

A world map with a light blue background and green landmasses. Several colored dots are scattered across the map: two yellow dots in the Pacific region, one yellow dot in the Indian Ocean, one pink dot in Southeast Asia, one pink dot in East Africa, one purple dot in Australia, and one yellow dot in the Pacific near the Philippines. A large white rectangular box is centered over the map, containing the title text.

# Language comparison through massively parallel texts

A world map with a light blue background and green landmasses. Several colored dots are scattered across the map: two yellow dots in the Pacific region, one yellow dot in the Indian Ocean, one pink dot in Southeast Asia, one pink dot in East Africa, one purple dot in Australia, and one yellow dot in the Pacific near the Philippines. A large white rectangular box is centered over the map, containing the title text.

Michael Cysouw  
Philipps-Universität Marburg

Philipps



Universität  
Marburg

# Lessons from worldwide language diversity

- Possible vs. impossible languages
- Universal categories
- The problem of comparing languages

# Lessons from worldwide language diversity

- **Possible vs. impossible languages**
- Universal categories
- The problem of comparing languages

# Possible vs. Impossible

- Traditional hope: some structures are possible, other are impossible in human language
- This idea cannot be maintained
- Assumed impossible structures always turn up after the investigation of more languages



Home

Features

Chapters

Languages

References

Authors

## Feature 81A: Order of Subject, Object and Verb



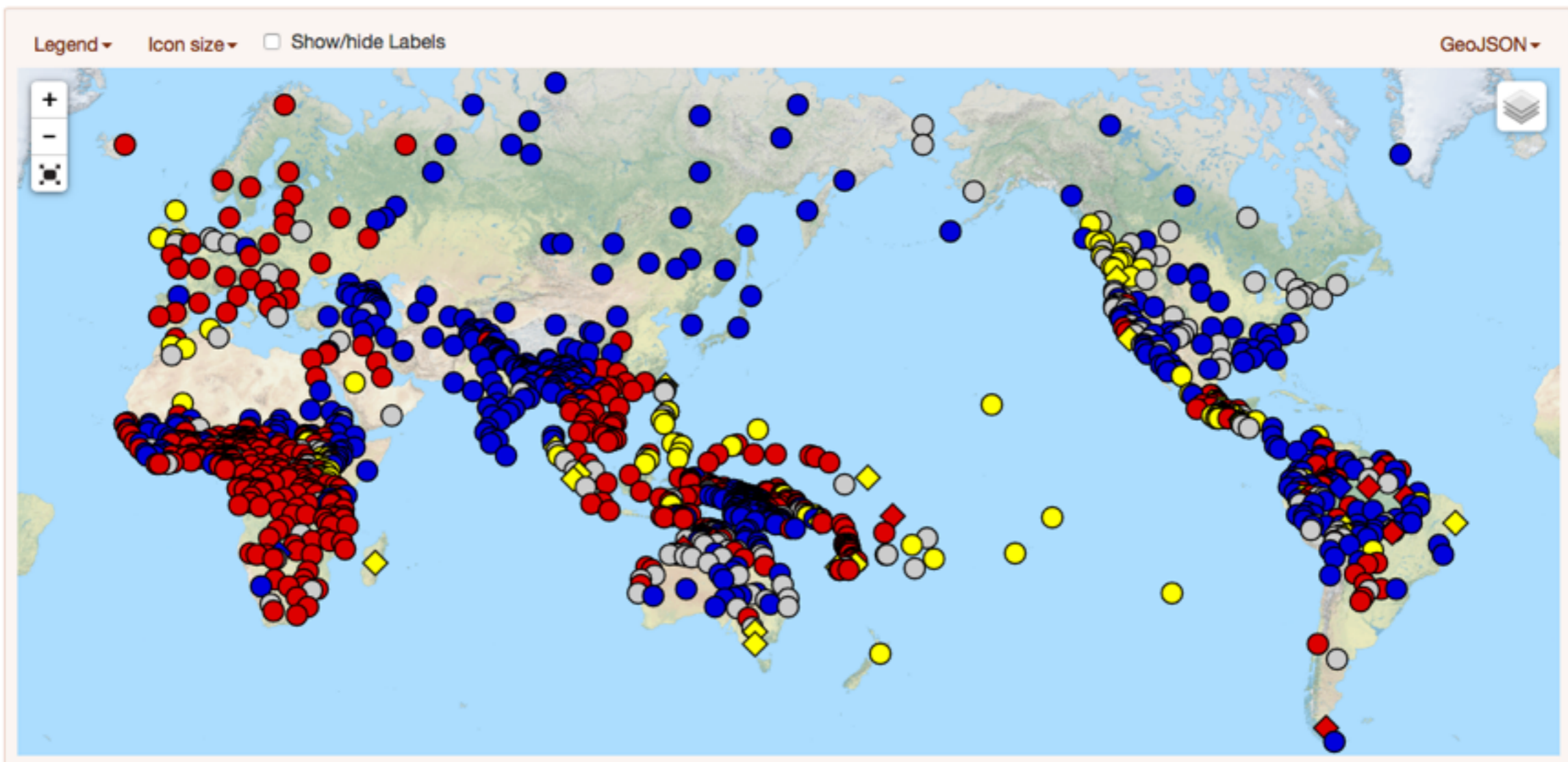
This feature is described in the text of chapter 81 [Order of Subject, Object and Verb](#) by [Matthew S. Dryer](#) [cite](#)

You may combine this feature with another one. Start typing the feature name or number in the field below.

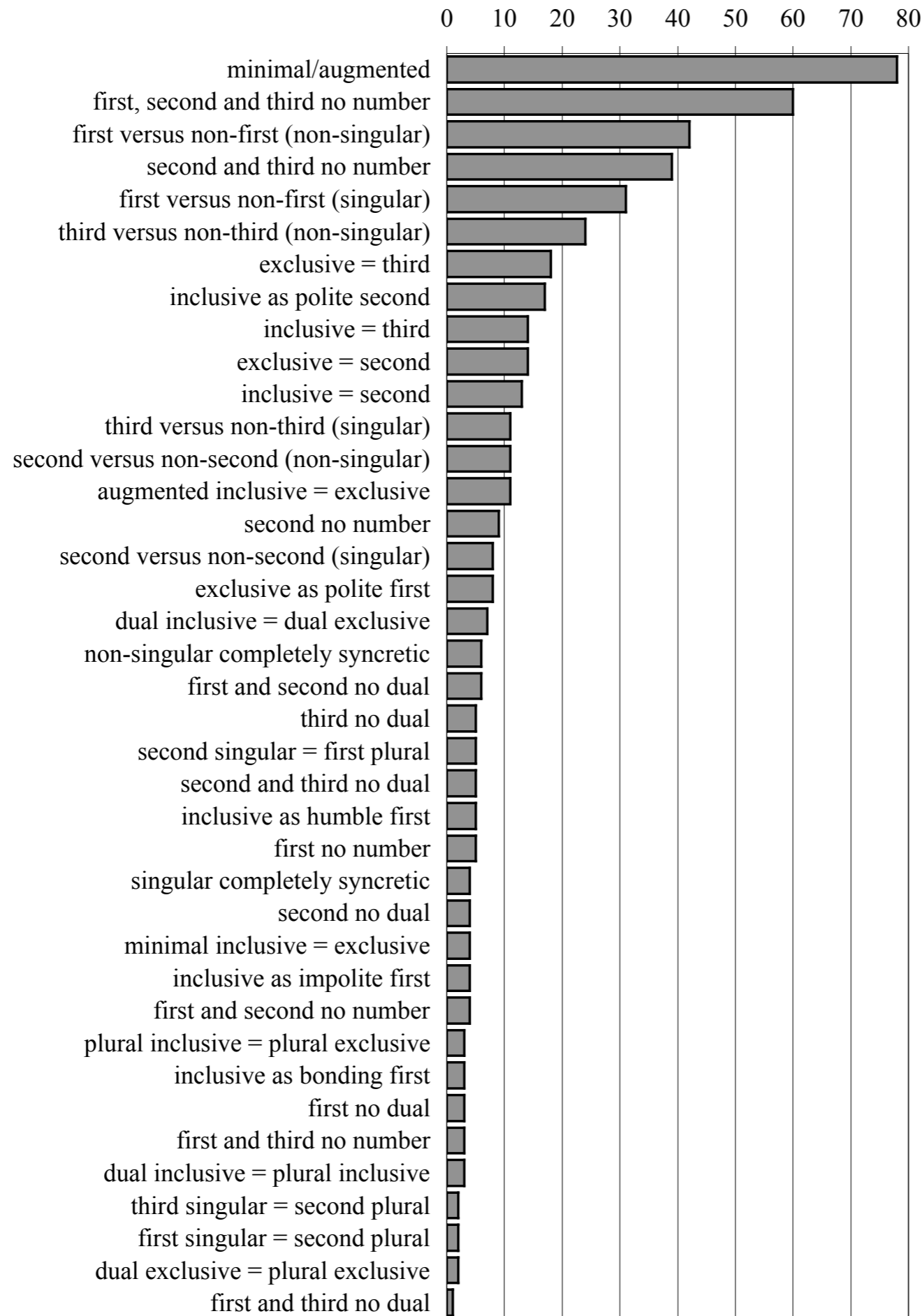
× 81A: Order of Subject, Object and Verb

### Values

<span style="color: blue;">●</span>	SOV	565
<span style="color: red;">●</span>	SVO	488
<span style="color: yellow;">●</span>	VSO	95
<span style="color: yellow;">◆</span>	VOS	25
<span style="color: red;">◆</span>	OVS	11
<span style="color: blue;">◆</span>	OSV	4
<span style="color: gray;">○</span>	No dominant order	189

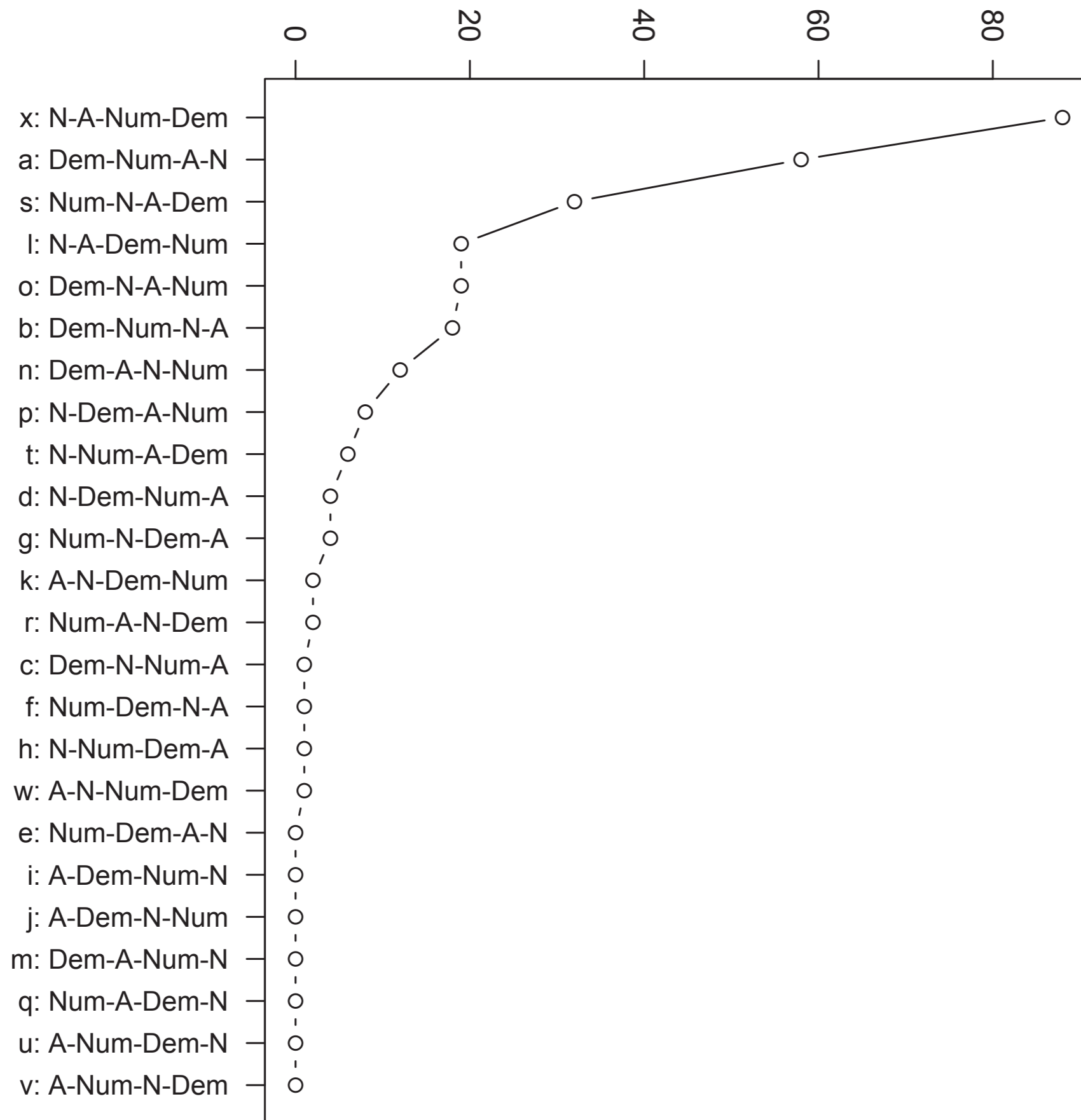


# Person marking Paradigms



Cysouw, Michael. 2005. What it means to be rare: the variability of person marking. In Zygmunt Frajzyngier, Adam Hodges & David S Rood (eds.), *Linguistic Diversity and Language Theories*, 235-258. (Studies in Language Companion Serie). Amsterdam: Benjamins.

