

# **Introducing the *Parallel Text Corpus***

*Michael Cysouw*  
*Philipps-Universität Marburg*

# Parallel Text Corpus

- Universal Declaration of Human Rights
- Pamphlets of the Jehovas Witnesses
- Bible
- currently a private github repo: but I invite everybody to become a member

# 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>

***Demo***

yor-x-bible.txt  
 vie-x-bible-2002.txt  
 npl-x-bible.txt  
 nhy-x-bible.txt  
 nhw-x-bible.txt  
 nhi-x-bible.txt  
 nhg-x-bible.txt  
 nhe-x-bible.txt  
 ngu-x-bible.txt  
 ncl-x-bible.txt  
 ncj-x-bible.txt  
 nch-x-bible.txt  
 haw-x-bible.txt  
 hat-x-bible-1985.txt  
 -x-bible-common.txt  
 azz-x-bible.txt

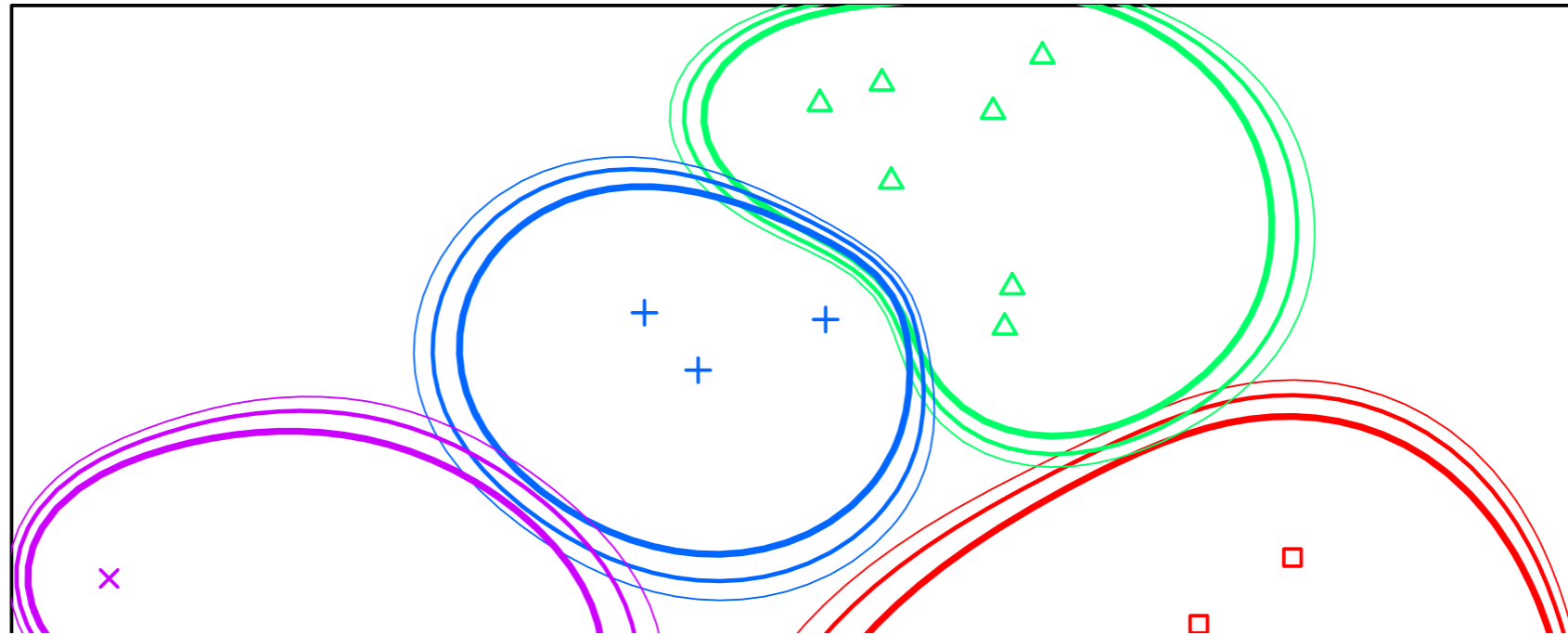
nigbati	herodu	oba	si	gbọ	ara	rẹ	kò	lele	ati	gbogbo	awon	ará	jerusalemu	pelu	rẹ				
nghe	vậy	vua	hê-rôt	và	cà	giê-ru-sa-lêm	đều	bói	rói										
ijkuak	n	rey	herodes	okikak	ni	omomojti	iwa	nochteh	tlakah	jerusalén	noiwa	omomojtjkeh							
	rey	herodes	tlapololistle	nijkuak	ijkón	okikakke	iwan	nochtin	tlakaj	jerusalén	nojiki	okinmakak	tlapololistle						
huan	quema	tanahuatiquetl	herodes	quicajqui	nopa	momajmati	huan	nojquiya momajmatijque	nochi	masehualme	altepetl	jerusalén							
ihcuac	in	ueyitiquiahqui	herodes	oquicac	non	tlailiuis	oyolpahsoliu	iuán	nochi	jerusalén									
hua	cuõc	oquecac	nõnca	inu	rey	herodes	omotequepacho	lalebes	hua	noche	giente	de	jerusalén	nuyejquomotequepachojque					
huan	quema	tanahuatiquetl	herodes	quicajqui	nopa	tlacame quijtohuayayaj	momajmati	huan	nojquiya momajmatijque	nochi	tlacame	altepetl	jerusalén						
ijcuac	on	rey	herodes	ijcon	ocac	sanoyej	onomojtij	niman	nochi	tlacatl	jerusalén	no	onomojtij						
quiman	quimatic	herodes	hué	in	quijtohuayayaj	huéyoyenten	yajmo	huil	mosehuiaya	yihual	huan	noje	niman	ca	yehuanten pa	jerusalén			
ihcuac	huéyixtoc	herodes	oquimat	inon	yehuatl	omomouti	ihuan	nochin	tlacamen	jerusalén									
huan	quema	tanahuatiquetl	herodes	quicajqui	nopa	quijtohuayayaj	momajmati	huan	nojquiya momajmatijque	nochi	masehualme	altepetl	jerusalén						
a	lohe	la	o	herode	alii	apoapo	ae	la	kona	oili	oia	a	me	ko	jerusalem	a	pau		
lè	we	ewod	pran	nouvèl	la	sa	te	boulvèse	tèt	li	sa	te	boulvèse	tout	moun	laviil	jerizalèm	yo	tou
when	king	herod	heard	this	he	was	troubled	and	everyone	in	jerusalem	was	troubled	with	him				
huan	cuac	rey	herodes	quicayic	nijfn	pehuacmoyolcuejmoloa	huan	hasta	nochi	in	tagayot	jerusalén	no	moyolcuejmolojque					

# 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

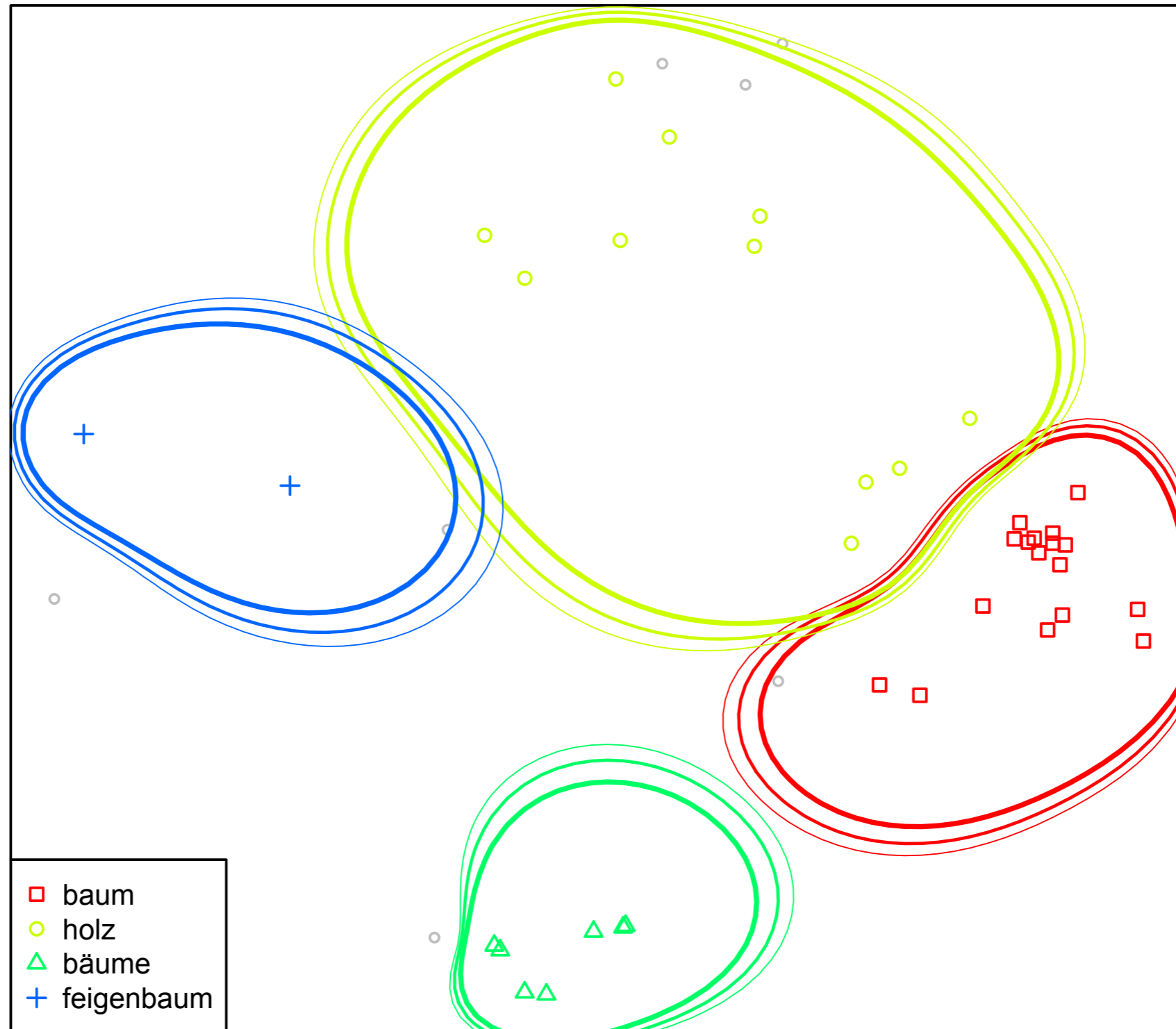


	<u>tree</u>	<u>wood (stuff)</u>	<u>firewood</u>	<u>small forest</u>	<u>large forest</u>
German	<i>Baum</i>	<i>Holz</i>		<i>Wald</i>	
Danish	<i>træ</i>			<i>skov</i>	
French	<u><i>arbre</i></u>	<u><i>bois</i></u>		<u><i>forêt</i></u>	
Spanish	<u><i>árbol</i></u>	<u><i>madera</i></u>	<u><i>leña</i></u>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

