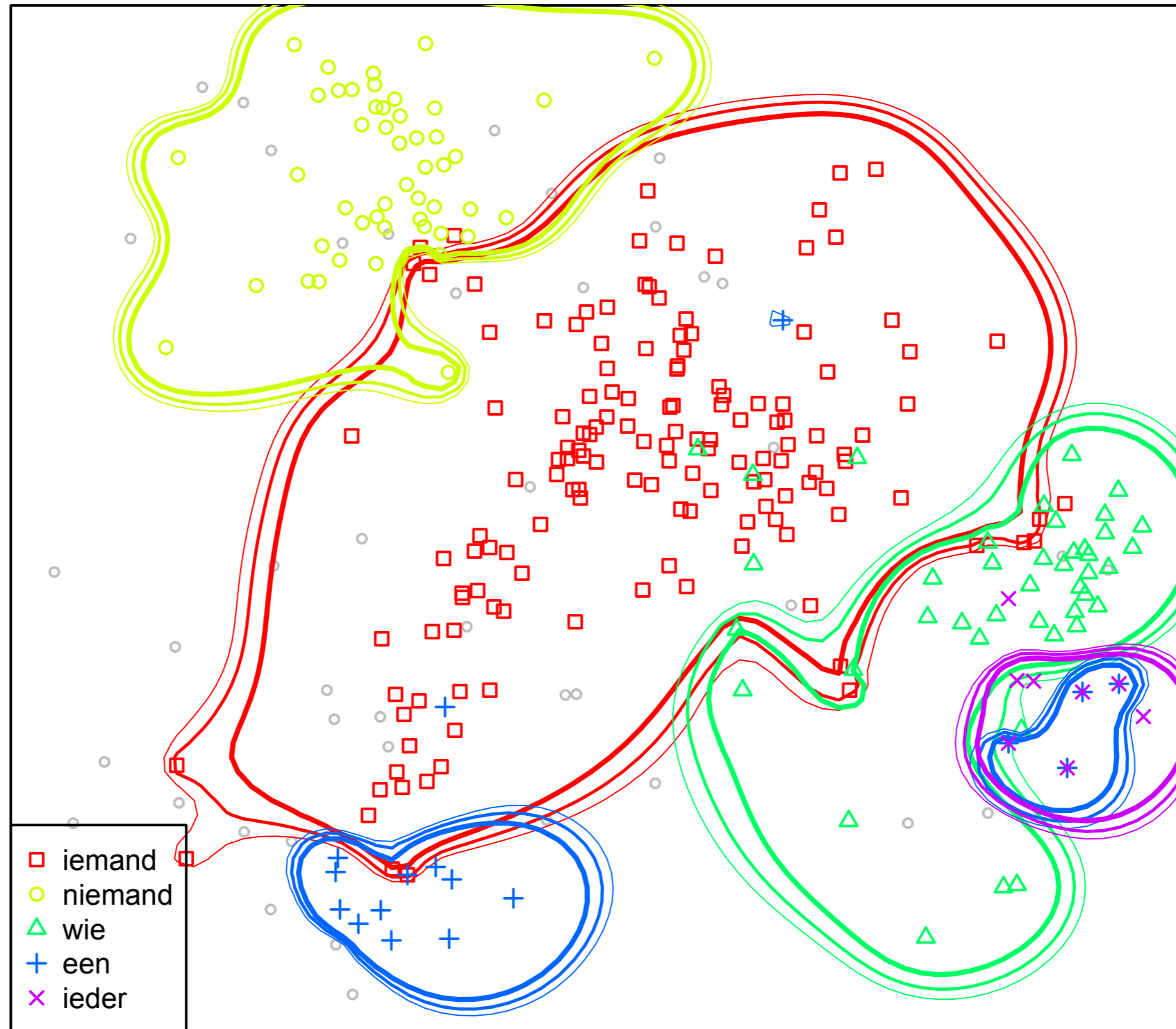


Swedish

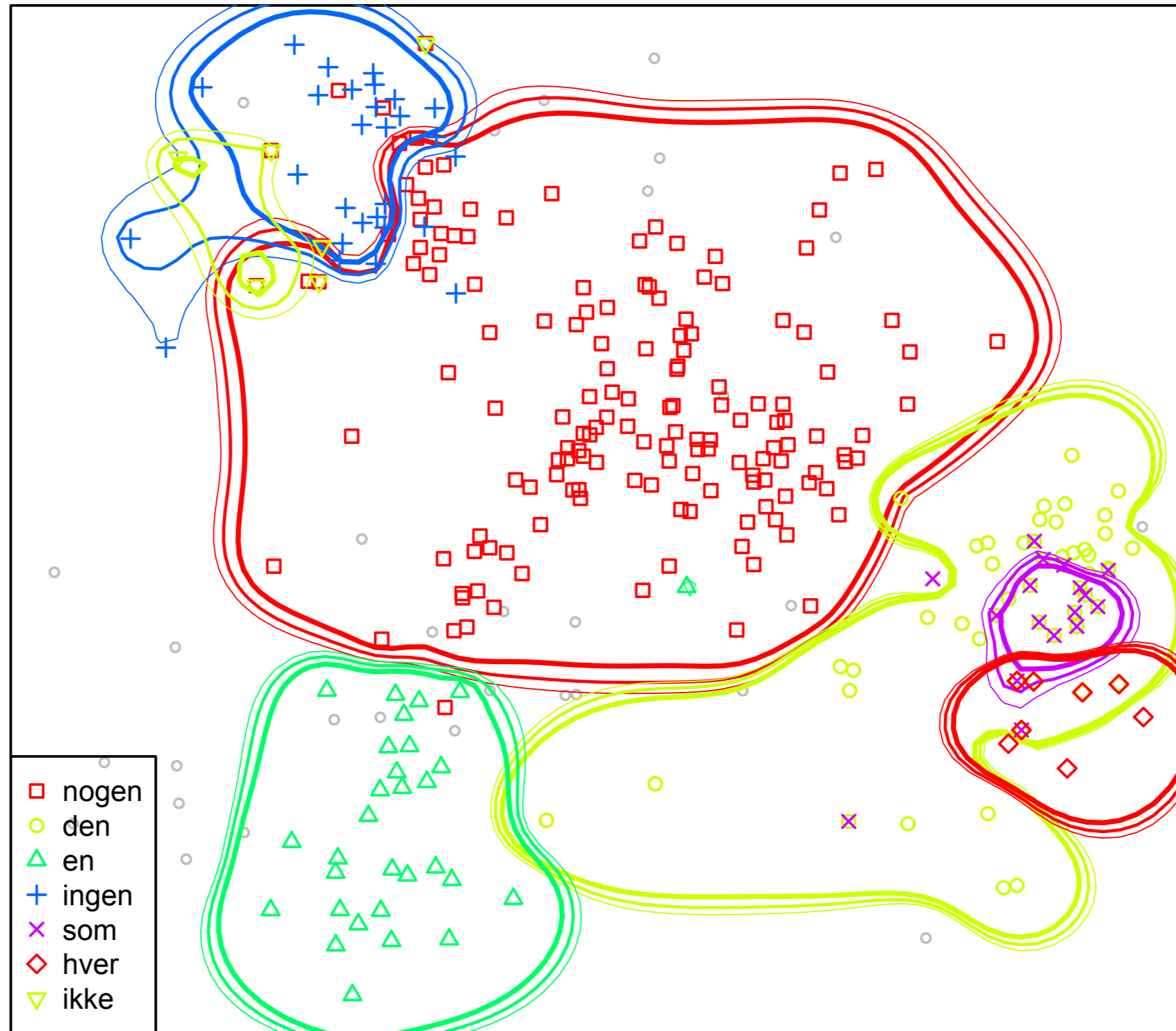


English

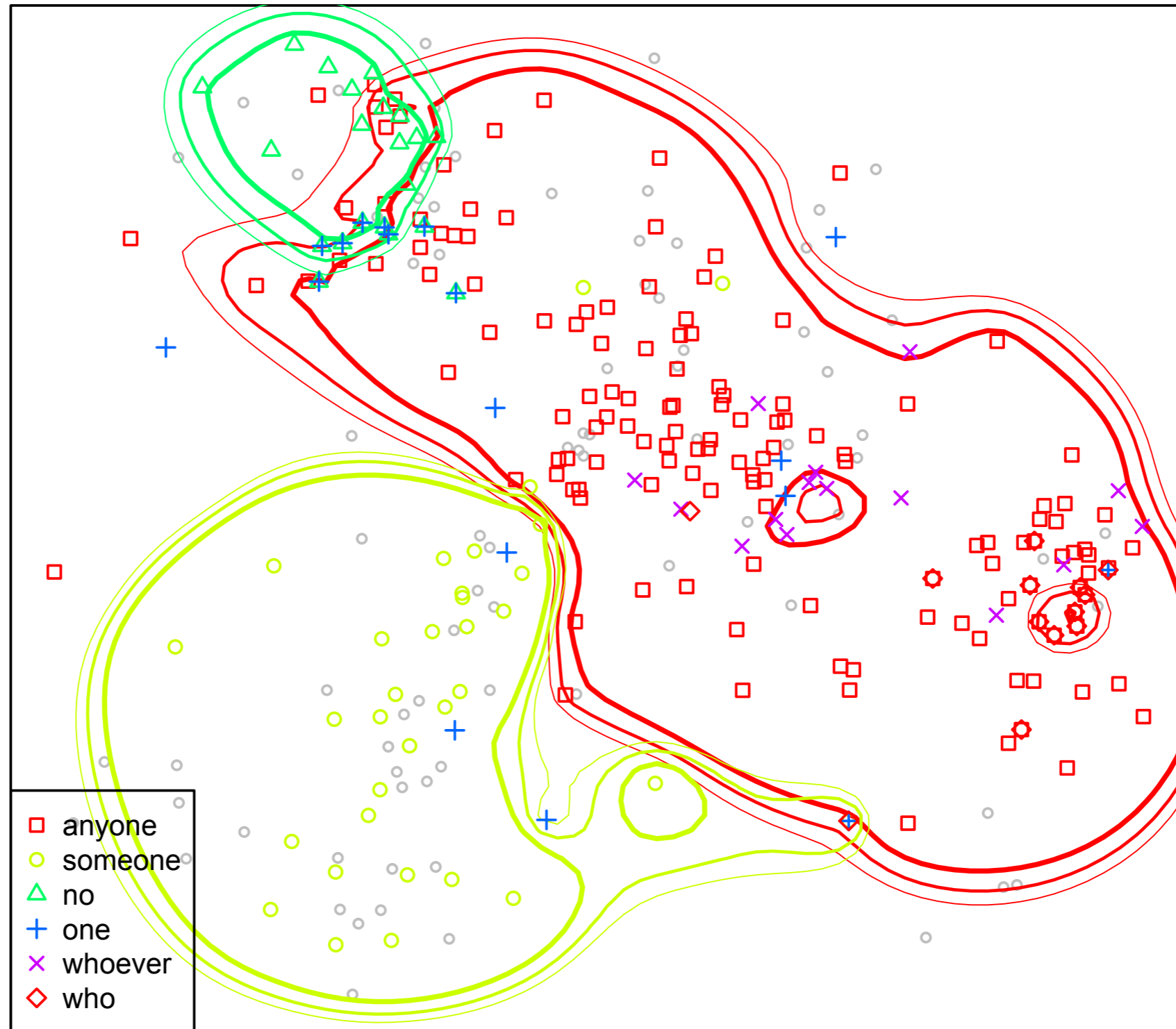
nld-x-bible-1951.txt



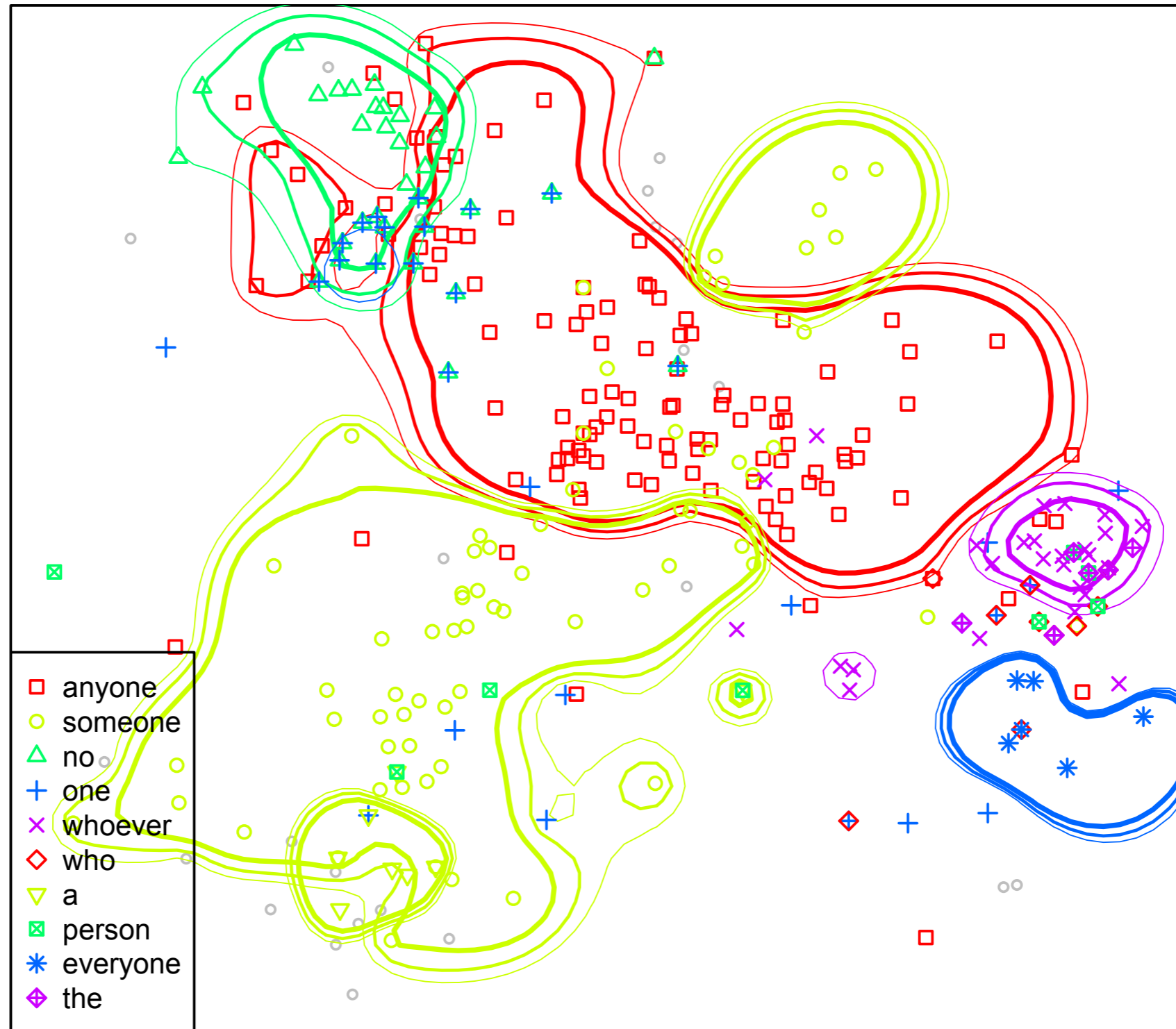
dan-x-bible-1931.txt