# Frequency



# Word Order within NP

Cysouw, Michael. 2010. Dealing with diversity: towards an explanation of NP word order frequencies. *Linguistic Typology* 14(2). 253-287.



## Feature 96A: Relationship between the Order of Object and Verb and the Order of Relative Clause and Noun



This feature is described in the text of chapter 96

Relationship between the Order of Object and Verb and the Order of Relative Clause and Noun by Matthew S. Dryer

cite

You may combine this feature with another one. Start typing the feature name or number in the field below.

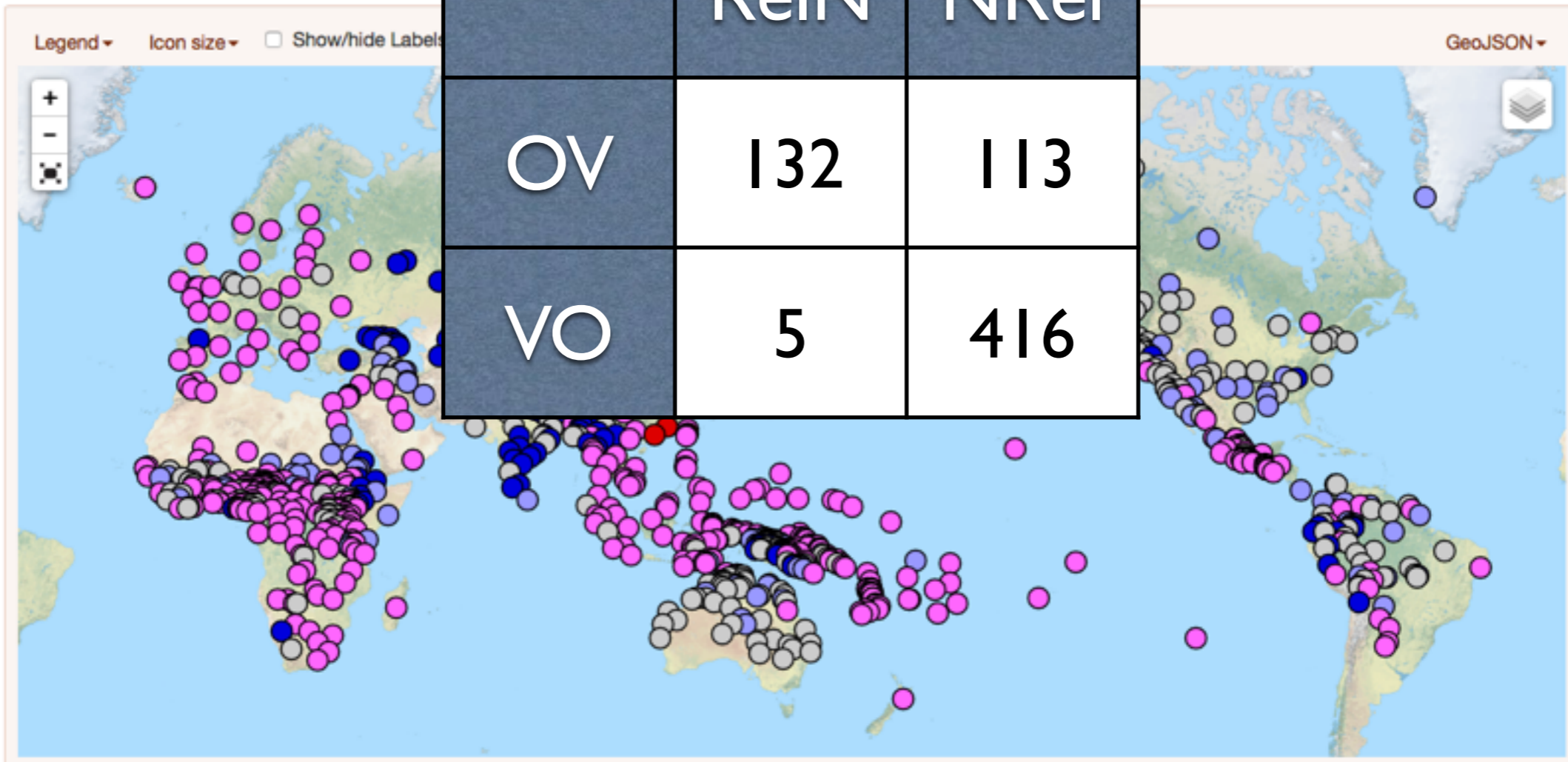
x 96A: Relationship between the Order of Object and Verb and the Order of Relative Clause and Noun

Submit

### Values

<span style="color: blue;">●</span>	OV and RelN	132
<span style="color: lightblue;">●</span>	OV and NRel	113
<span style="color: red;">●</span>	VO and RelN	5
<span style="color: magenta;">●</span>	VO and NRel	416
<span style="color: gray;">●</span>	Other	213

	ReIN	NRel
OV	132	113
VO	5	416







## Feature 96A: Relationship between the Order of Object and Verb and the Order of Relative Clause and Noun



This feature is described in the text of chapter 96

Relationship between the Order of Object and Verb and the Order of Relative Clause and Noun by Matthew S. Dryer

cite

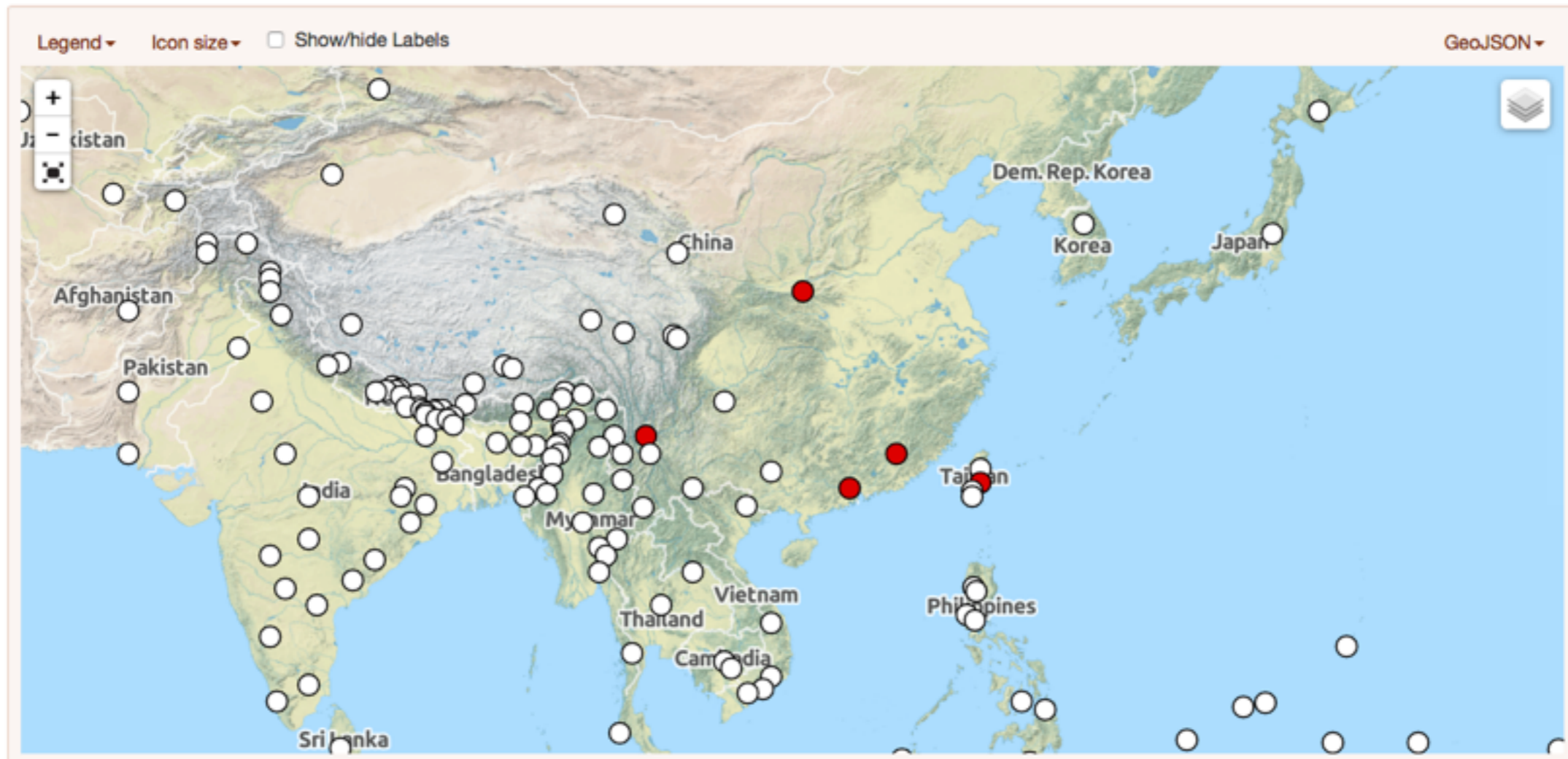
You may combine this feature with another one. Start typing the feature name or number in the field below.

x 96A: Relationship between the Order of Object and Verb and the Order of Relative Clause and Noun

Submit

### Values

<input type="radio"/>	OV and RelN	132
<input type="radio"/>	OV and NRel	113
<input checked="" type="radio"/>	VO and RelN	5
<input type="radio"/>	VO and NRel	416
<input type="radio"/>	Other	213



# Possible vs. Impossible

- The difference between attested and unattested is not a very robust observation
- Different samples will lead to different boundaries between ‘possible’ and ‘impossible’
- It is better to focus on the frequent phenomena: whether something is frequent or not is a much more robust observation

# Lessons from worldwide language diversity

- Possible vs. impossible languages
- **Universal categories**
- The problem of comparing languages

# Franz Boas

The necessary categories  
to describe a language  
“depend entirely on the  
inner form of each language”  
(*Handbook of American Indian  
languages*, 1911: 81)

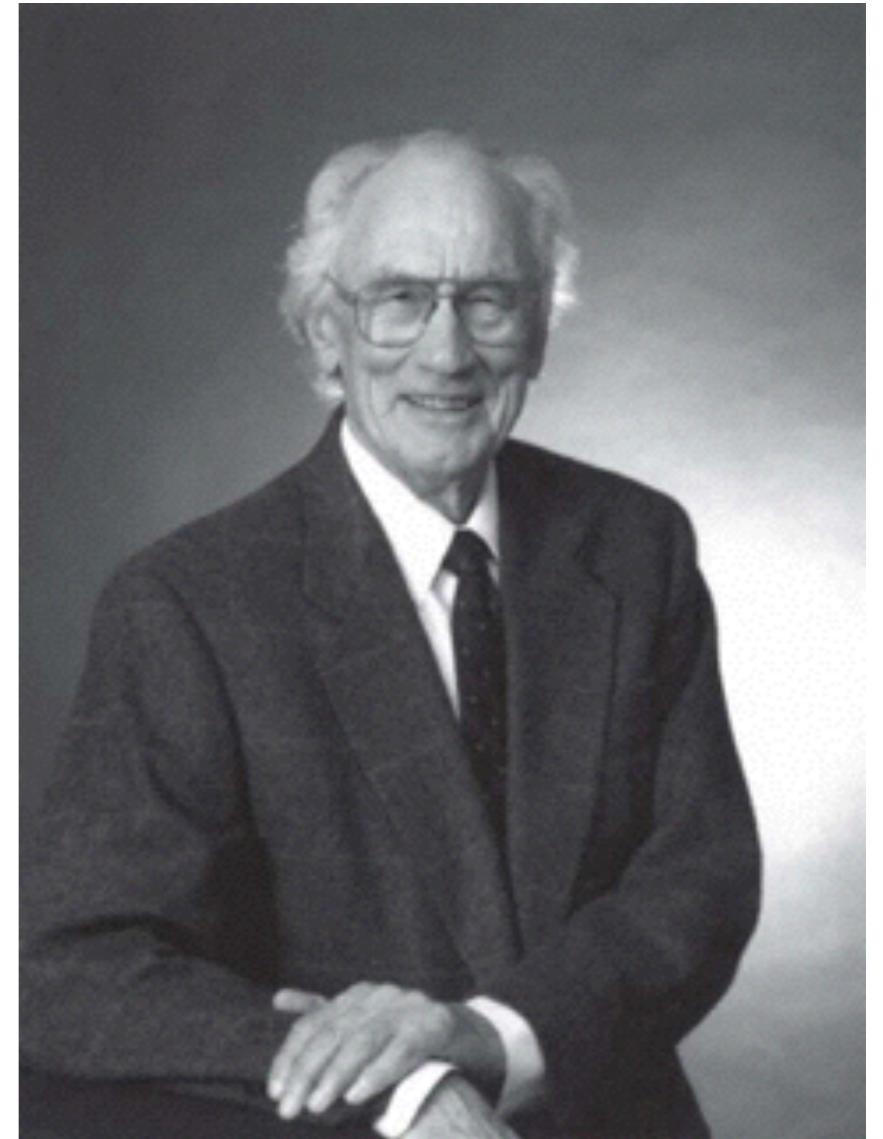


# Kenneth Pike

**Etic - Emic distinction**  
(phonetic - phonemic)

etic: universal/comparative level

emic: language-specific level



# Be aware when naming things!

- etic ~ comparative concepts
  - ▶ use lower-case (“the perfect”)
- emic ~ descriptive categories
  - ▶ use upper-case, like proper names (“the Perfect”)
  - ▶ add language names (“the English Perfect”)

Comrie, Bernard. 1976. Aspect: an introduction to the study of verbal aspect and related problems. Cambridge; New York: Cambridge University Press.  
Haspelmath, Martin. 2010. Comparative concepts and descriptive categories in crosslinguistic studies. *Language* 86(3). 663-687.