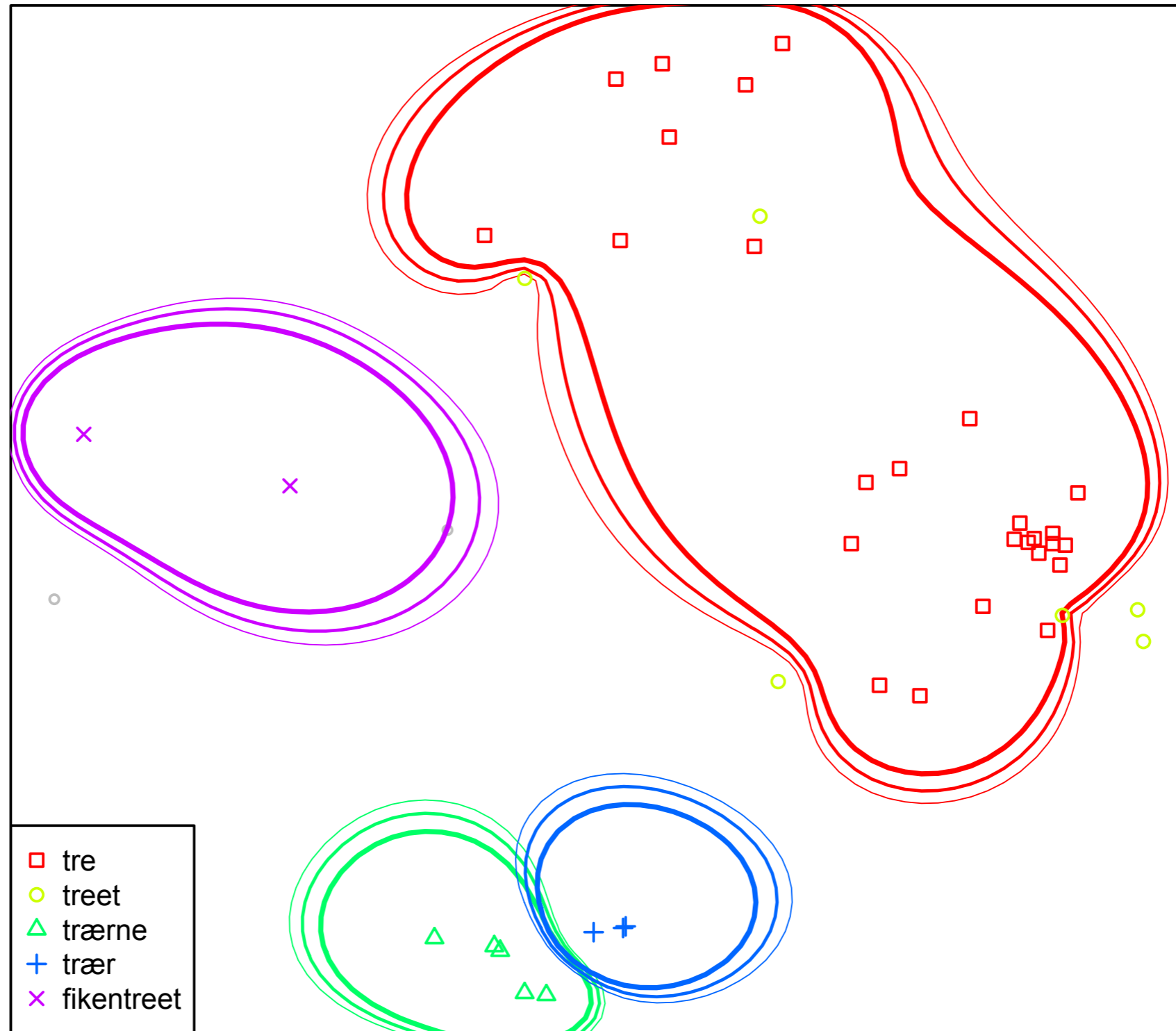
Louis Hjelmslev  
*Prolegomena to a Theory of Language* (1963)