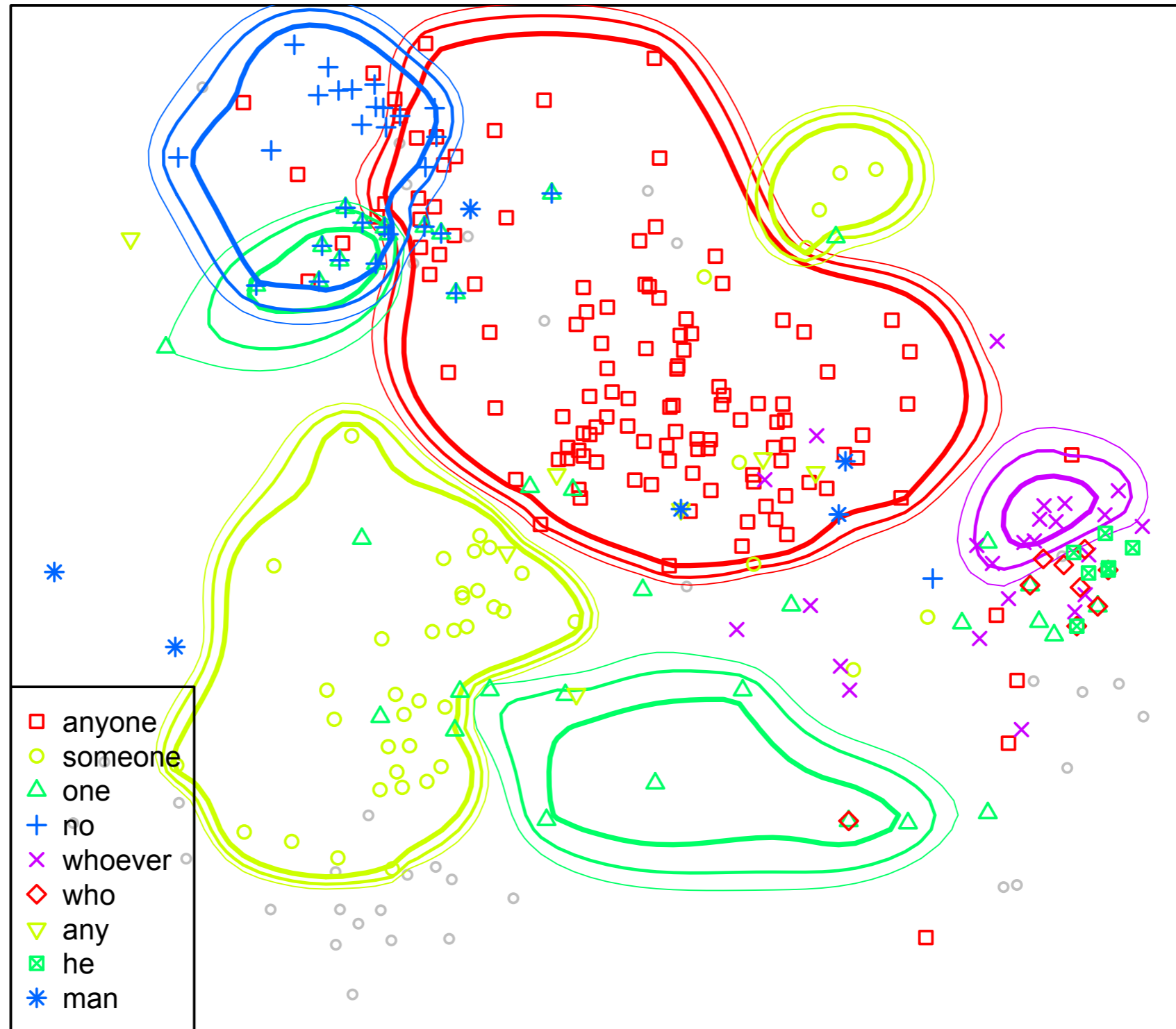
eng-x-bible-books.txt



eng-x-bible-new2007.txt



eng-x-bible-treeoflife.txt



Multialignment of Sounds

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter	178 'vɪntə ^h ʁ	ʁ = kontinuierlich ɔ bereit velarisieren
* 2	fliegen	56 'flaɪ→ə ^h θə	'fliegen die', Sequenzierung unklar - kein geminiertes [t]
3	Blätter	23 'blɛ: dɛ ^h ʁ	
4	Luft	103 lu ^h ft ^h	ʁ = kontinuierlich
5	hört	89 hɪ ^h t ^h	ʁ = kontinuierlich
6	gleich	78 k ₊ lɛ ₊ ɪk ₊	folgt P
7	schneien	130 ʃnɛɪən	
8	Wetter	174 /	statt dem 'vɪtəʀɔŋə
9	tu	151 dɛ→v	

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	ʏf:
Breddin	ɒ:ph
Brelingen	ɑfβə
Bremscheid	ɒ':ph̥ə
...	...