- Claim about the world's languages, e.g.
  - ▶ *“In all languages with a dative and an accusative case, the dative case marker is at least as long as the accusative case marker.”* (Haspelmath 2010: 665)
- Needs **etic** definition, e.g.
  - ▶ *“A dative case is a morphological marker that has among its functions the coding of the recipient argument of a physical transfer verb (such as ‘give’, ‘lend’, ‘sell’, ‘hand’), when this is coded differently from the theme argument.”*
- This includes:
  - ▶ German Dative, Russian Dative
  - ▶ Finnish Allative, Tsez Lative
- This does not include:
  - ▶ Nivkh Dative-Accusative  
(used e.g. in causative constructions)

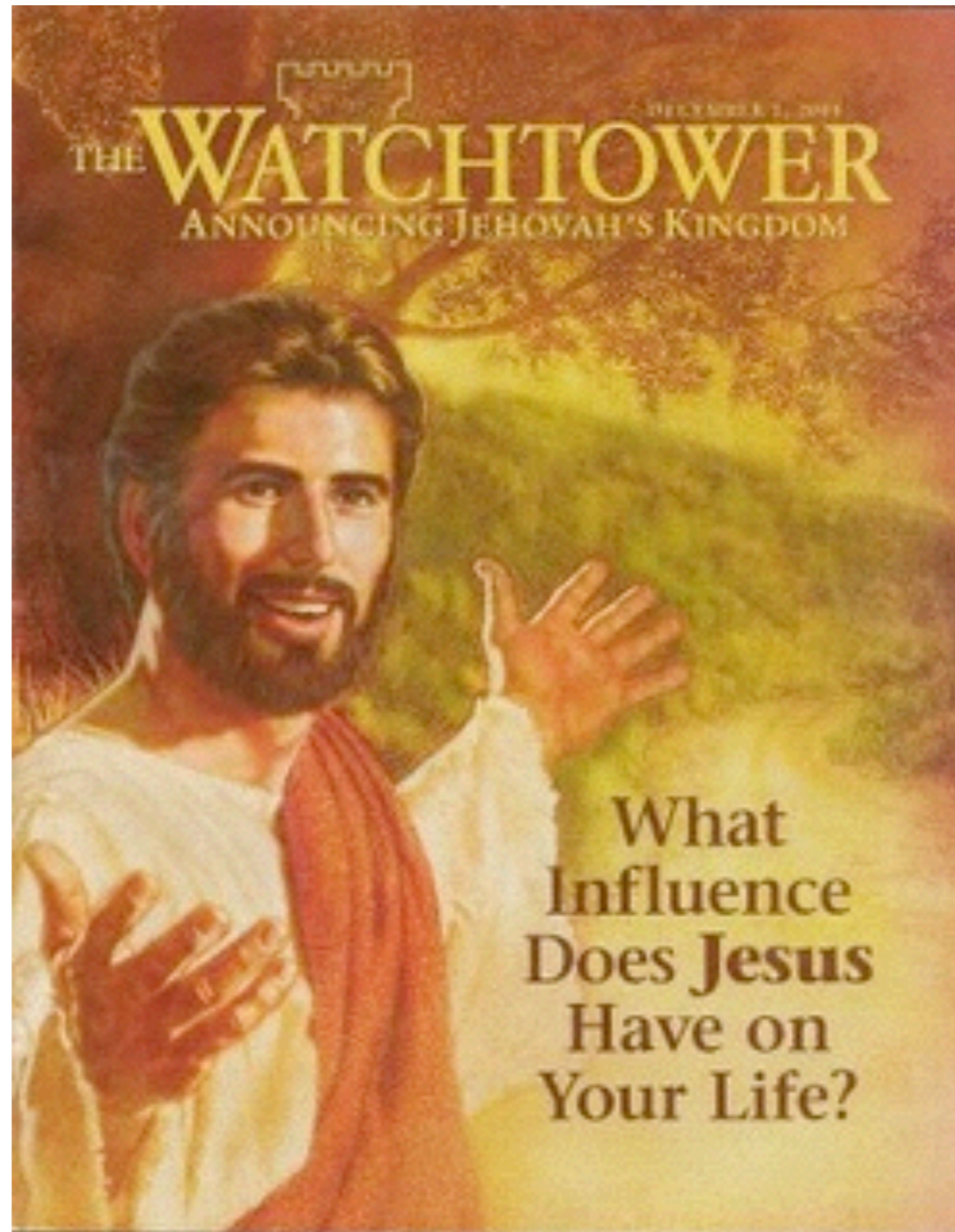
# Lessons from worldwide language diversity

- Possible vs. impossible languages
- Universal categories
- **The problem of comparing languages**

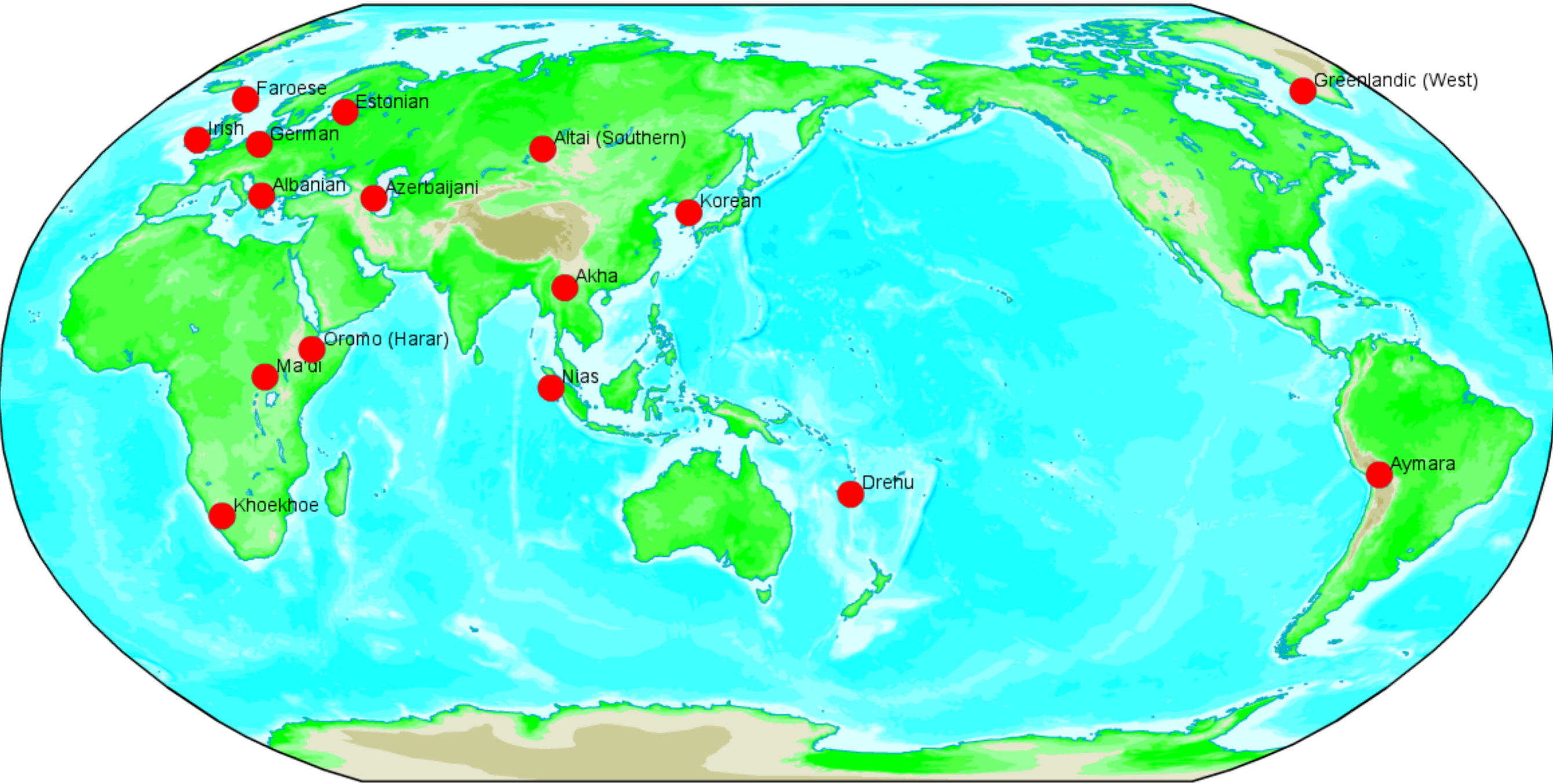


# How to compare languages

- How to compare like with like ?
  - ▶ Solution:  
use **etic definitions**
  - ▶ Extreme etic:  
use **contextually situated utterances**
  - ▶ To get the same contexts in many languages:  
use **parallel data**



- 1 What important information is contained in the Bible?
- 2 Who is the Bible's author?
- 3 Why should you study the Bible?
- 4 The Bible is a precious gift from God.
- 5 The Bible alone tells us what we must do to please God.
- 6 The Bible was written by some 40 different men over a period of 1,600 years, beginning in 1513 B.C.E.
- 7 So God in heaven, not any human on earth, is the Author of the Bible.
- 8 God made sure that the Bible was accurately copied and preserved.
- 9 More Bibles have been printed than any other book.
- 10 Not everyone will be happy to see you studying the Bible, but do not let that stop you.
- 11 But the Bible tells us that there is only one TRUE God.
- 12 But when the Bible was written, the name Jehovah appeared in it some 7,000 times
- 13 God is a Spirit, says the Bible.
- 14 The Bible reveals Jehovah's personality to us.
- 15 The Bible tells us that he is also merciful, kind, forgiving, generous, and patient.
- 16 We learn about God from creation and from the Bible.
- 17 Another way we can learn about God is by studying the Bible.
- 18 By disobeying God's command, the first man, Adam, committed what the Bible calls sin.
- 19 This is what the Bible refers to as the ransom.
- 20 Some of your loved ones may become very angry because you are studying the Bible.
- 21 What is the Bible's view of separation and of divorce?
- 22 The Bible says that a husband is the head of his family.
- 23 Parents need to spend time with their children and study the Bible with them,
- 24 When marriage mates have problems getting along together, they should try to apply Bible counsel.
- 25 The Bible urges us to show love and to be forgiving.
- 26 But God does not approve of them if they come from false religion or are against Bible teachings.
- 27 The only two birthday celebrations spoken of in the Bible were held by persons who did not worship Jehovah.
- 28 The Bible teaches that only a few people are on the narrow road to life.
- 29 The Bible foretold that after the death of the apostles, ...
- 30 True Christians love one another, respect the Bible, and preach about God's Kingdom.
- 31 Another mark of true religion is that its members have a deep respect for the Bible.
- 32 They try to live by the Bible in their everyday life.
- 33 The Bible is the basis for what is taught.
- 34 By now you have learned many good things from the Bible.



Albanian

*bibla*

Nominative

*biblën*

Accusative

*biblës*

Genitive/Dative

...

Faroese

*biblian*

Nominative

*bibliuna*

Accusative

*bibliunnar*

Genitive

*bibliuni*

Dative

...

Estonian

*piibel*

Nominative

*piiblit*

Partitive

*piibli*

Genitive

*piiblis*

Inessive

*piiblist*

Elativ

...

Greenlandic

*biibilip*

Ergative

*biibli*

Absolutive

*biibilmik*

Instrumental

*biibilmi*

Locative

...