# deu-x-bible-erben.txt

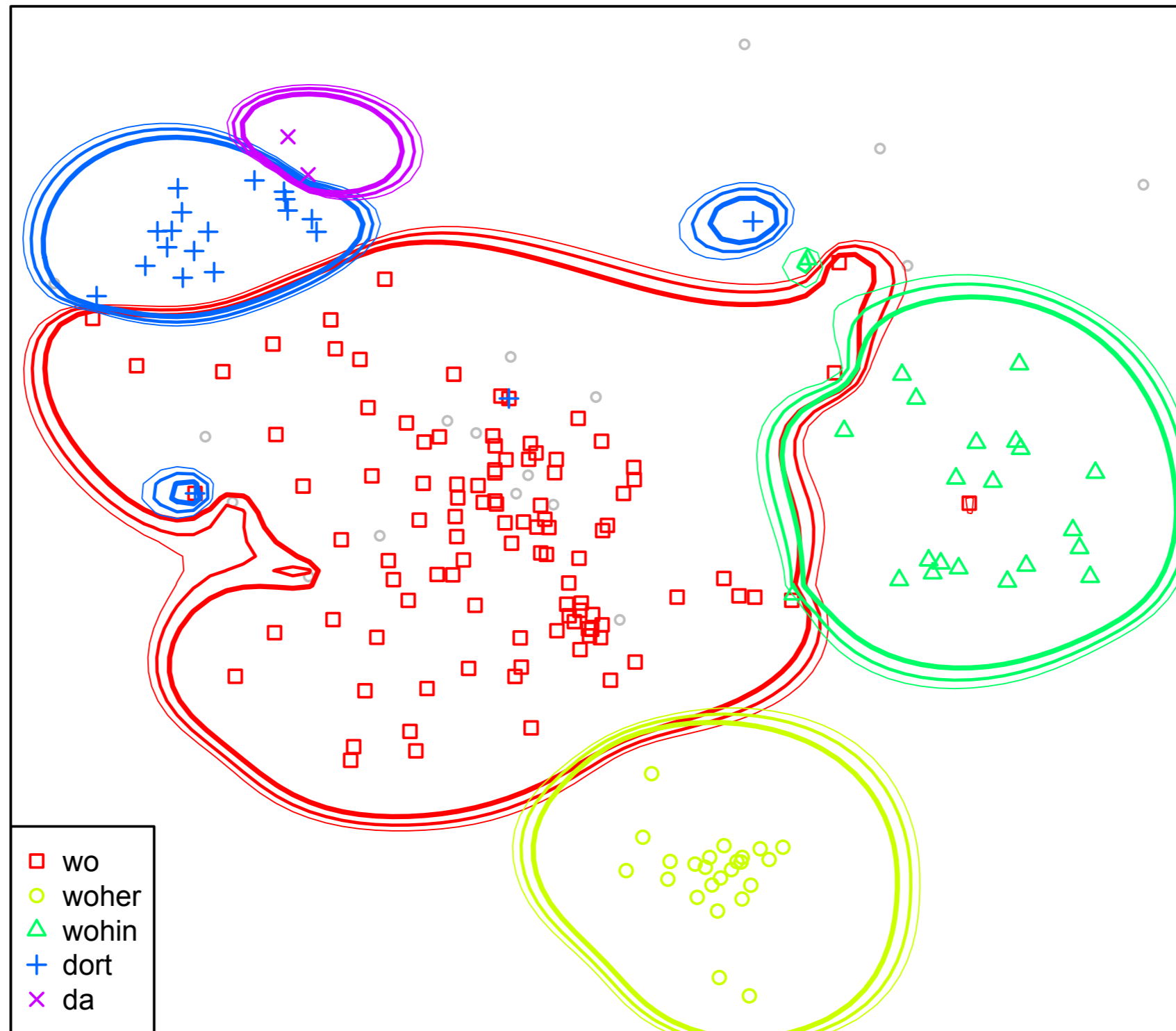


# nob-x-bible-2007.txt

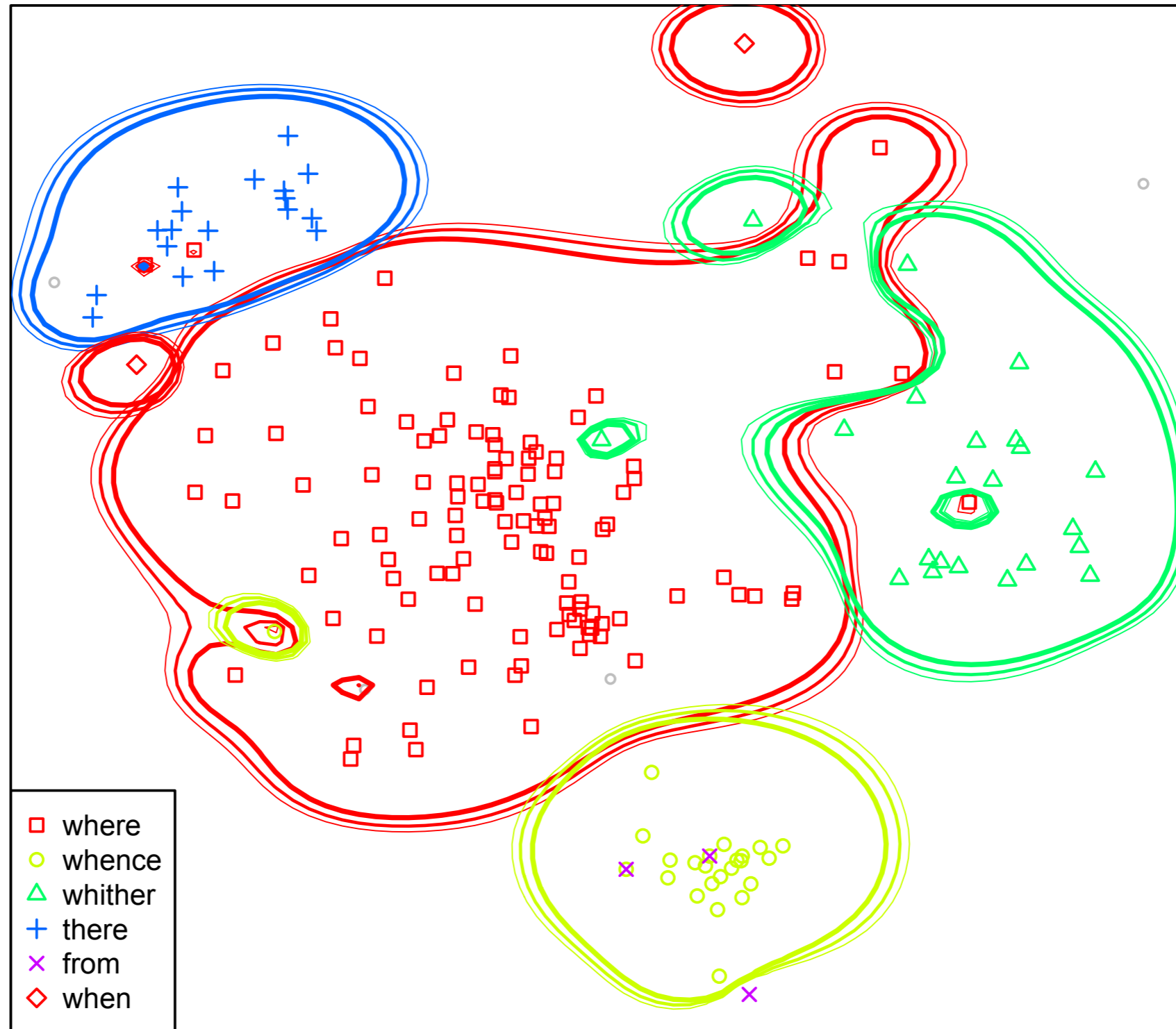


where

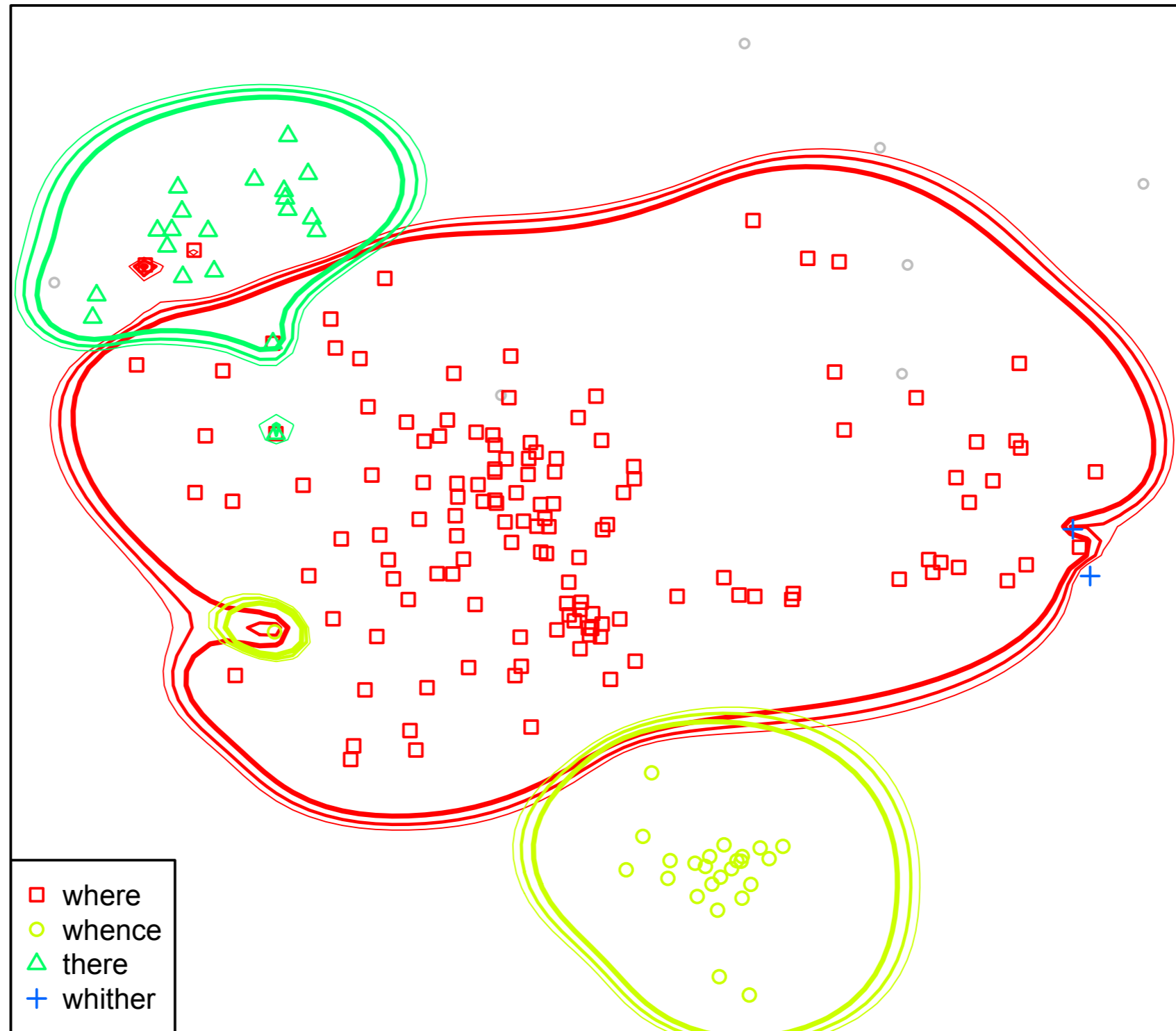
# deu-x-bible-pattloch.txt



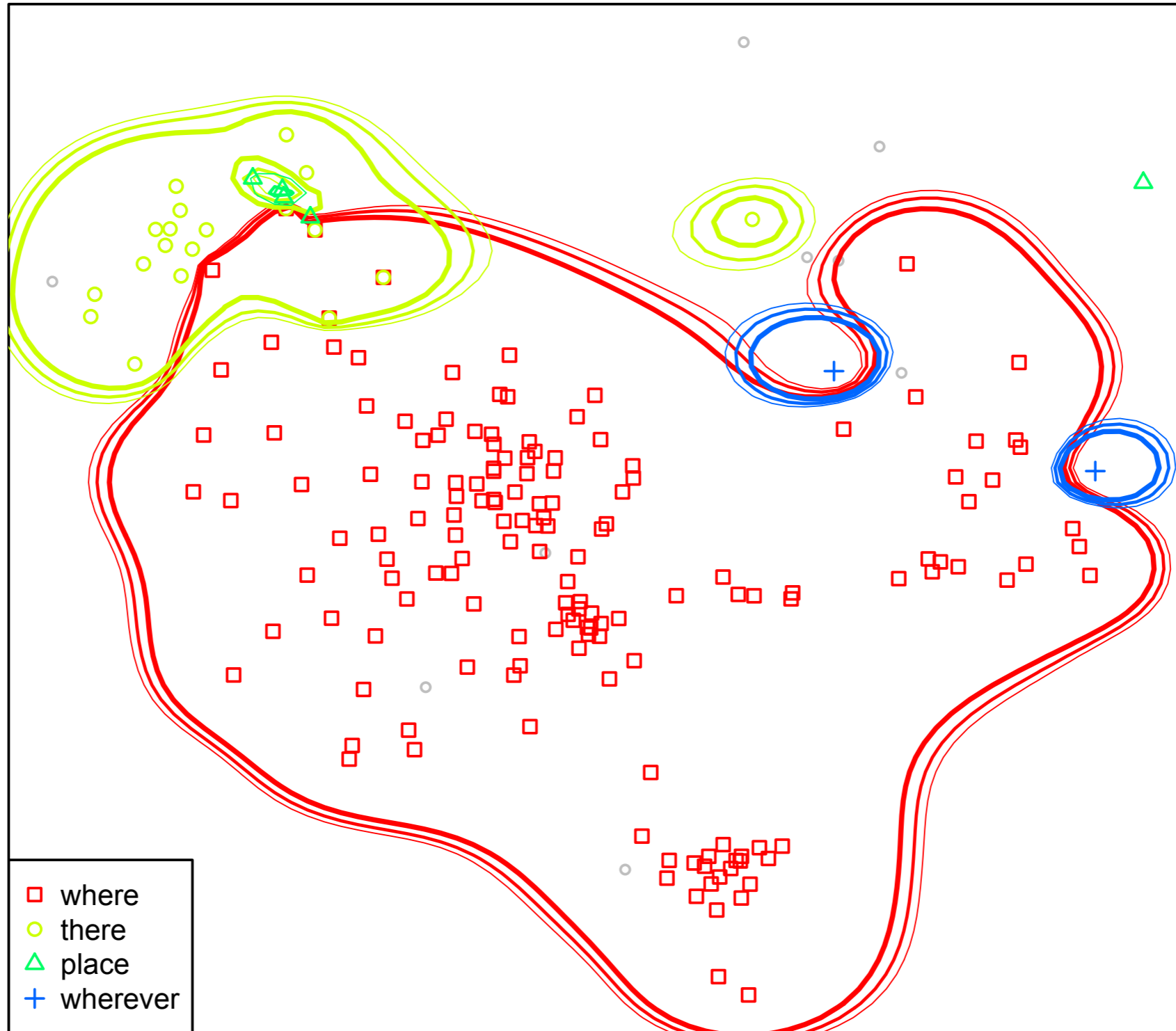
# eng-x-bible-kingjames.txt



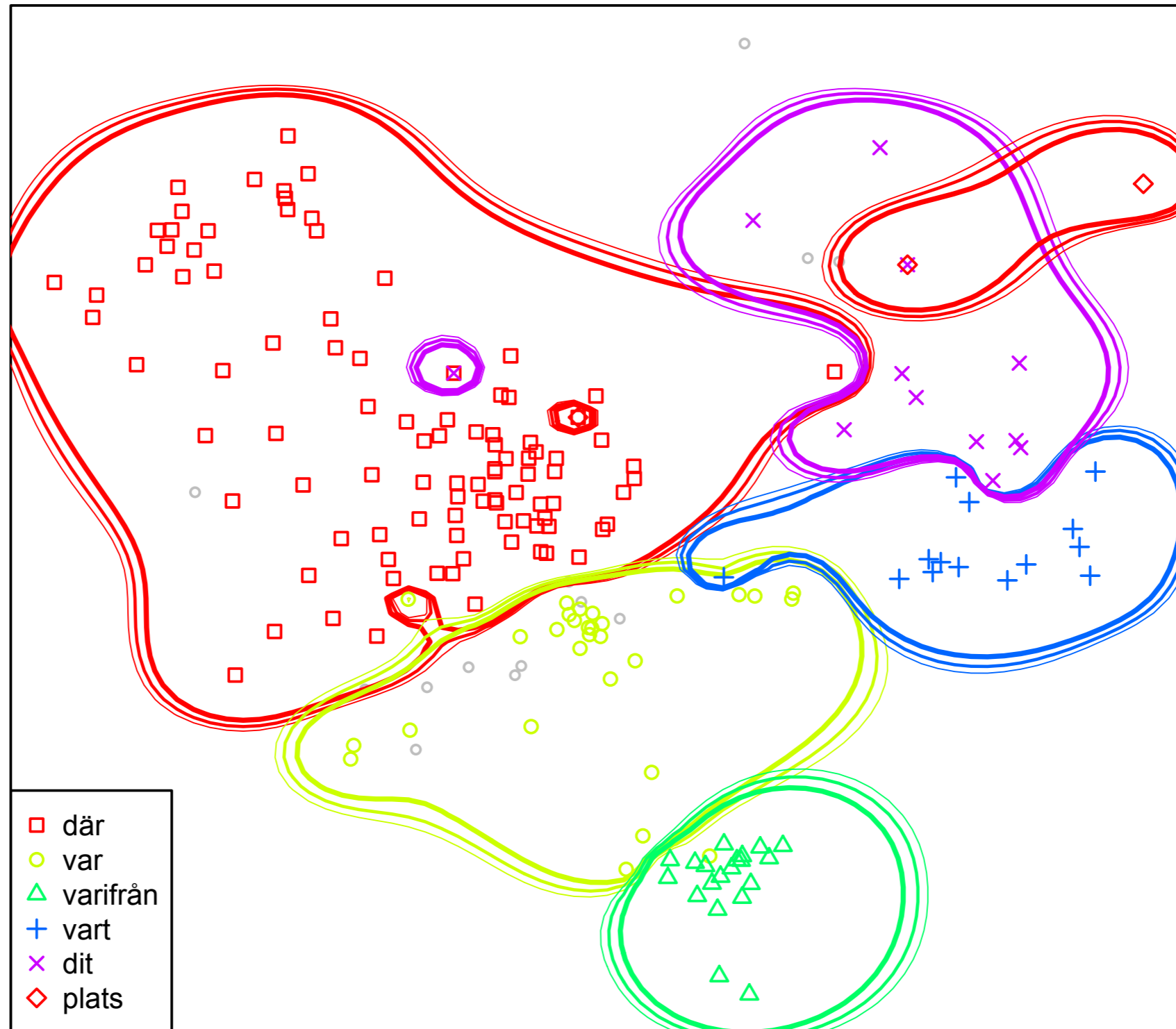
# eng-x-bible-darby.txt



# eng-x-bible-treeoflife.txt



# swe-x-bible-folk1998.txt

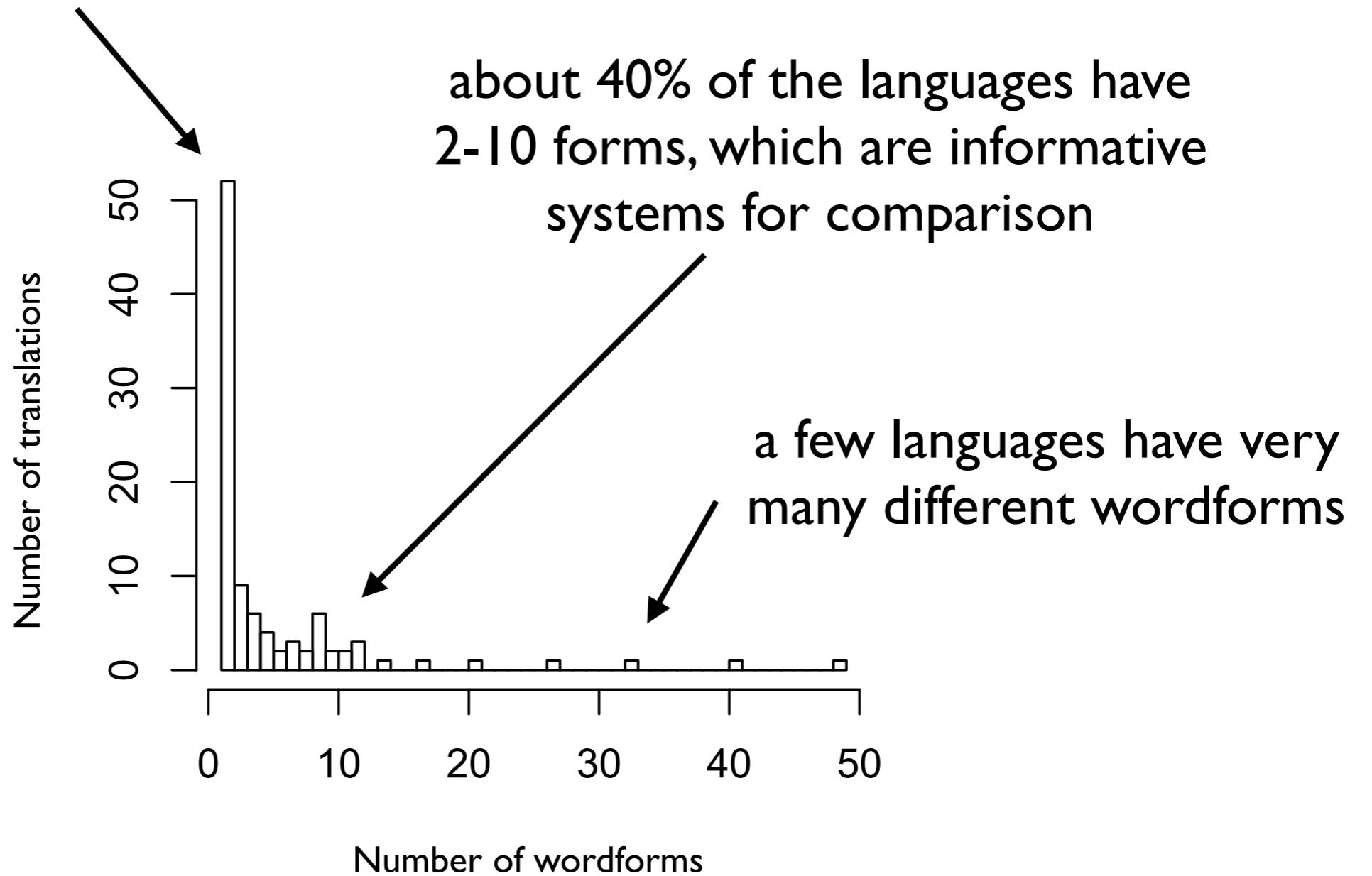




# quick test: translation of “*Jerusalem*”

- Identify wordforms that correspond to the English name *Jerusalem* in 100 translations
  - ▶ mostly automatic, minor manual work still needed
- Many languages will just have one wordform, but some will have more than one
- These different wordforms might give us information about local case functions

more than half of the languages  
have just one wordform



# Bamanankan

(“Bambara”, ISO 639-3 bam, spoken in Mali)

Word ID	Bamanankan	Best English back-translation	Best German back-translation	Relevance
1441	jerusalem	jerusalem	jerusalem	1
1443	jerusalemkaw	jerusalem	jerusalem	0,67
1442	jerusalemka	---	---	0,6