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

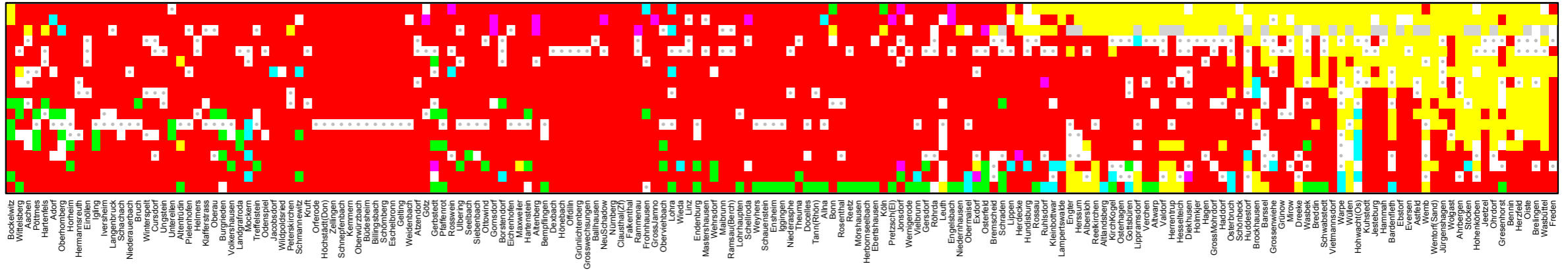
COLUMNID	1	2	3	4	5	6	7	8	9
STANDARD	Au	g	e	n	b	l	i	ck	(e)
Adorf	a→u	ɸ	-	ŋ	b	l	ɛ	k ^h	-
Ahrbergen	ə→ō	ɣ	ə	m	b	l	ɪ	k'	-
Albersloh	a→u	-	-	m	β	l	ɪ	k	-
Allna	ɔɪ	-	-	-	p	l	æ ^c	x	-
Altenberg	ɤɣ	ɸ	ã	-	b	l	ɪ	k	-
Altentrüdin	a→u	ɸ	ə	-	p	l	ɪ	g	-
Altlandsberg	a'→u	ɣ	-	ŋ	b	l	ɪ	k̄x	-
Altwarp	ōu	-	-	ŋ	b	l	ɪ	k	-
Astfeld	ʊɪ	ɣ	ə	m	b	l	ɪ	k ^h	ə
Ballhausen	a→u	ɣ	-	ŋ	p	l	ɪ	k	-
Bardenfleth	oɪ	g	-	ŋ	b	l	ɛ	k̄x ₊	-
Barssel	oɪ	g	-	ŋ	p	l	ɪ	k ₊	-
Bempflingen	a→u	g	ə	-	b	l	ɪ	c ^h	-
Bennin	oɪ	-	-	ŋ	b	l	ɪ	x	-
Billingsbach	aɪ	x	ə	-	p	l	ɪ	k̄x	-
Bockelwitz	a→u	ɣ	-	ŋ	b	l	ɪ	k	-
Borstendorf	ɔɪ	ɣ̄x	-	ŋ	p	l	ɛ	k	-