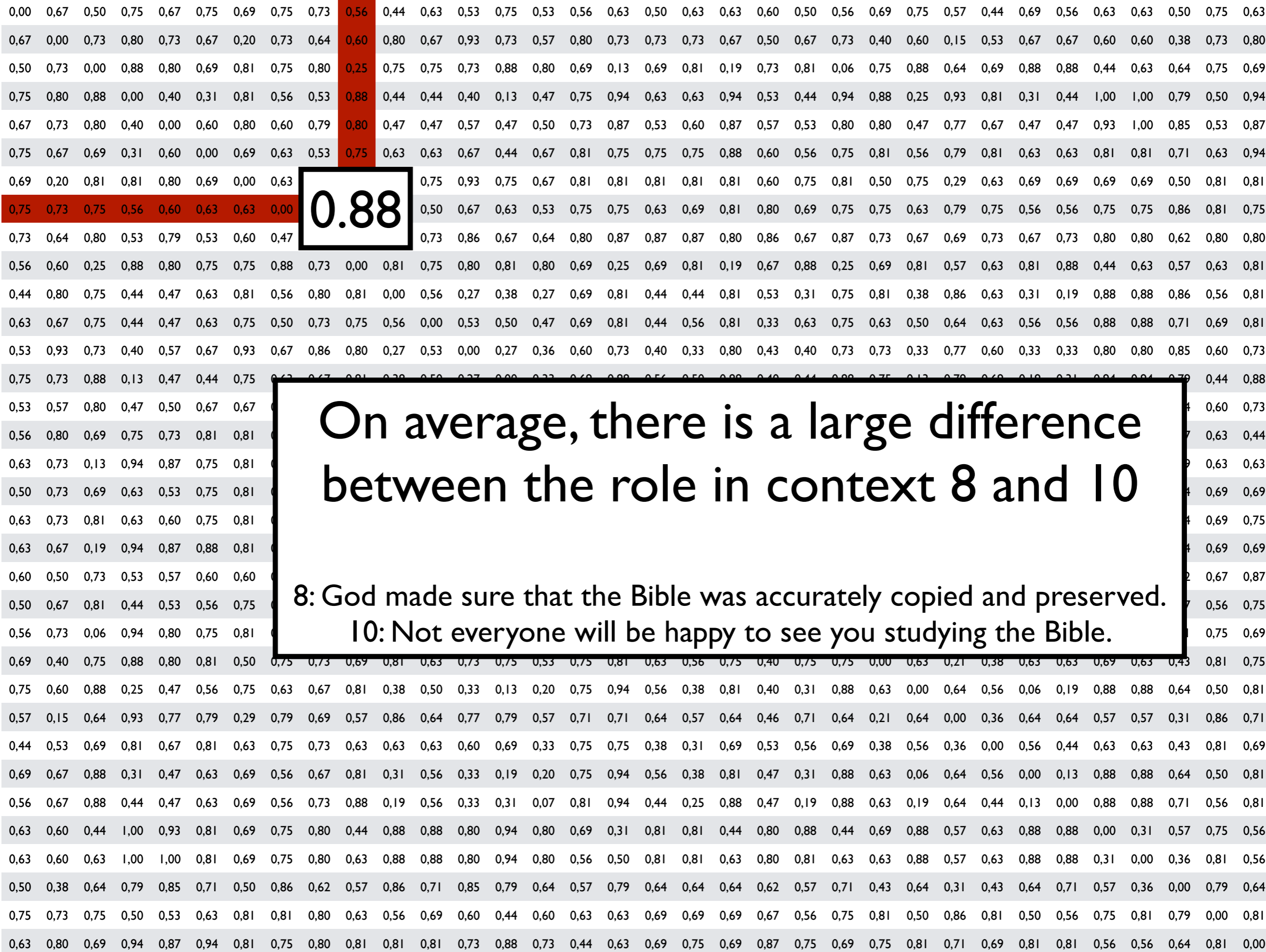
Context	Albanian	Faroese	Estonian	Greenlandic
1	bibla	bíbliuni	piibel	biibili
2	biblës	bíbliunnar	piibli	biibilimik
3	biblën	bíbliuna	piiblit	biibili
4	bibla	bíblían	piibel	biibili
5	bibla	bíblían	piibel	biibilip
6	bibla	bíbliuna	piibli	biibili
7	biblës	bíbliunnar	piibli	biibilimut
8	bibla	bíblían	piiblit	biibilip
9	bibla	NA	piiblit	biibili
10	biblën	bíbliuna	piiblit	biibilimik
11	bibla	bíblían	piibel	biibilimili
12	bibla	bíblían	piibel	biibilili
13	bibla	bíblían	piibel	biibilimi
14	bibla	bíblían	piibel	biibilimi
15	bibla	bíblían	piibel	biibilimi
16	bibla	bíbliuni	piibli	biibililu
17	biblën	bíbliuna	piiblit	biibilimik
18	bibla	bíblían	piiblis	biibilip
19	bibla	bíblían	piiblis	biibilimi
20	biblën	bíbliuna	piiblit	biibilimik
21	NA	bíblían	piibel	biibilimi
22	bibla	bíbliuni	piibel	biibili
23	biblën	bíbliuna	piiblit	biibilimillu
24	biblike	bíblían	piibli	biibilimi
25	bibla	bíblían	piibel	biibilimi
26	biblës	bíbliunnar	piibli	biibilimi
27	bibla	bíblían	piiblis	biibilimi
28	bibla	bíblían	piibel	biibilimi
29	bibla	bíblían	piibel	biibilimi
30	biblën	bíbliuna	piiblist	biibilimik
31	biblën	bíbliuni	piibli	biibilimik
32	biblës	bíbliuni	piibli	biibili
33	bibla	bíbliuna	piibel	biibilimik
34	bibla	bíbliuni	piiblist	biibilimeersunik