# Angaataha

(ISO 639-5 agm, spoken in Papua New Guinea)

- jerusaremthanda
- jerusaremthandaahapɿ
- jerusaremthandɿ
- jerusaremthandaahiyai
- jerusaremthandaahɿ
- jerusaremthandaahapɿhiyauntɿ
- jerusaremthandaahiyaisangi
- jerusaremthandaahapɿhiya
- jerusaremthandaahɿraapɿ
- jerusaremthandaahɿhɿ
- jerusaremthandaahɿhe
- jerusaremthandamɿ
- jerusaremɿmanda
- jerusaremthandapɿ
- jerusaremɿndɿ
- jerusaremthandaahapɿto
- jerusaremthandaahapaahɿhɿ
- jerusaremthandi
- jerusaremɿmandaahapɿ
- jerusaremthandaahuntɿ
- jerusaremthandaahapuntɿ
- jerusaremthandaahiya
- jerusaremthandamɿhɿntɿ
- jerusaremthandaahapɿhiyaatihɿ
- jerusaremthandaahapɿhiyaate
- jerusaremthandaahiyauntɿ
- jerusaremosthiyaate

# Amharic

(ISO 639-3 amh, spoken in Ethiopia)

- ኢየሩሳሌም
- በኢየሩሳሌም
- ከኢየሩሳሌም
- ኢየሩሳሌምም
- በኢየሩሳሌምም
- ኢየሩሳሌምን
- ከኢየሩሳሌምም
- የኢየሩሳሌም
- ለኢየሩሳሌም
- ለኢየሩሳሌምም
- የኢየሩሳሌምንም
- የኢየሩሳሌምምም

# Amharic

(ISO 639-3 amh, spoken in Ethiopia)

- ኢየሩሳሌም
- በኢየሩሳሌም
- ከኢየሩሳሌም
- ኢየሩሳሌምም
- በኢየሩሳሌምም
- ኢየሩሳሌምን
- ከኢየሩሳሌምም
- የኢየሩሳሌም
- ለኢየሩሳሌም
- ለኢየሩሳሌምም
- የኢየሩሳሌምንም
- የኢየሩሳሌምም

# Highland Puebla Náhuatl

(ISO 639-3 azz)

- Jerusalén
- Jerusaléncopaca

# Huasteca Náhuatl

(ISO 639-3 nch)