github.com/digitallinguist/msa-editor

Correspondence Sets

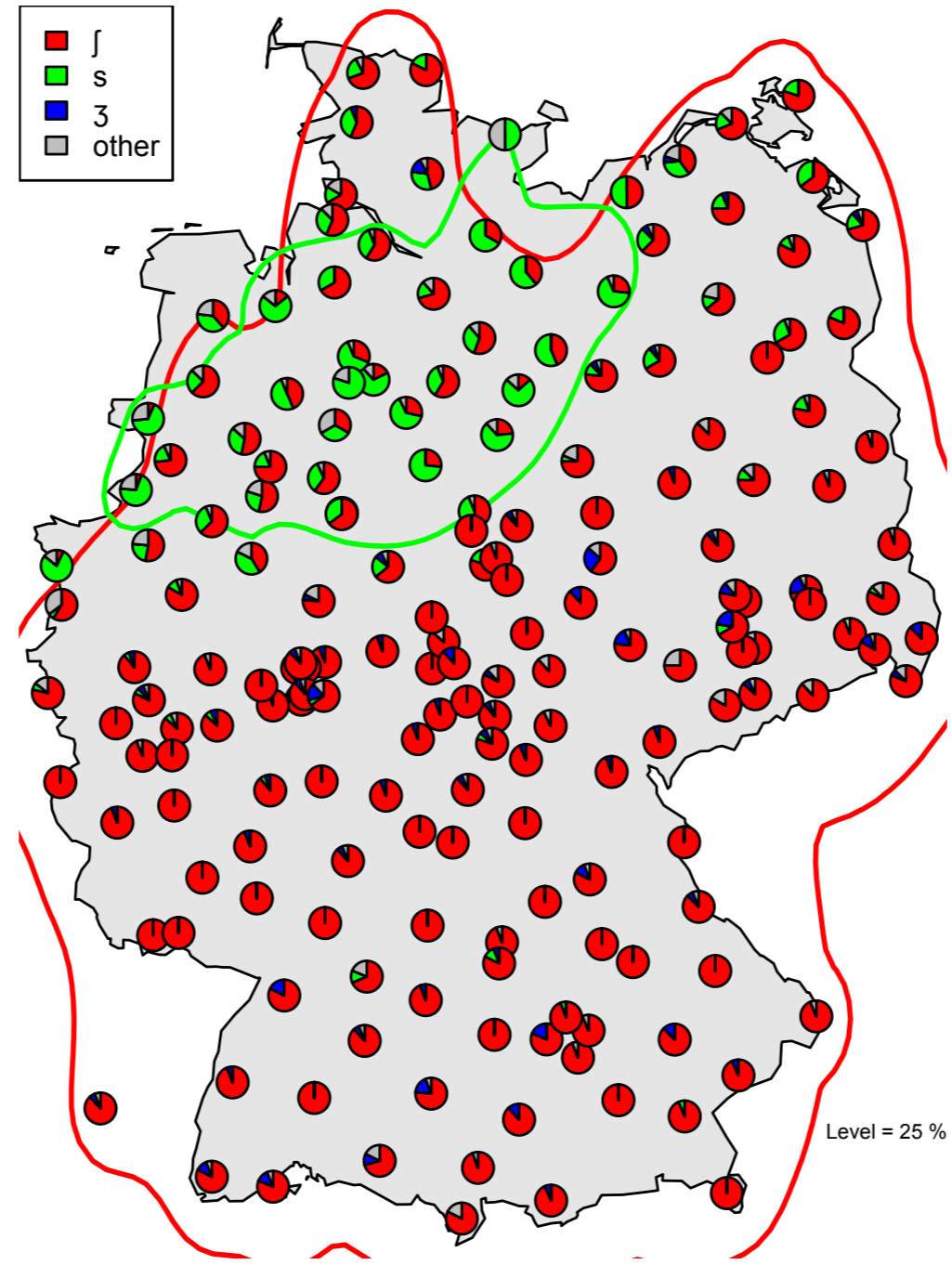
- more than 700 columns of aligned segments (“correspondences”)
- Comparative-historical linguistics uses clusters of correspondences (“correspondence sets”)
- Automatic clustering of columns is a good start, but needs correction
- Visualisations in R
github.com/cysouw/qlcVisualize

700 Wurst 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 Fleisch 4 sch
 587 Tisch 3 sch
 156 Dreschen 7 sch



┘
 s
 3
 6
 -
 j:
 other
 NA

Correspondences "sch"