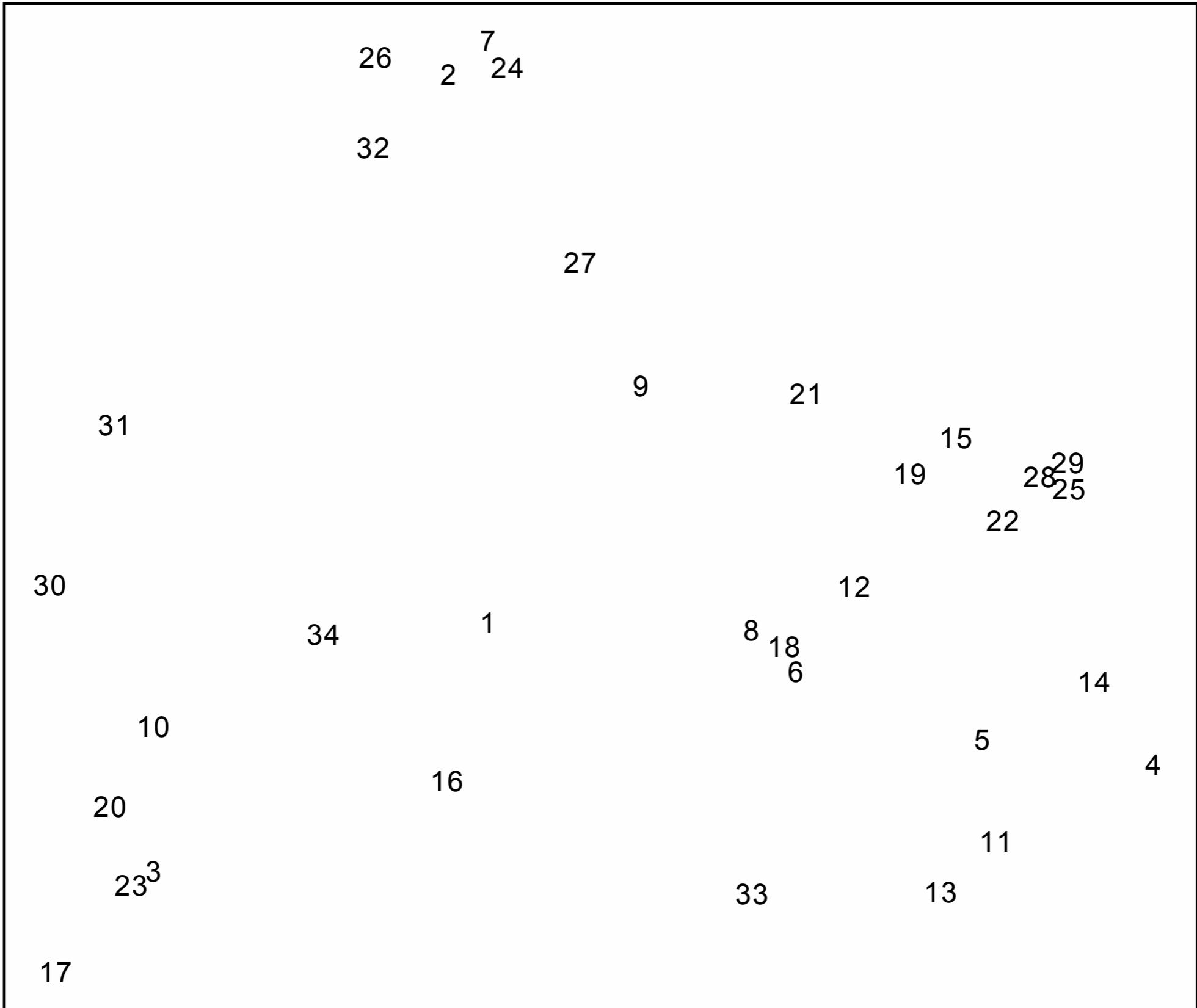




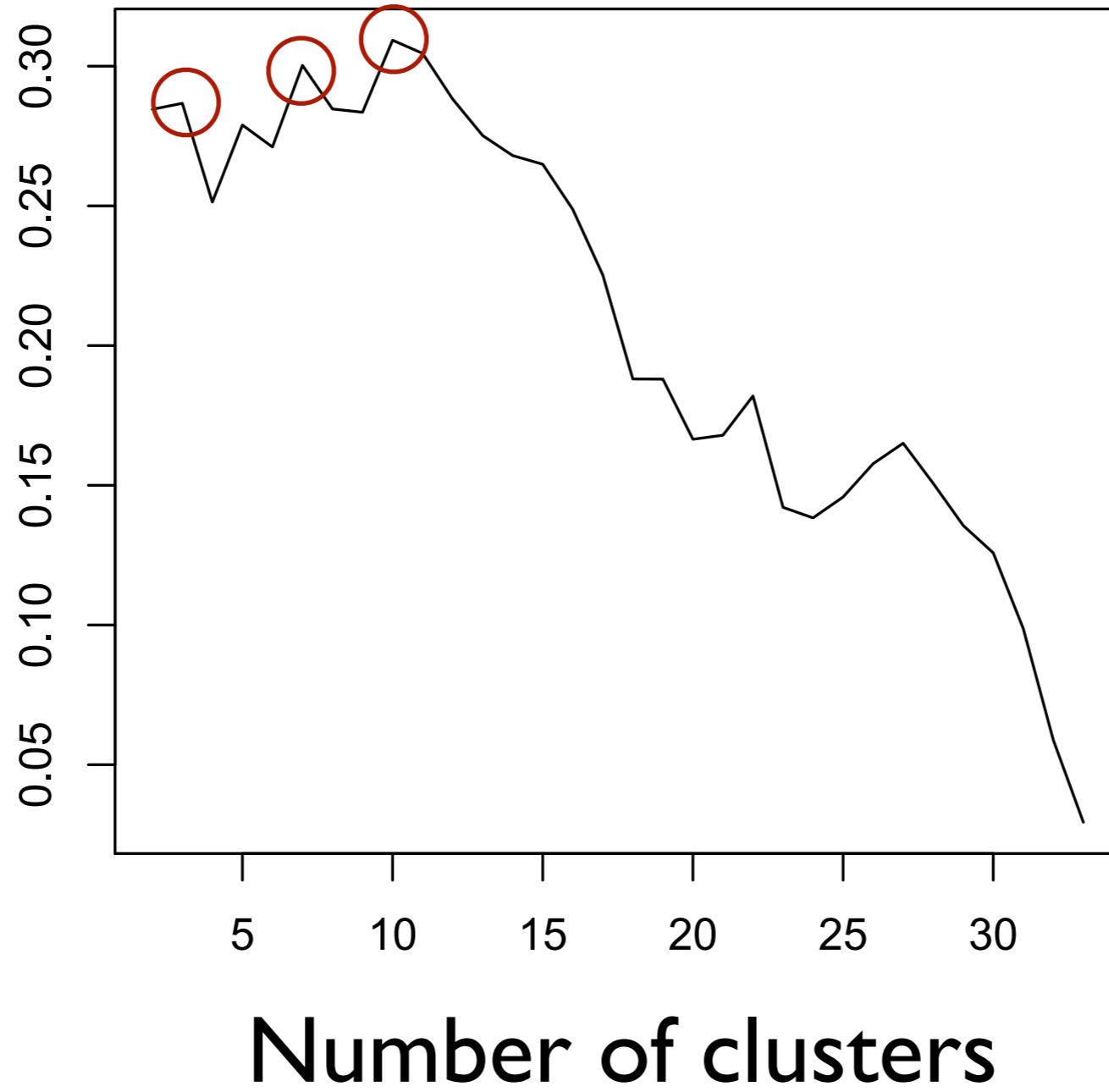
0.88

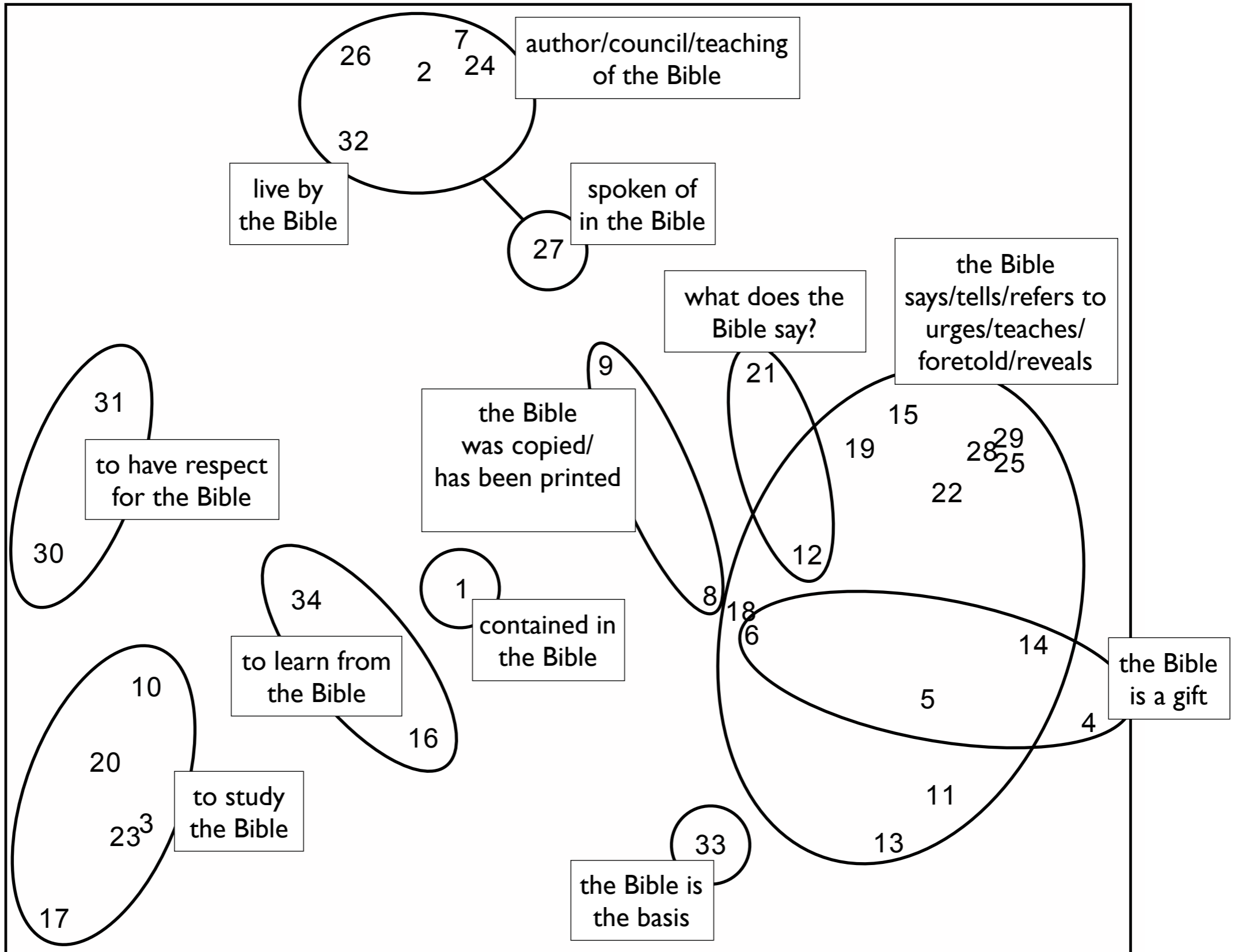
On average, there is a large difference  
 between the role in context 8 and 10

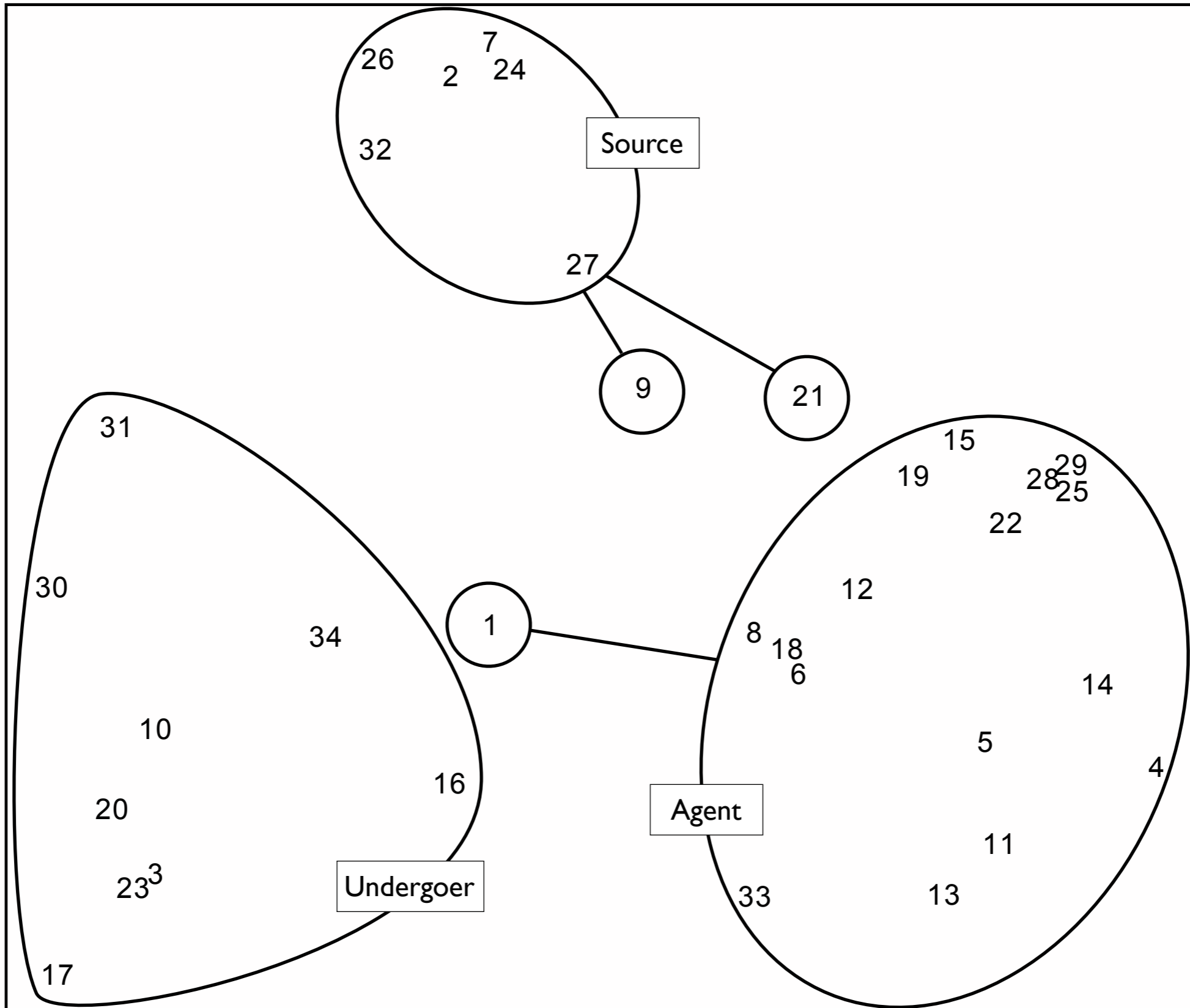
8: God made sure that the Bible was accurately copied and preserved.  
 10: Not everyone will be happy to see you studying the Bible.