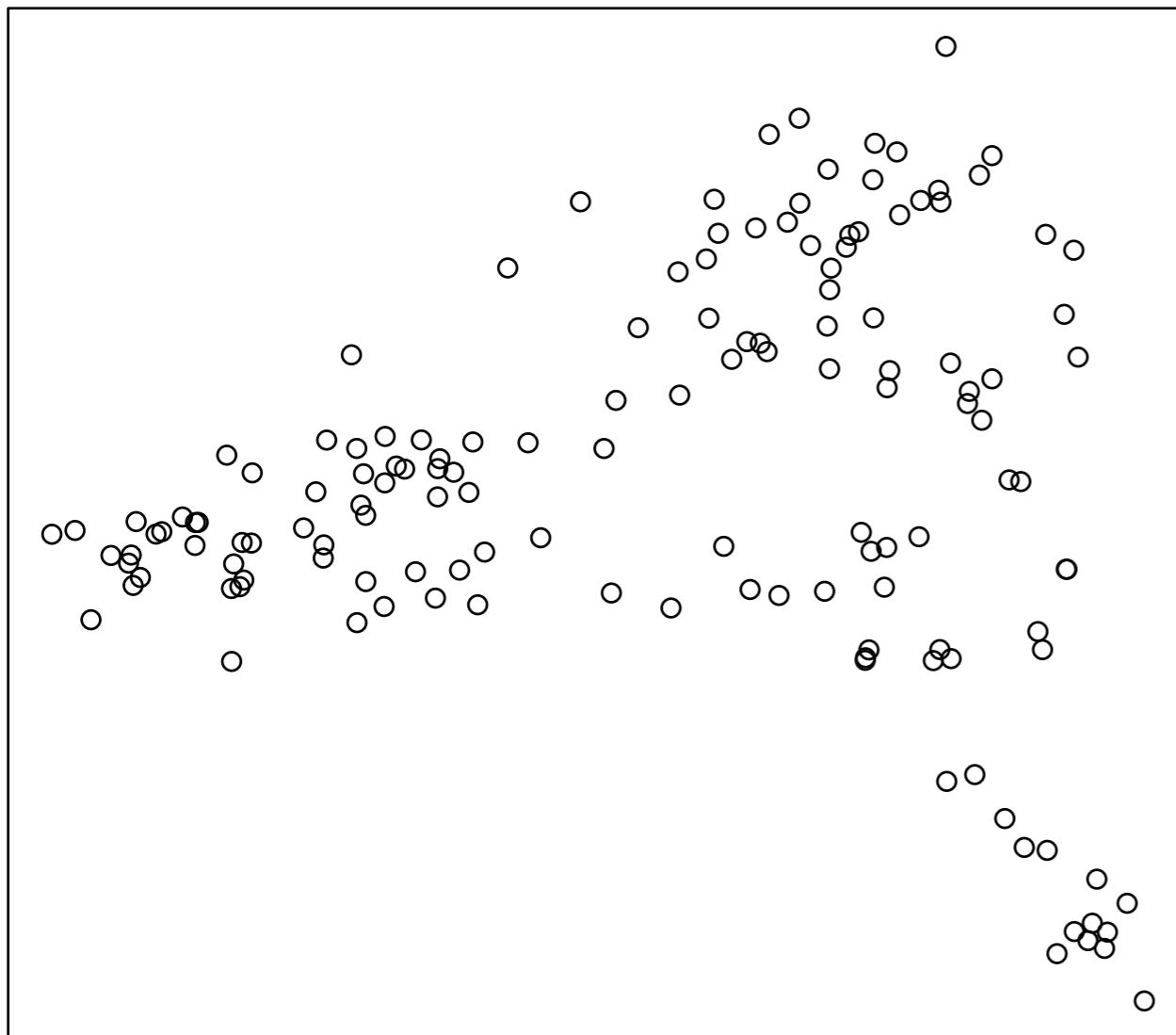
- Jerusalén
- anJerusalén



98 languages, in total 520 different wordforms

I selected 167 verses including *Jerusalem* only once in more than 40 languages

Matrix of size 520 x 167 coding the distribution of wordforms over verses



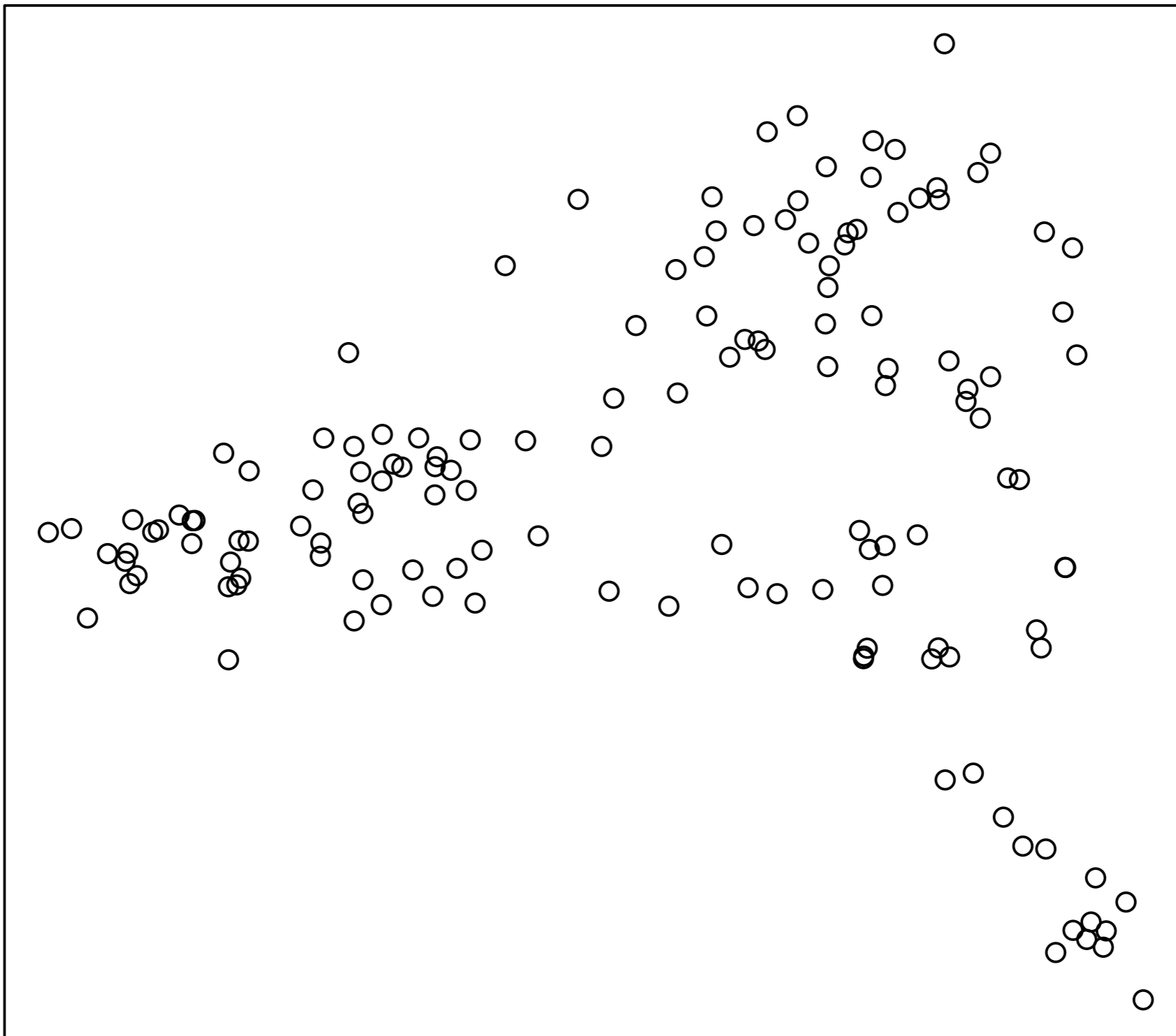
Two contexts of *Jerusalem* are similar when they often share the same wordform in language after language

here showing two main dimensions of variation of 167 contexts

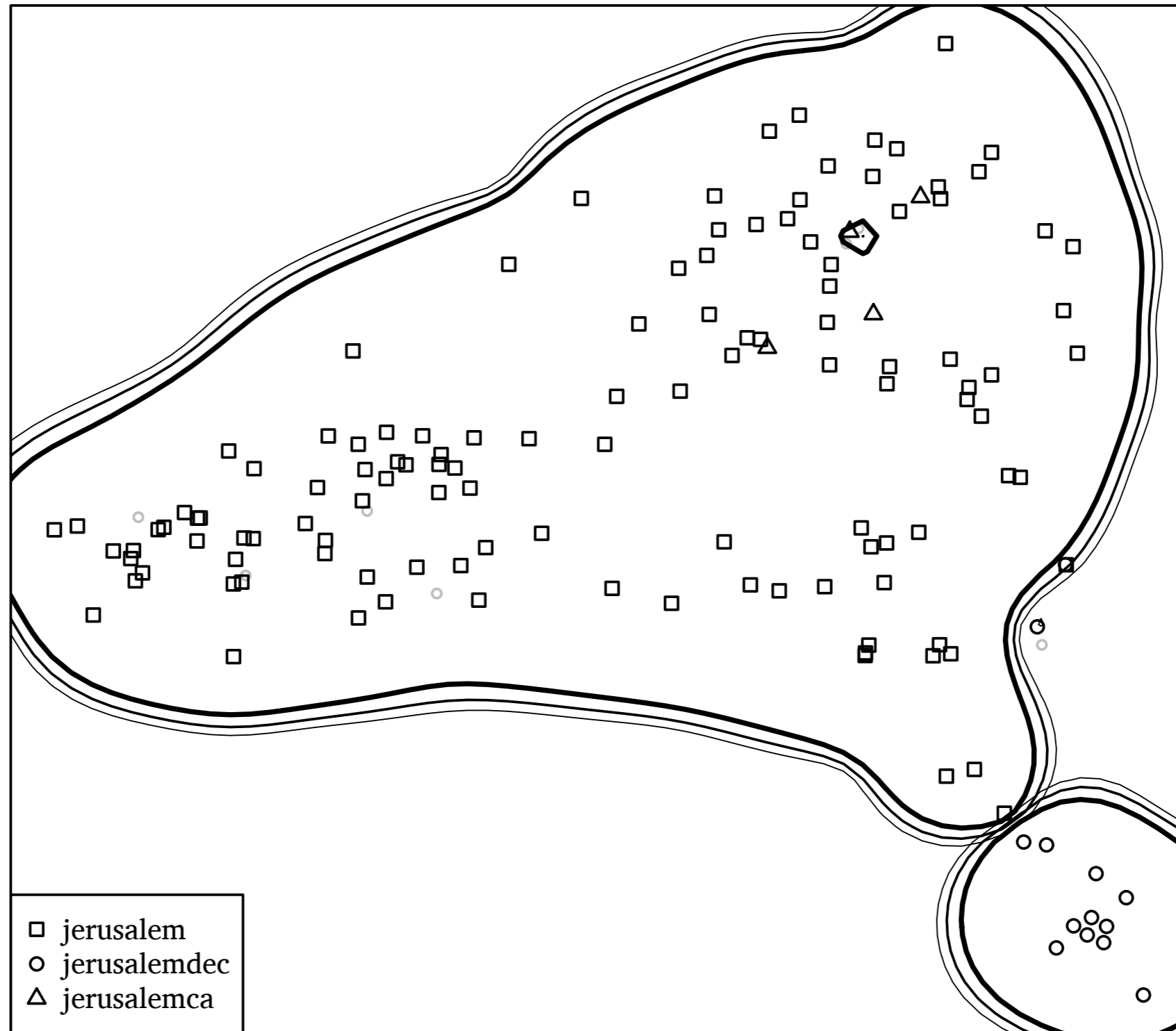
- the importance of dimensions depend strongly on the content of the corpus, which we cannot control
- only the first two dimensions are discussed here because of easy visualisation

# System Typology

- Taking the different wordforms in a language as a system of marking
- How similar are the systems functionally?
- We can compare functions through distribution in the text