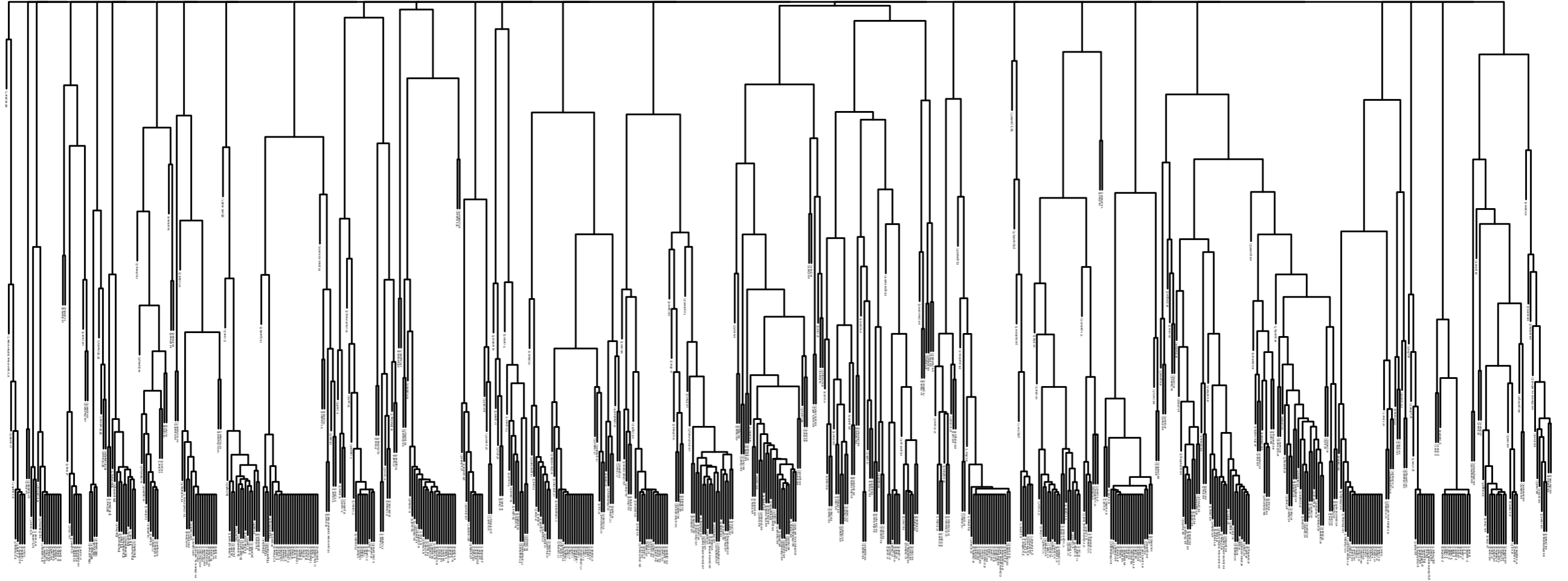
Bulgarian dialect data

- 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



v, t, (when:1)
 't/ (this:1)
 't/ (yours:1)
 t/ t (such:1)
 t (Saturday:5)
 t/- (fear:3)
 t/- (the middle:7)
 t (old man:2)

t/ti (son-in-law; brother-in-law:3)

t/ti (ten)

't/'tj (dark:1)
't/'tj (thin:1)

l/- (wool:3)
l/i (yellow:3)

l/i (deep:3)

l/- (apple:5)

l/'l (enter:3)

l/'lj (rooster:5)

l (down:4)
l (hungry:2)
l (head:2)
l (pay:3)
l/'l (white:3)

l/- (ox:3)
l/'lj (whole:3)
l/- (has come:6)

'l/'li (easy:1)
'l/'li (lentils:1)

li/'l (key:2)
li/'j (Sunday:5)

'li/'l (bread:2)
l/'li (the milk:2)

'li/'l (iron:3)

l/'li (salt:4)
l/'li (apples:5)

k/c (tong
k/c (bar

k/- (today:6)

k/c (wolf:5)
k/c (Monday:10)

k/c (person:5)
k/'d (where:1)
k/'c (key:1)

k/'k (hand:3)
k/'k (river:3)
k/'k (such:3)

k/c (sand:5)
k/c (apple:7)
k/'k (which:1)

k/'k (the milk:4)
k (blood:1)
k/'k (horse:1)
k/'k (stone:1)
k (deep:7)
k/- (outside:4)
k (wait:3)
k/- (Friday:5)
k/'k (thin:4)

k/- (today:6)

k/c (wolf:5)
k/c (Monday:10)

k/c (person:5)
k/'d (where:1)
k/'c (key:1)

k/'k (hand:3)
k/'k (river:3)
k/'k (such:3)

k/c (sand:5)
k/c (apple:7)
k/'k (which:1)

k/'k (the milk:4)
k (blood:1)
k/'k (horse:1)
k/'k (stone:1)
k (deep:7)
k/- (outside:4)
k (wait:3)
k/- (Friday:5)
k/'k (thin:4)

p/b (bake:2)
p/b (pocket:3)

p/b (bread:4)

p/pi (rooster:1)