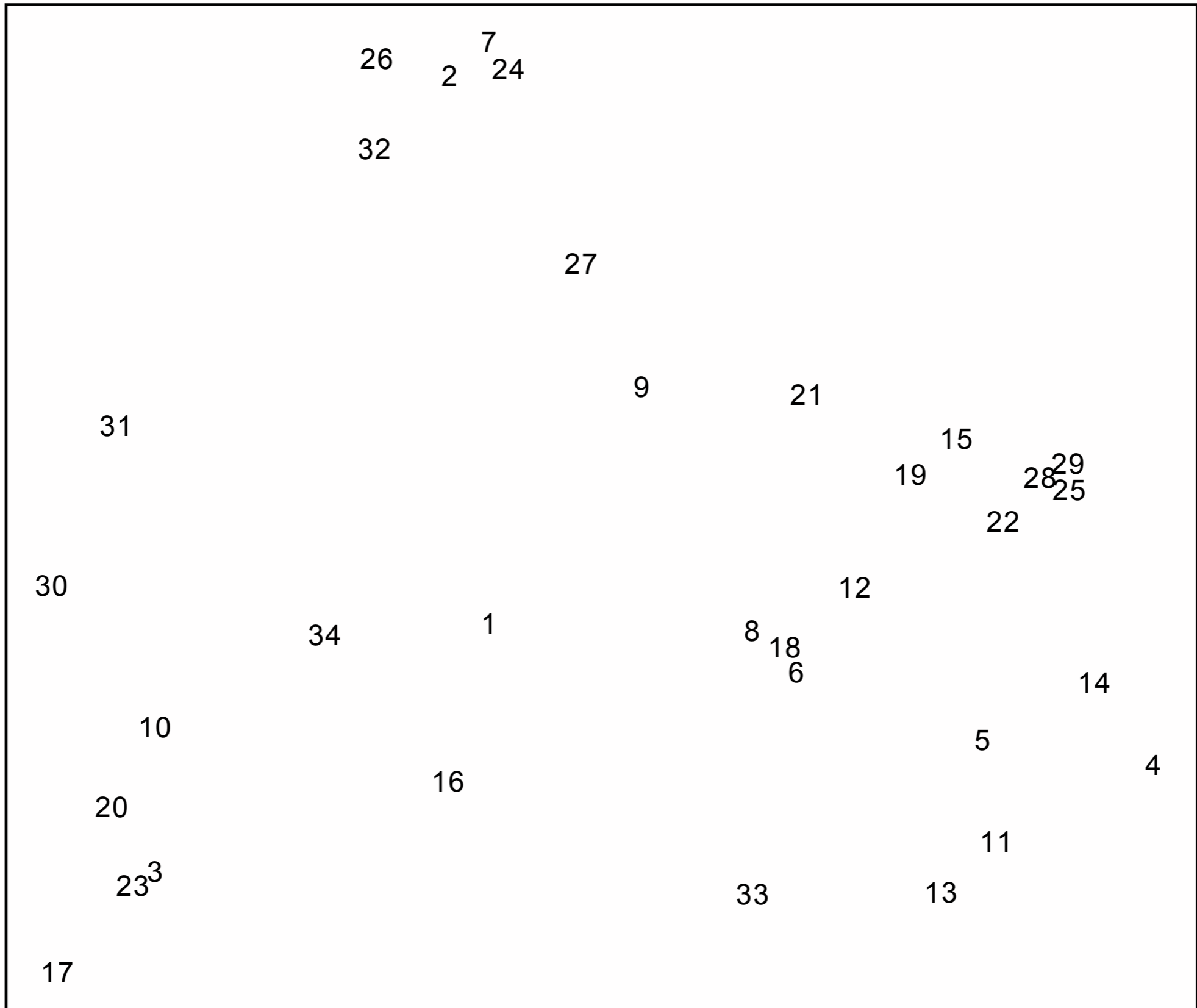


# “Fit” of clustering

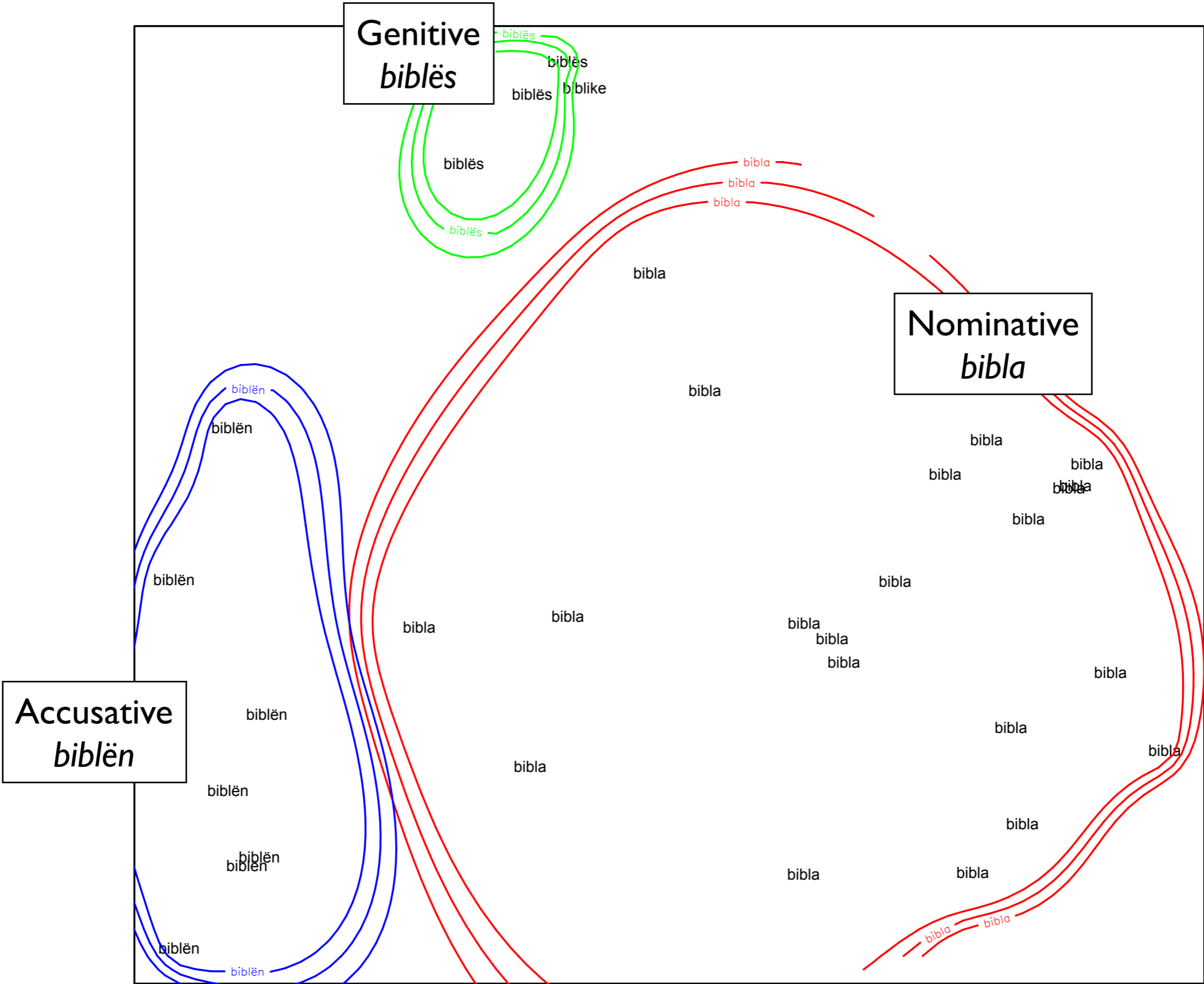






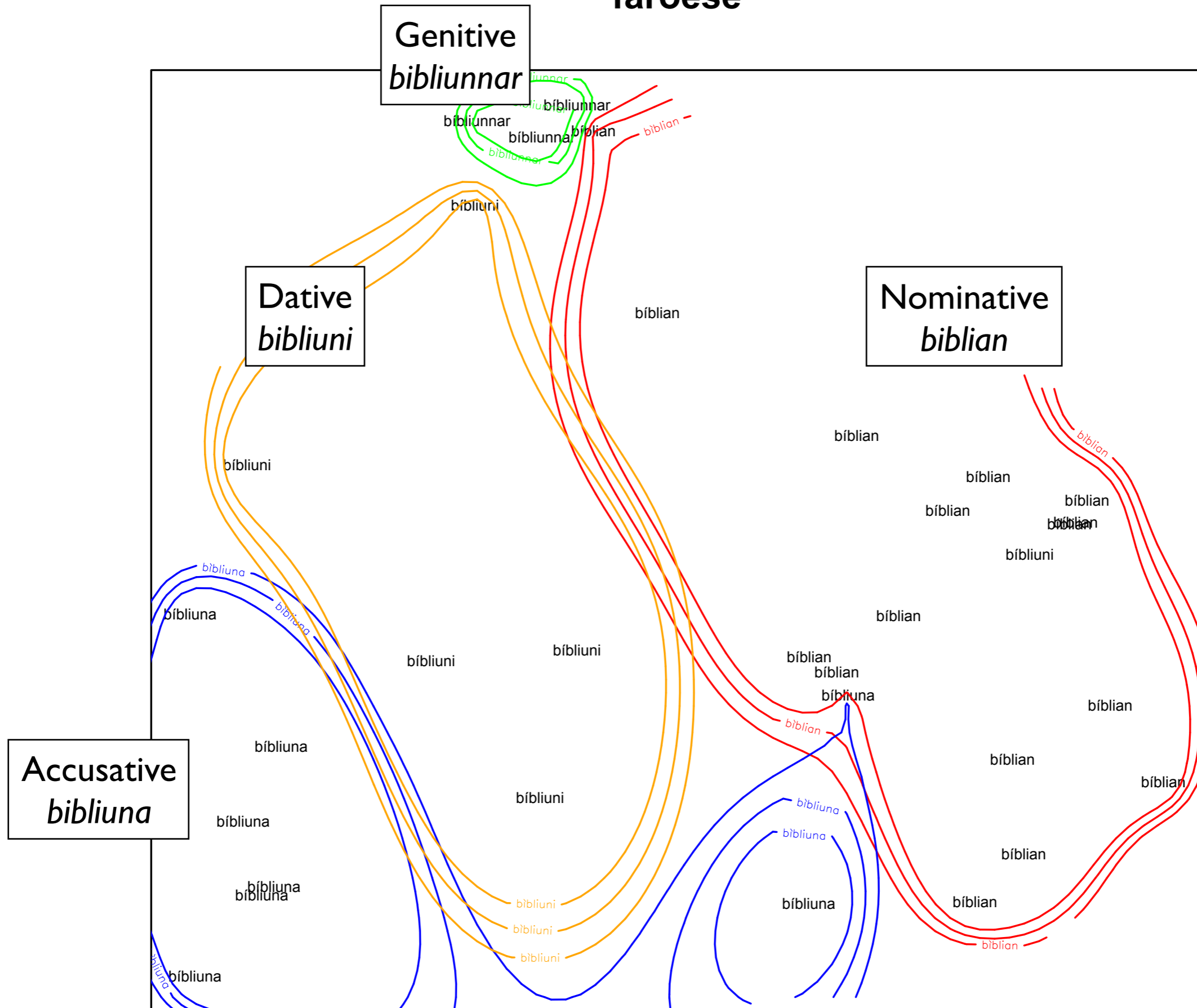


albanian





# faroese





A stylized map of Southeast Asia and Oceania. The landmasses are colored in shades of green and yellow, while the surrounding waters are light blue. Several colored dots are placed on the map: two yellow dots in the Philippines and Indonesia, one yellow dot in the Pacific, one red dot in the Malay Peninsula, one purple dot in Australia, and one yellow dot in New Zealand. A semi-transparent white rectangular box is centered over the map, containing the title text.

# Introducing the *Parallel Text Corpus*

# Parallel Bible Corpus

- 1169 translations
- 906 different ISO-639/3 codes
- In total more than 350 Million wordforms
- More than 17 Million different wordforms
- <http://paralleltext.info/data>

***Demo***

# 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/paralleltextprocessing>

# Multiple Alignment

- Based on sentence-by-sentence alignment, induce word-by-word alignment
- Translations can be (and often are!) quite different
- Bi-text alignment is widely researched problem
- Multitext alignment not so much (but multi-string alignment in bio-informatics is!)

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 .

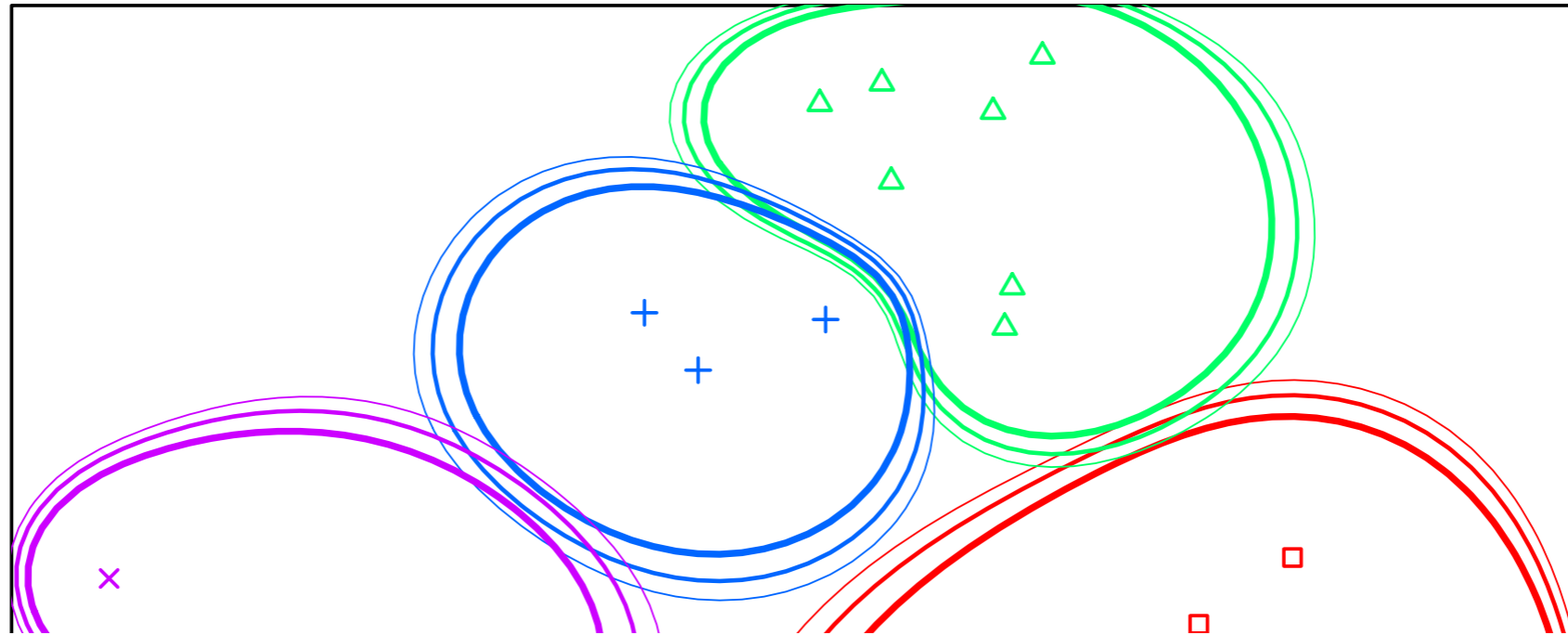
Der König Herodes war total aufgebracht , als er das hörte , und nicht nur er , alle in Jerusalem waren das .