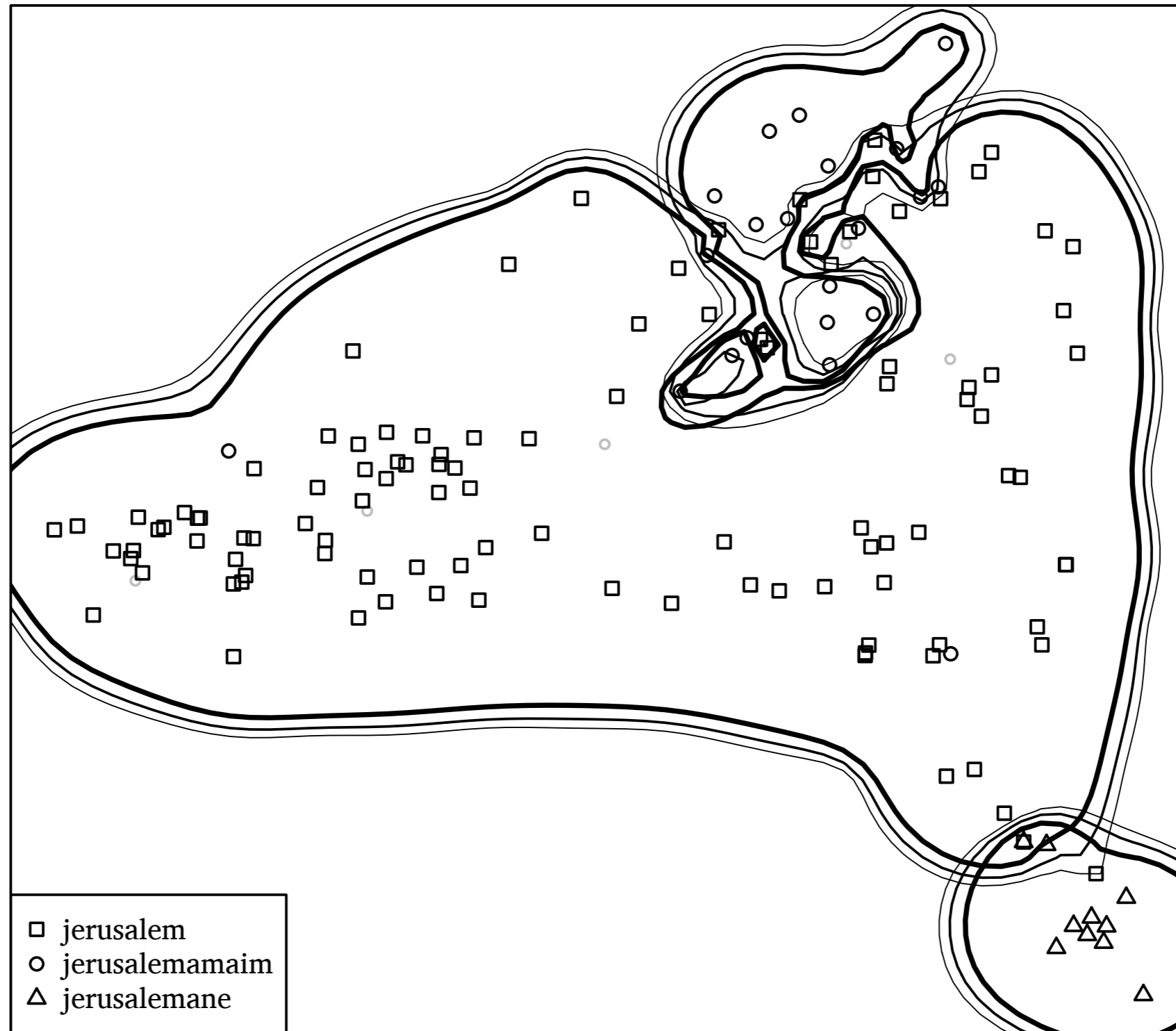


aey



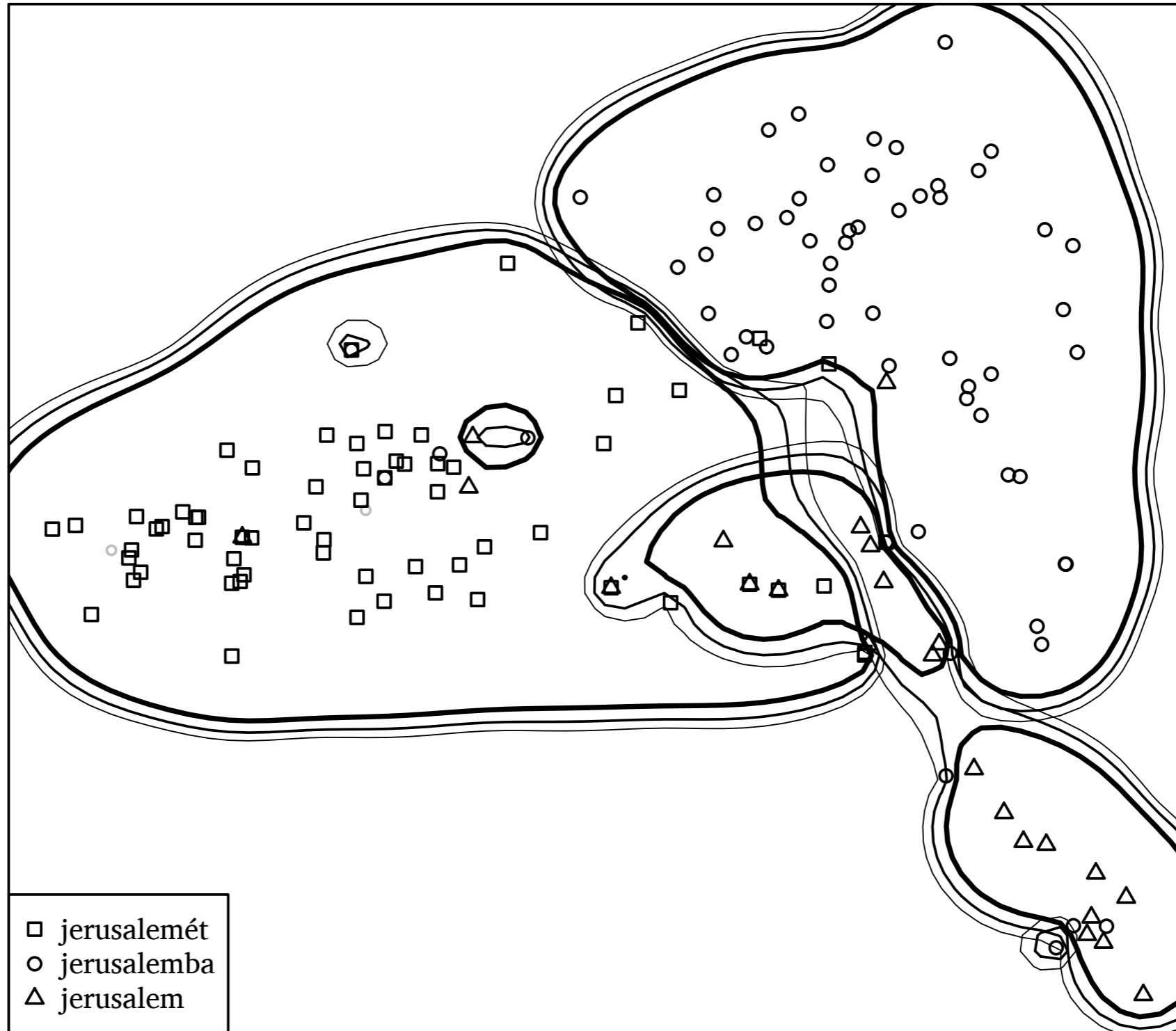
Amele (A language of Papua New Guinea)

# aai



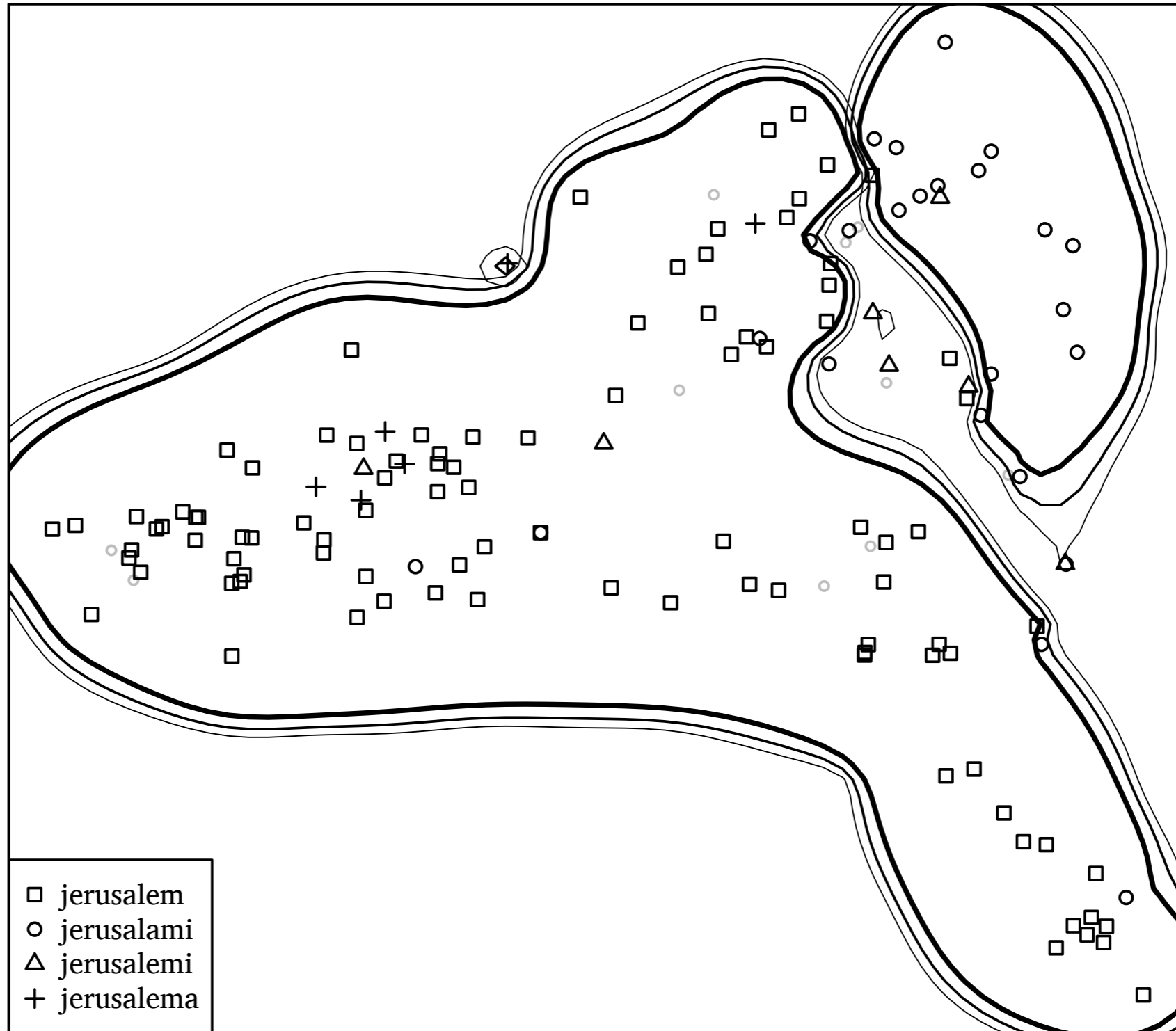
Arifama-Miniafia (a language of Papua New Guinea)