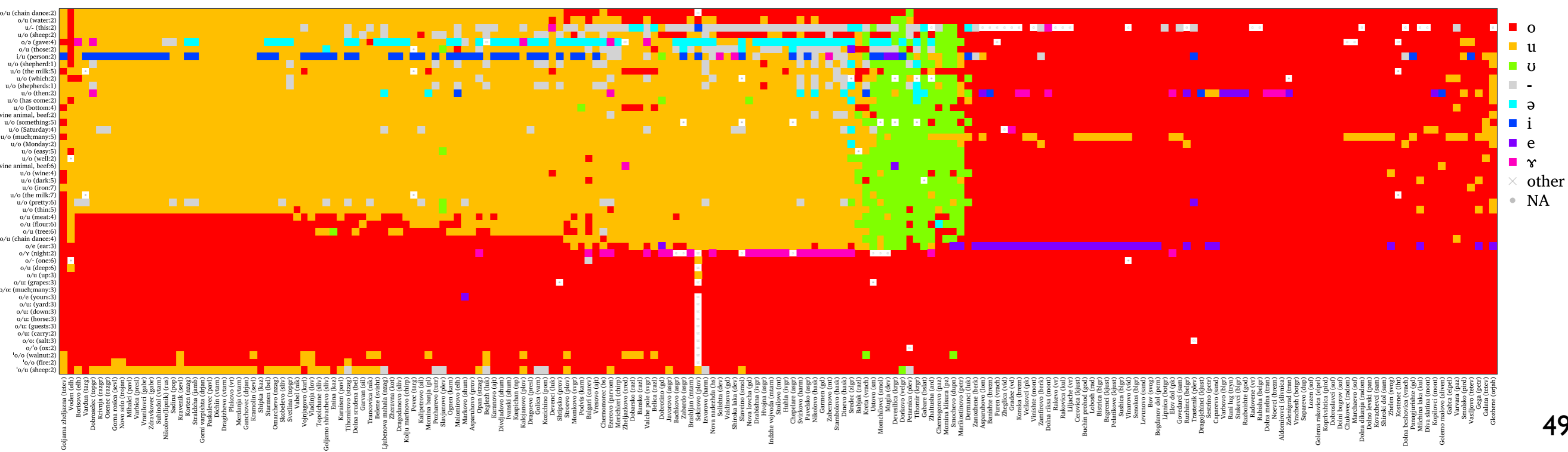
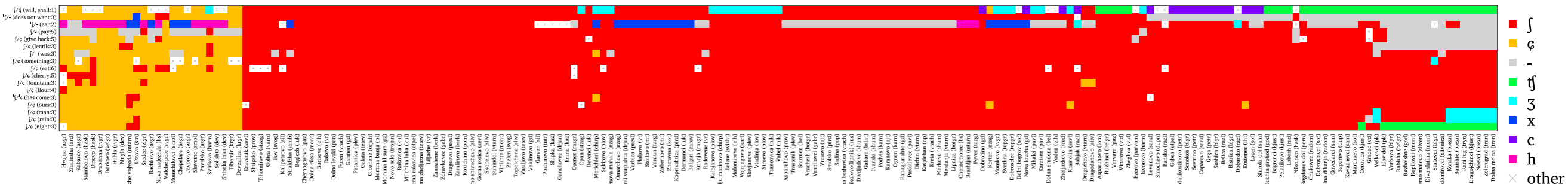
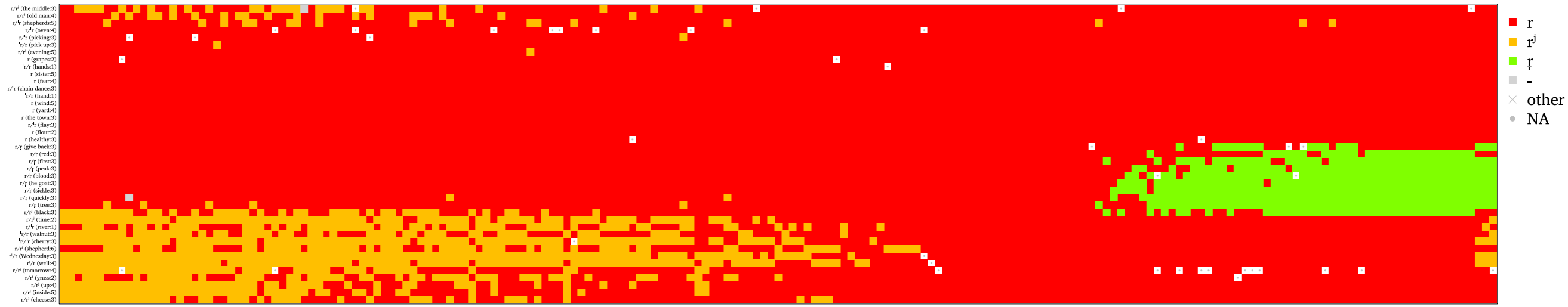
p/f (sickle:5)
'p/p (road:1)
p/p (first:1)
'p/p (he-goat:1)
p/p (pay:1)
p (Monday:1)

'pi/'p (sand:

'pi/'pj (ash:1)
'pi/'pj (Friday:1)

x/f (dry:3)

x/- (fear:6)
x/- (saw:5)

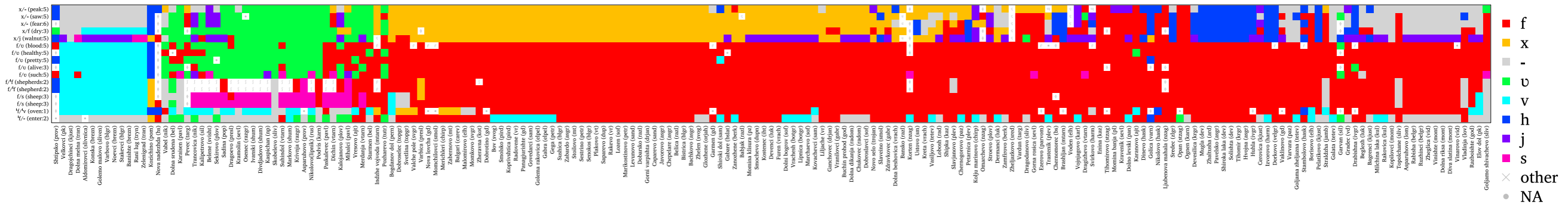


Proto-Bulgarian

	Labial	Dental	Palatal	Velar
Plosive	p b	t d		k g
Affricate		\widehat{ts}	$\widehat{tʃ}$	
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		\widehat{ts} \widehat{dz}	$\widehat{tʃ}$	
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

- Multialignments are an ideal data structure to compare language variation and change
- Extracting Multialignments from large data sets is not trivial, and bound to produce errors
- Multialignments can be profitably discussed and collectively improved