# 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**

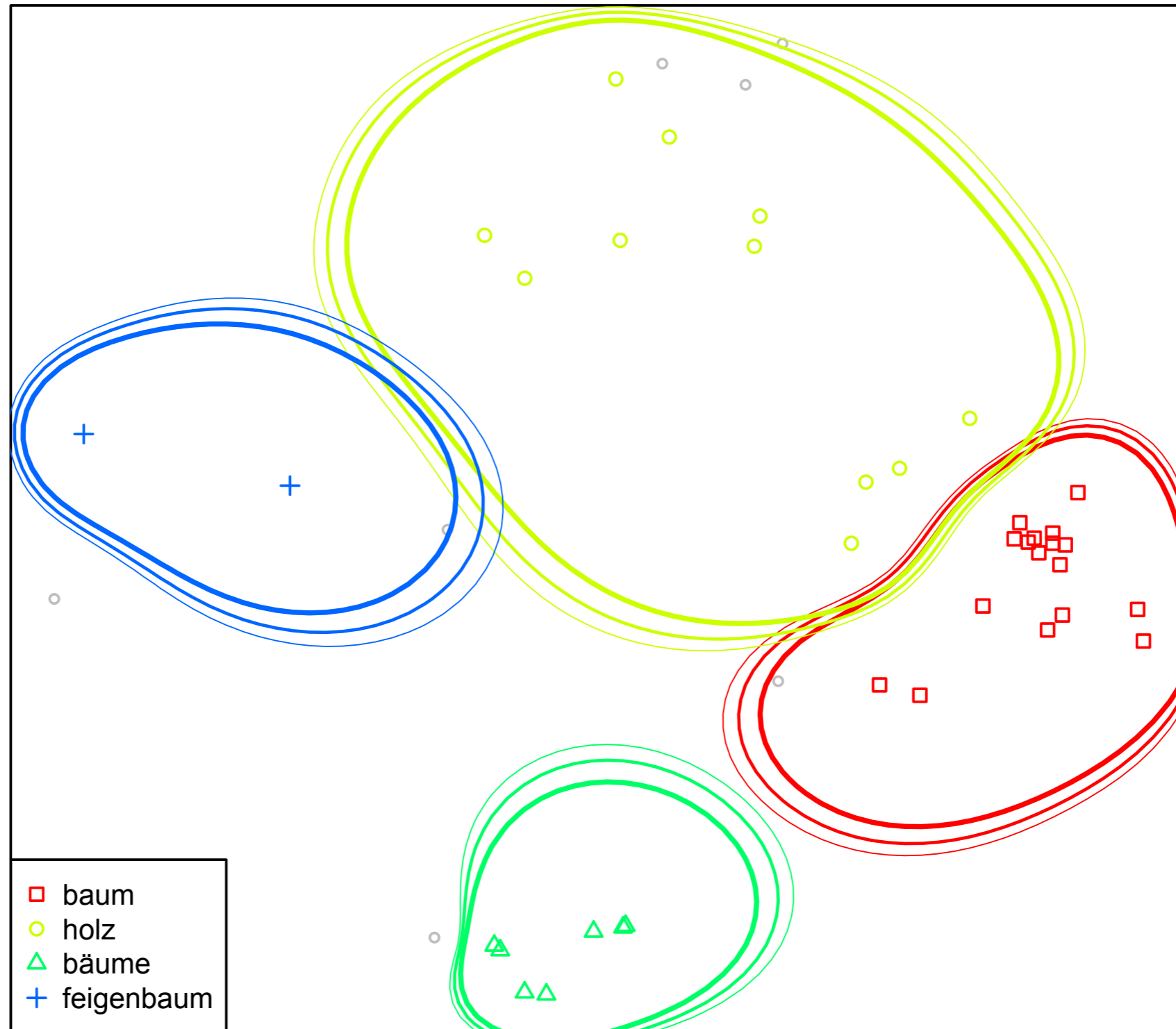
# afr-x-bible-1953.txt



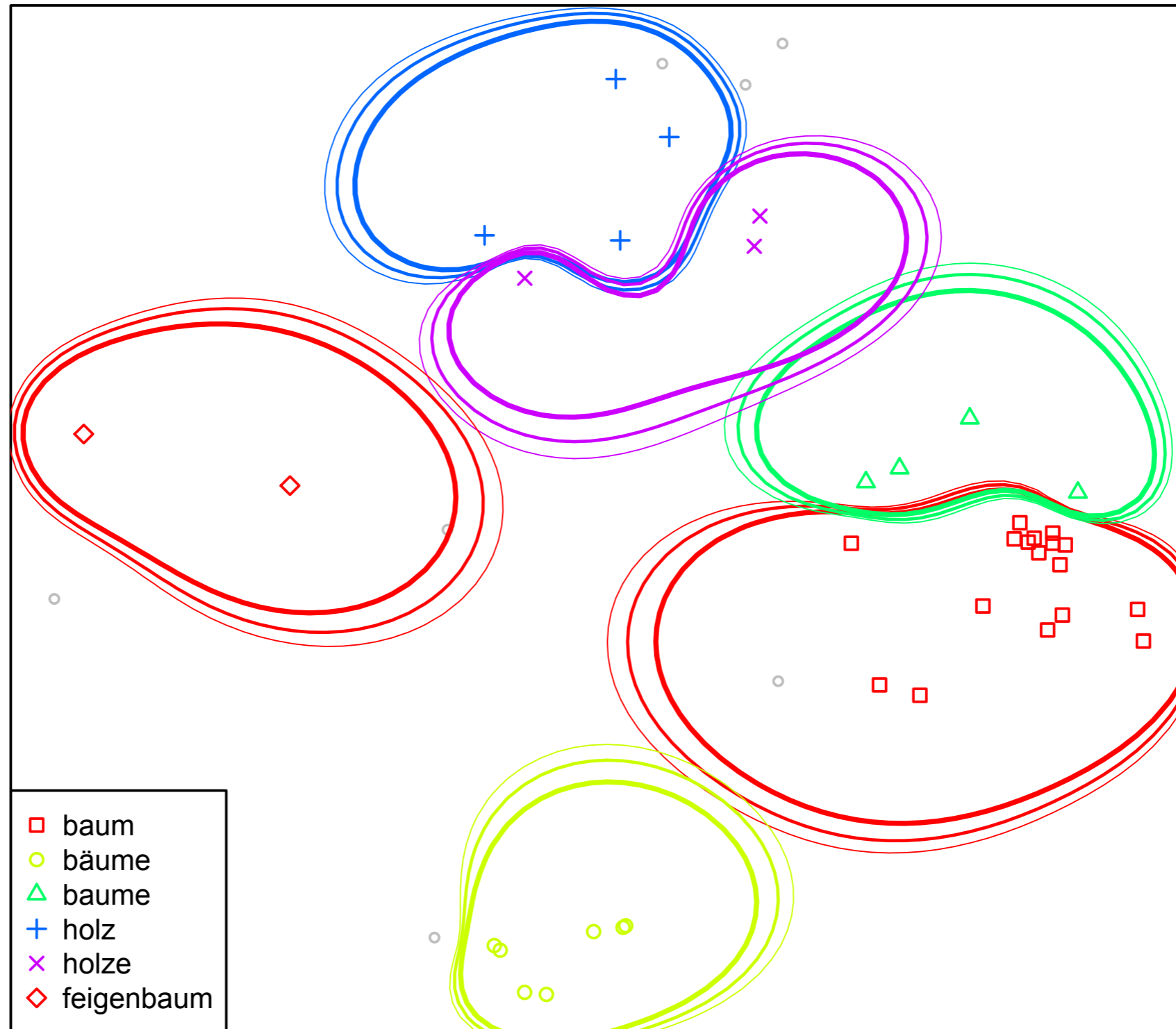
	tree	wood (stuff)	firewood	small forest	large forest
German	<u>Baum</u>	<u>Holz</u>		<u>Wald</u>	
Danish	<u>træ</u>			<u>skov</u>	
French	<u>arbre</u>	<u>bois</u>		<u>forêt</u>	
Spanish	<u>árbol</u>	<u>madera</u>	<u>leña</u>	<u>bosque</u>	<u>selva</u>