abt



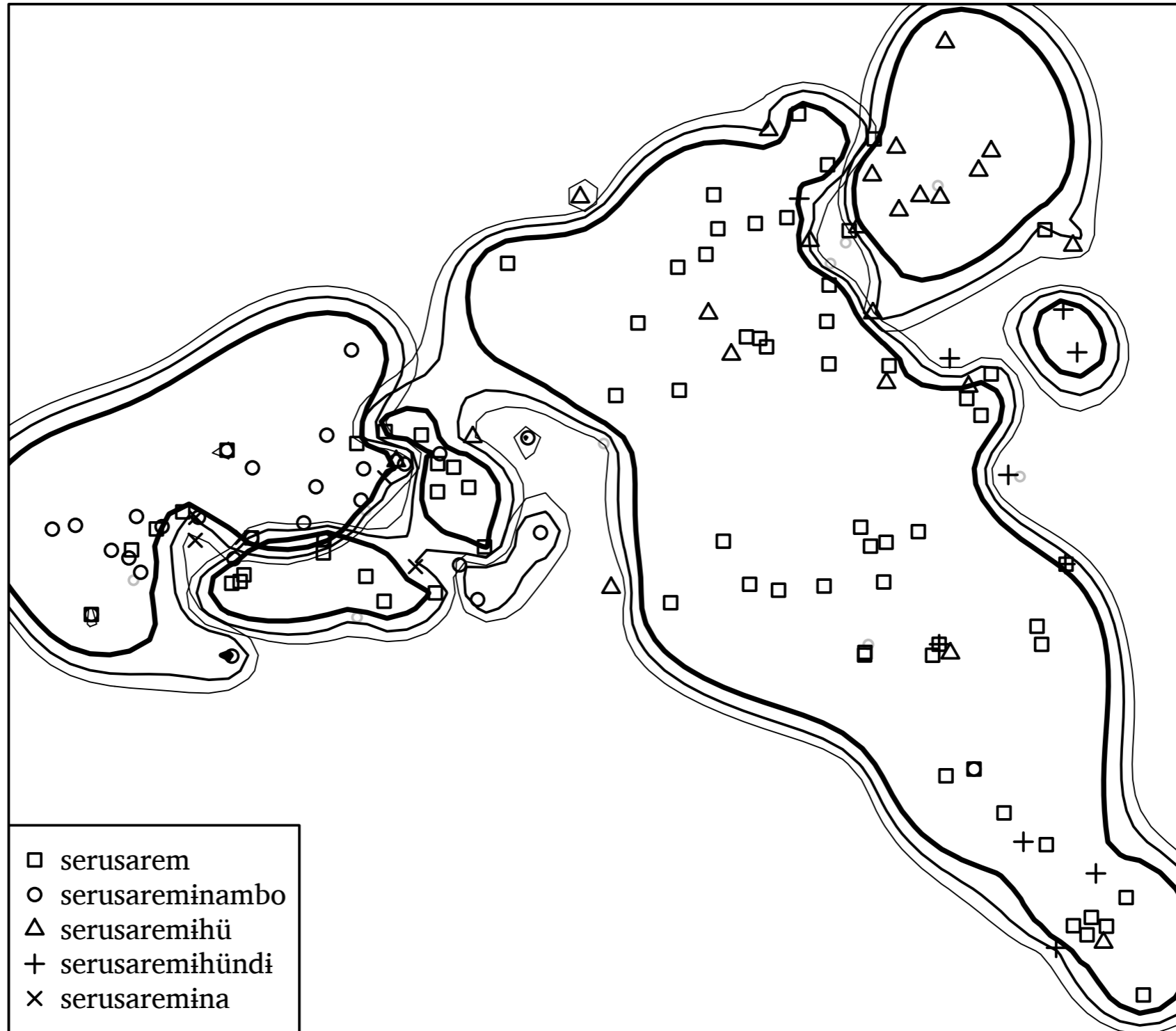
Ambulas (a language of Papua New Guinea)

aoj



Muffian (a language of Papua New Guinea)

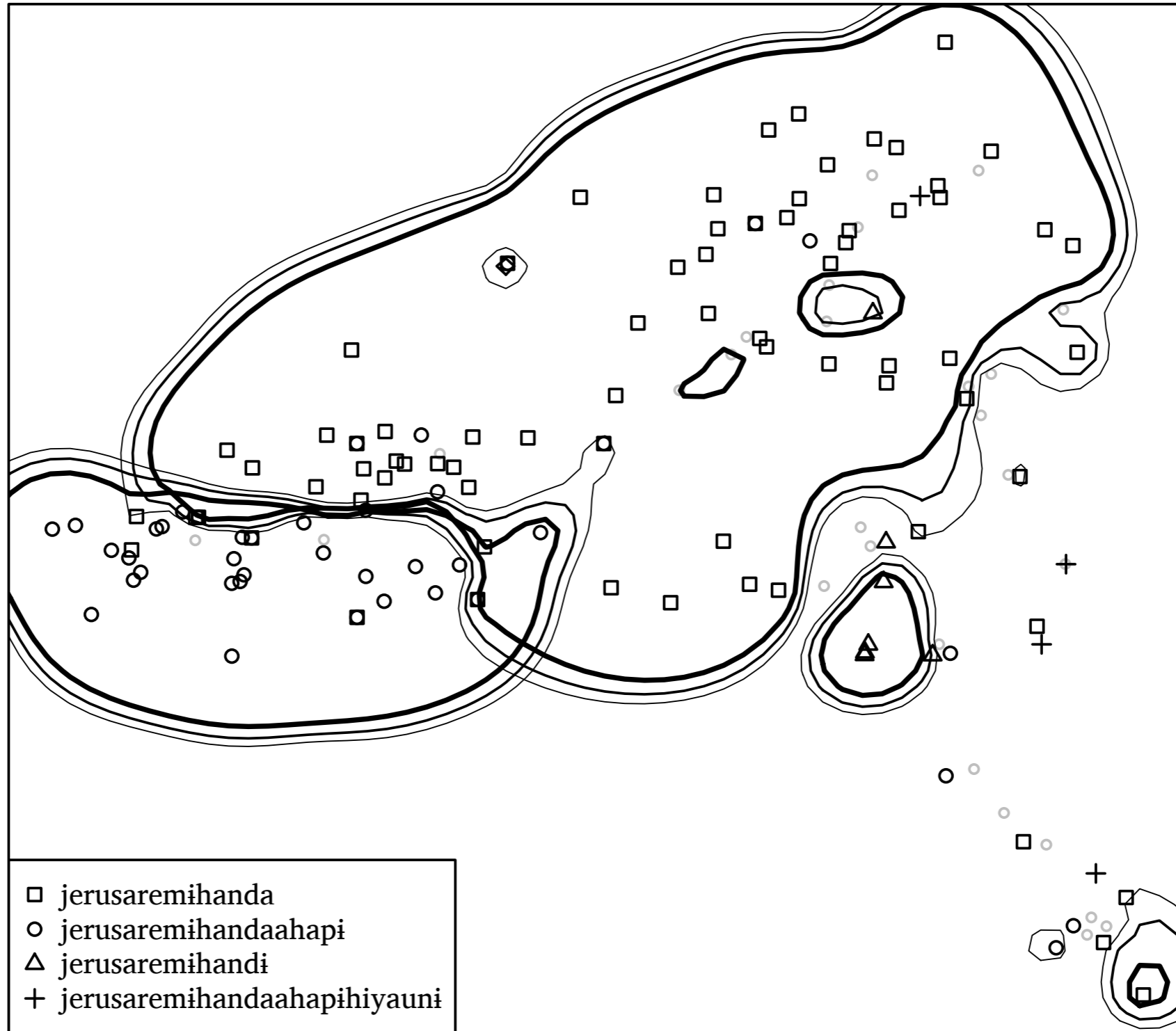
# agg



Angor (a language of Papua New Guinea)

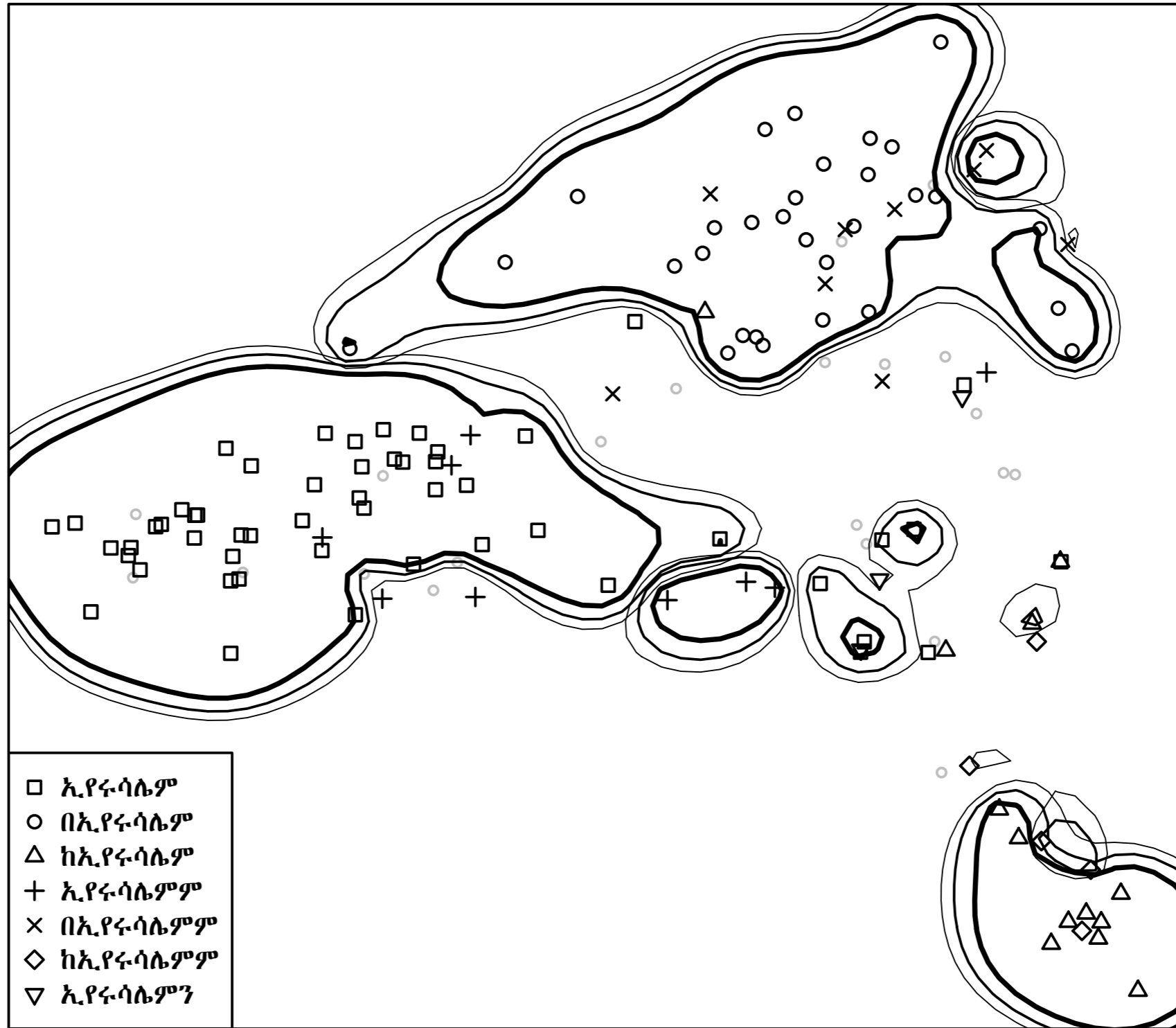


agm



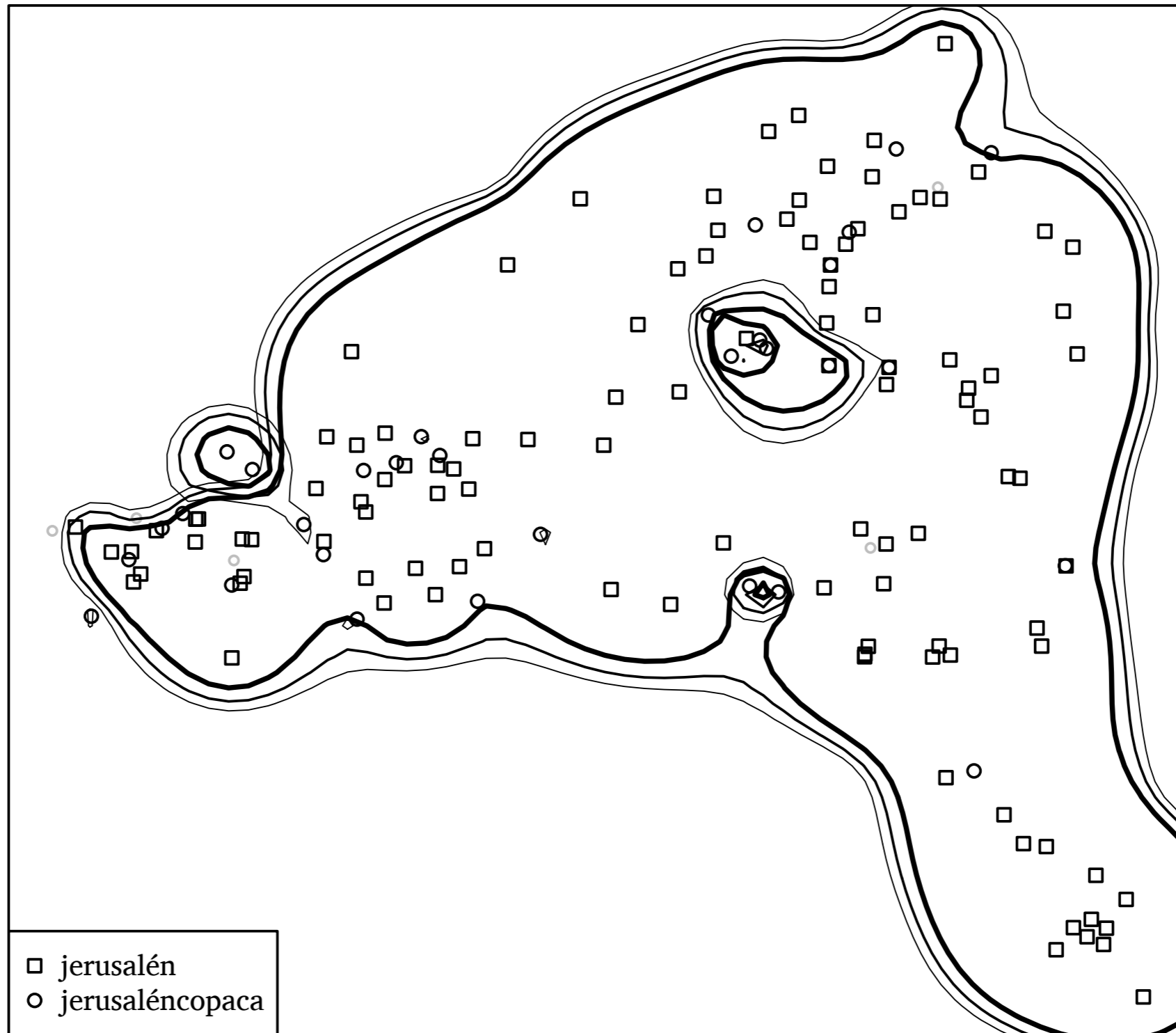
Angaatiha (a language of Papua New Guinea)

amh



Amharic (a language of Ethiopia)

azz



Highland Puebla Náhuatl

42002025: "... there was a man **in Jerusalem** whose name was Simeon ..."

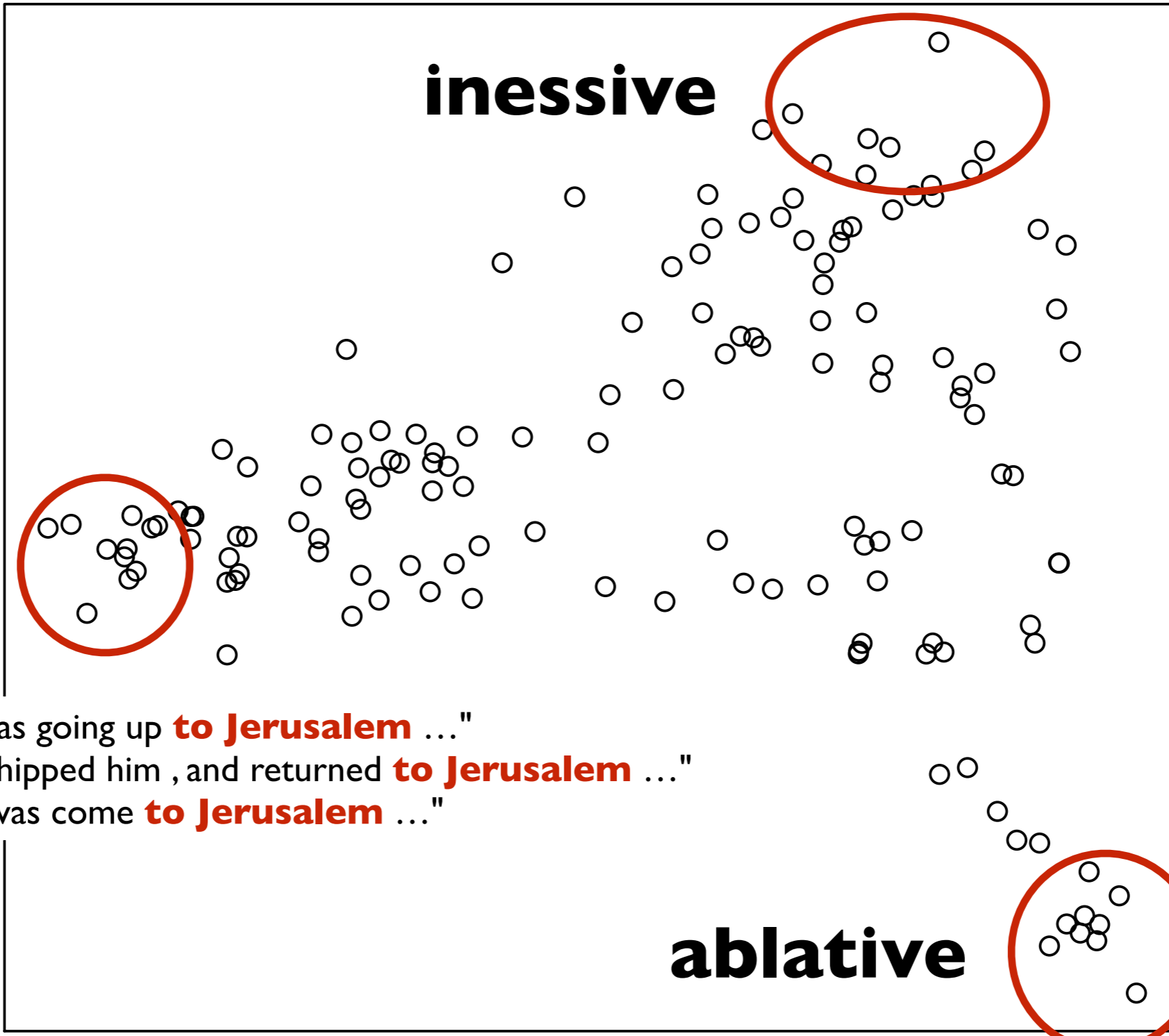
44004005: "... their rulers and elders and scribes were gathered together **in Jerusalem** ..."

44021011: "... So shall the Jews **at Jerusalem** bind the man that owneth this girdle ..."

**allative**

**inessive**

**ablative**



40020017: "... as Jesus was going up **to Jerusalem** ..."

42024052: "... they worshipped him , and returned **to Jerusalem** ..."

44009026: "... when he was come **to Jerusalem** ..."

41003022: "And the scribes that came down **from Jerusalem** said ..."

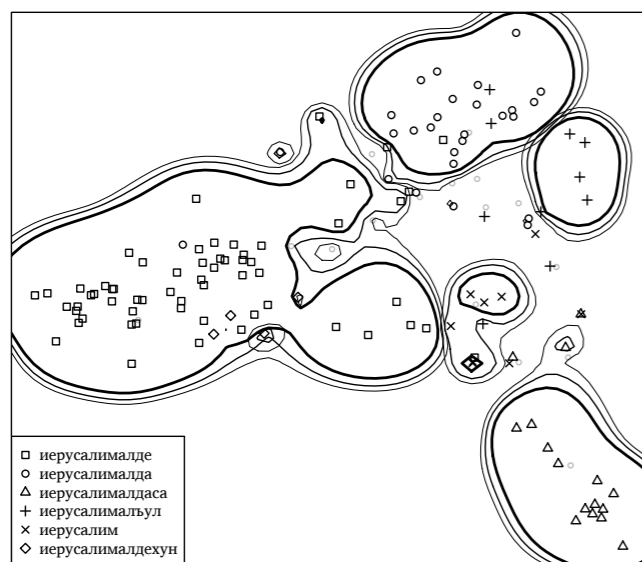
42005017: "... there were Pharisees and doctors of the law sitting by , who were come **out of** every village of Galilee and Judaea and **Jerusalem** ..."

42010030: "... A certain man was going down **from Jerusalem** to Jericho

...  
"

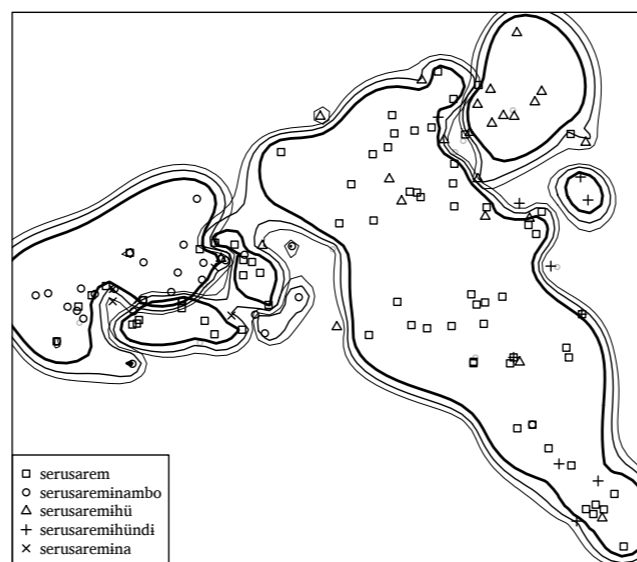
# allative dominated

ava



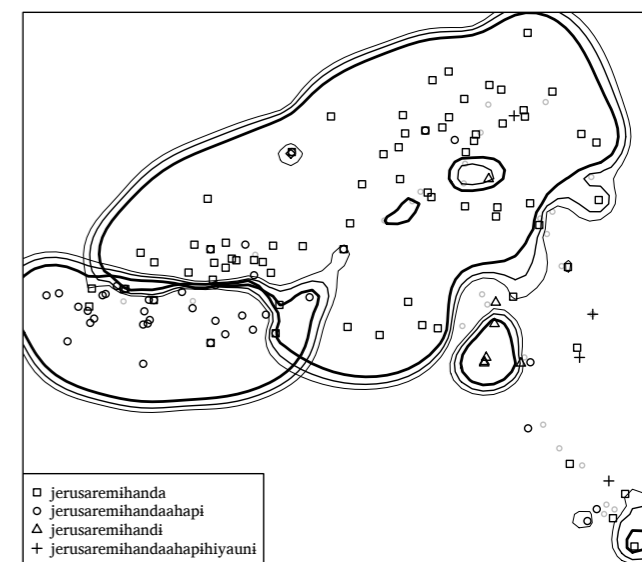
# ablative dominated

agg



# inessive dominated

agm



# Prospects

- Parallel texts offer the possibility for detailed functional comparison across languages
- The comparison is based on actual examples, so each typological generalisation can be scrutinised by specialists
- Algorithmic assistance is possible, so manual decisions