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

# deu-x-bible-erben.txt

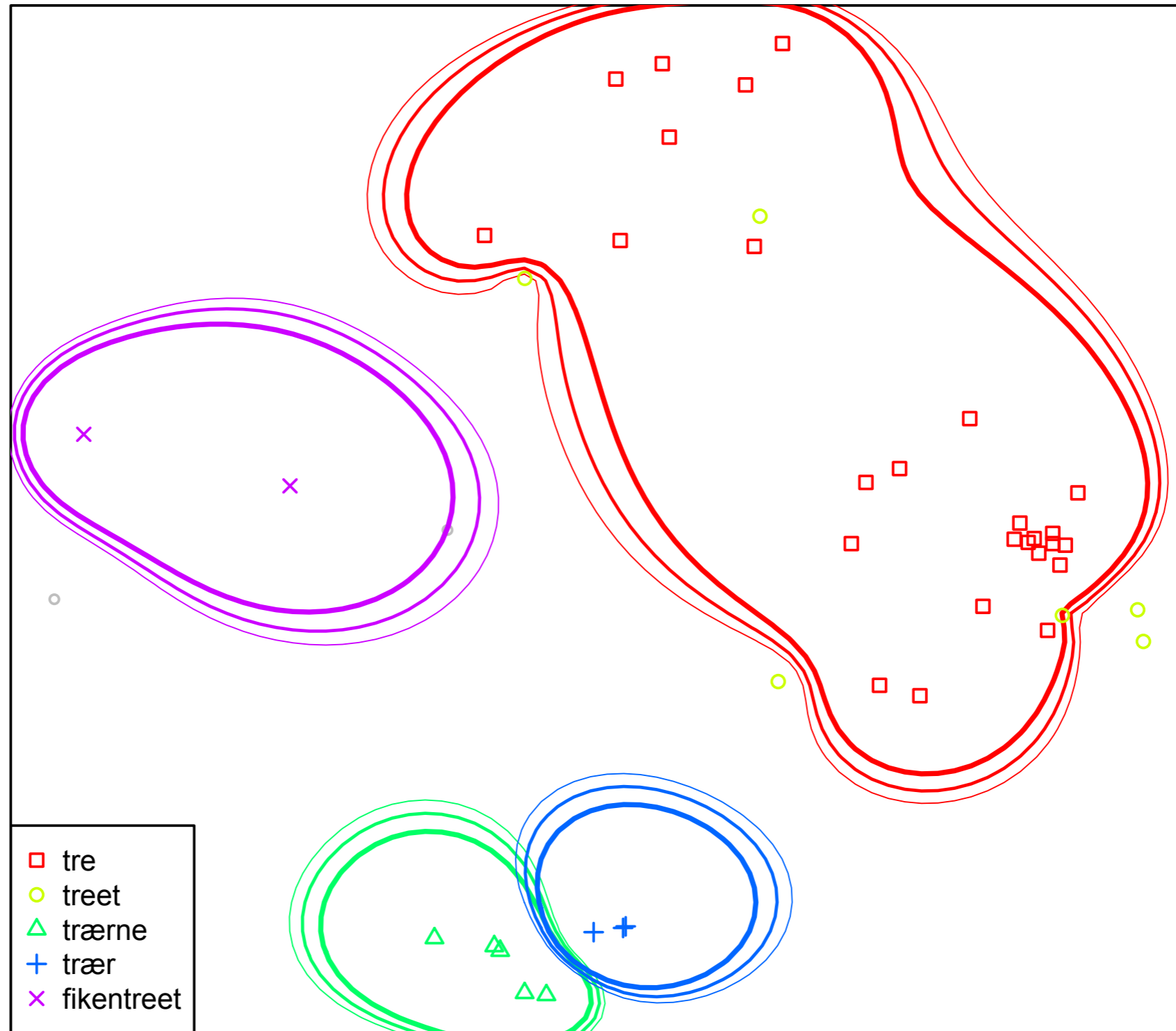


# deu-x-bible-freebible.txt



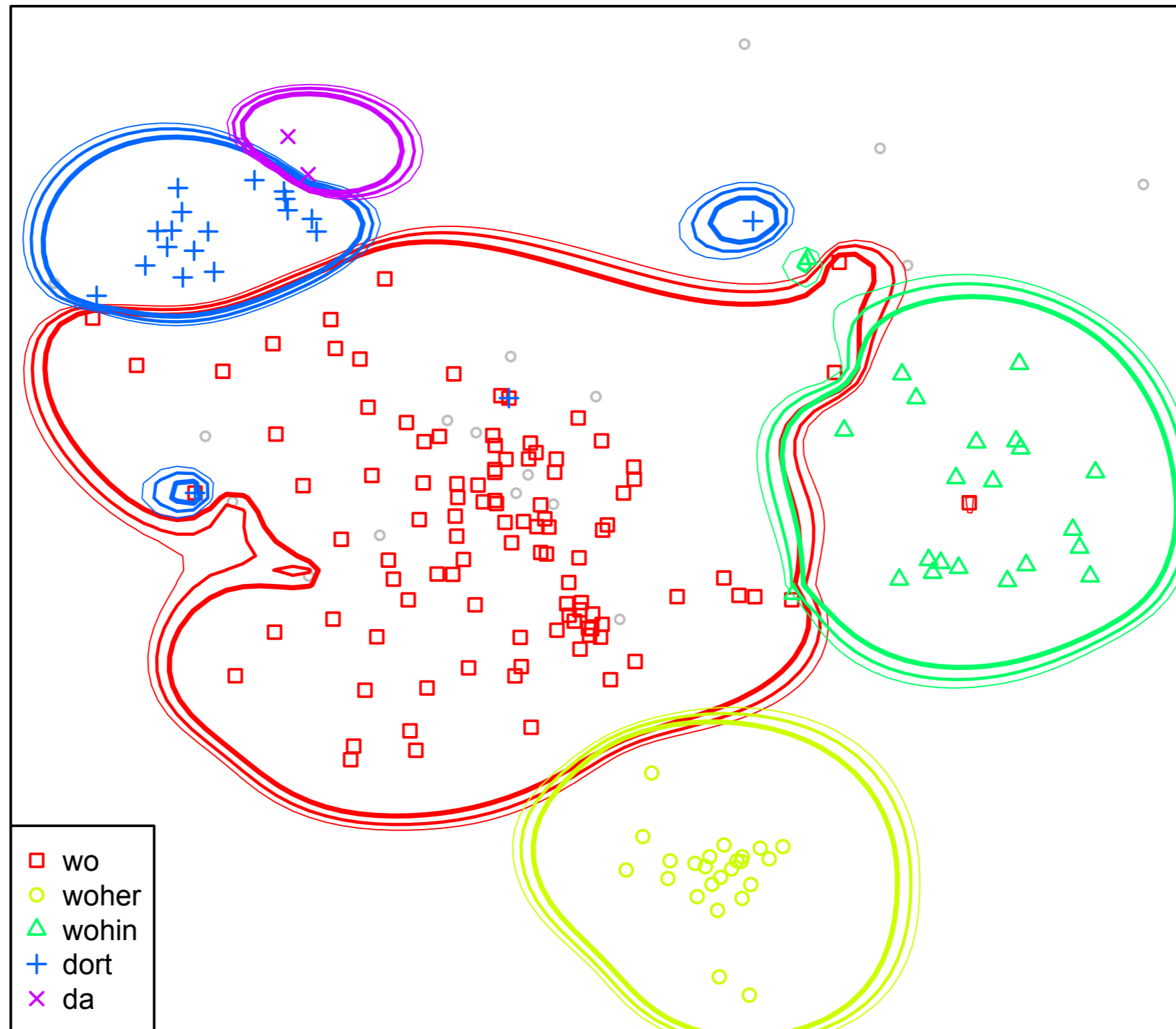


# nob-x-bible-2007.txt

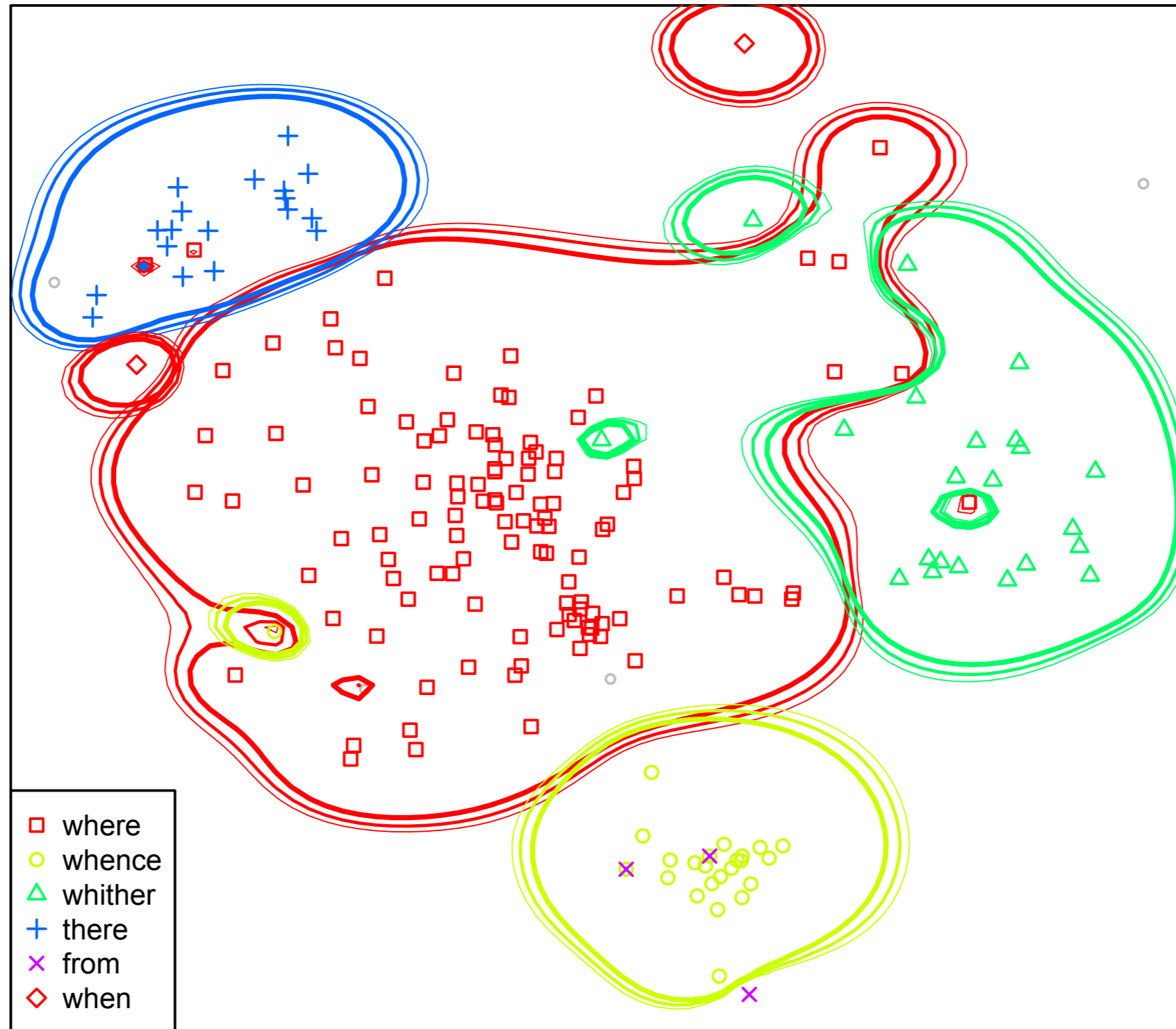


where

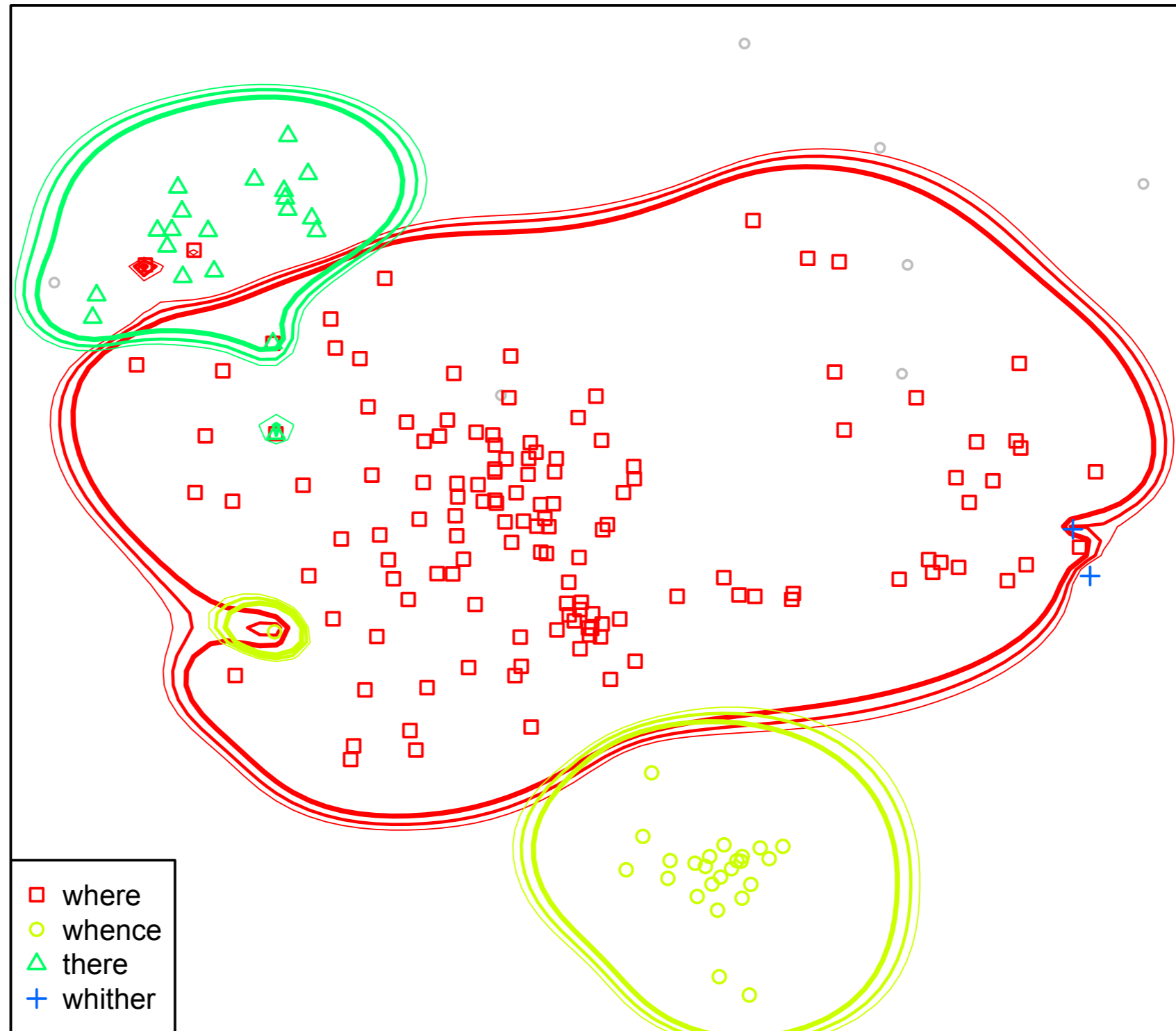
# deu-x-bible-pattloch.txt



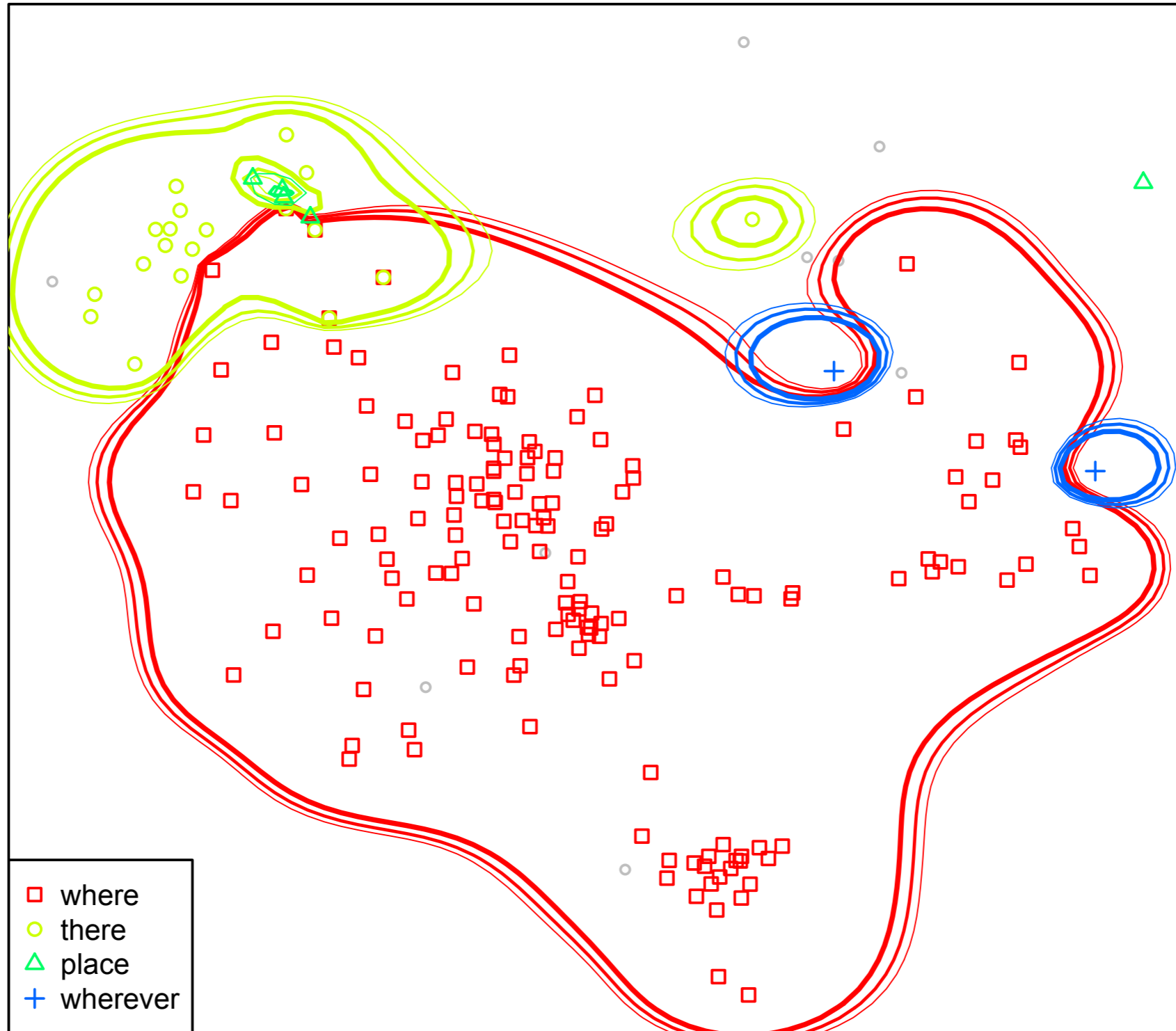
# eng-x-bible-kingjames.txt



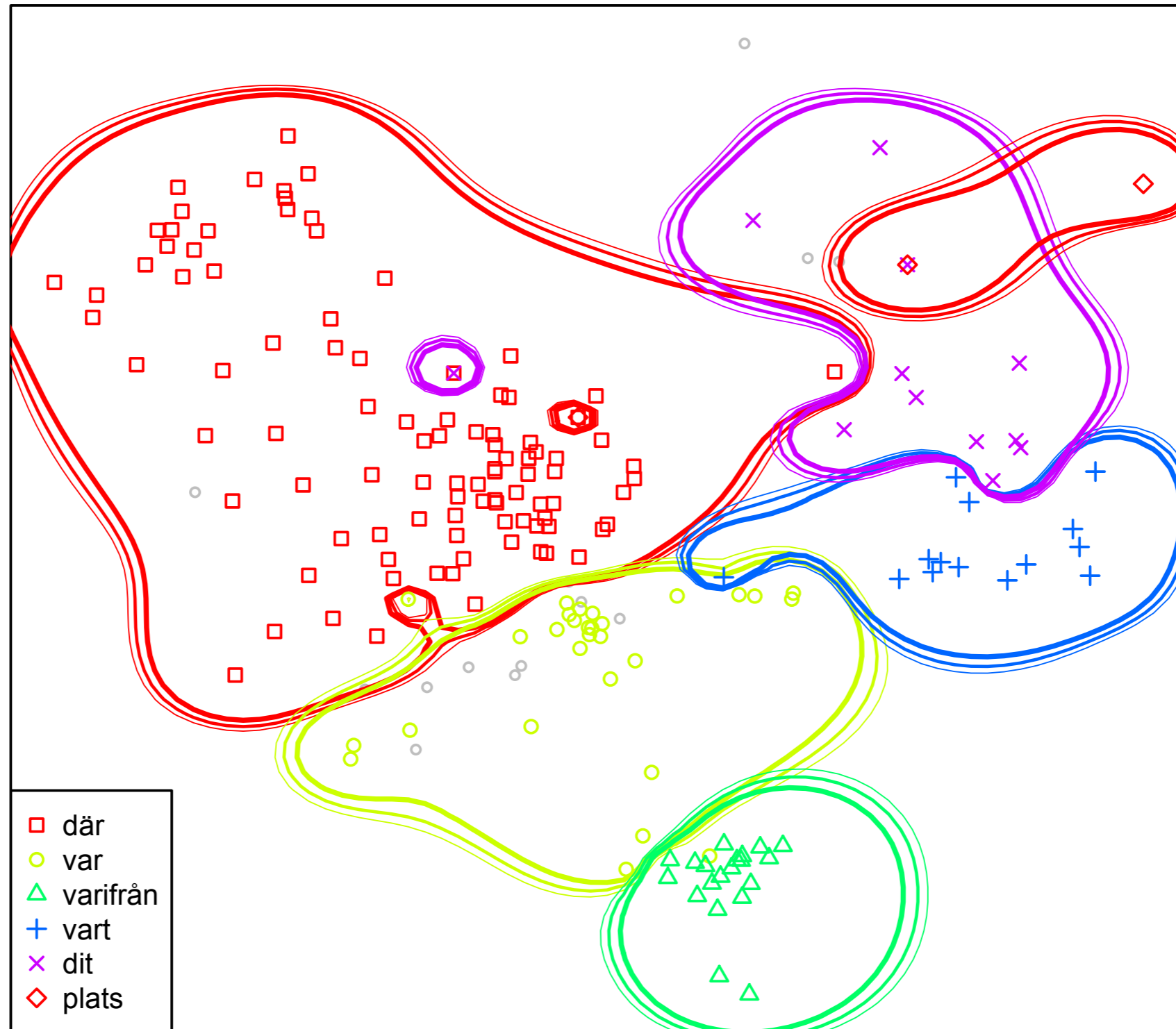
# eng-x-bible-darby.txt



# eng-x-bible-treeoflife.txt

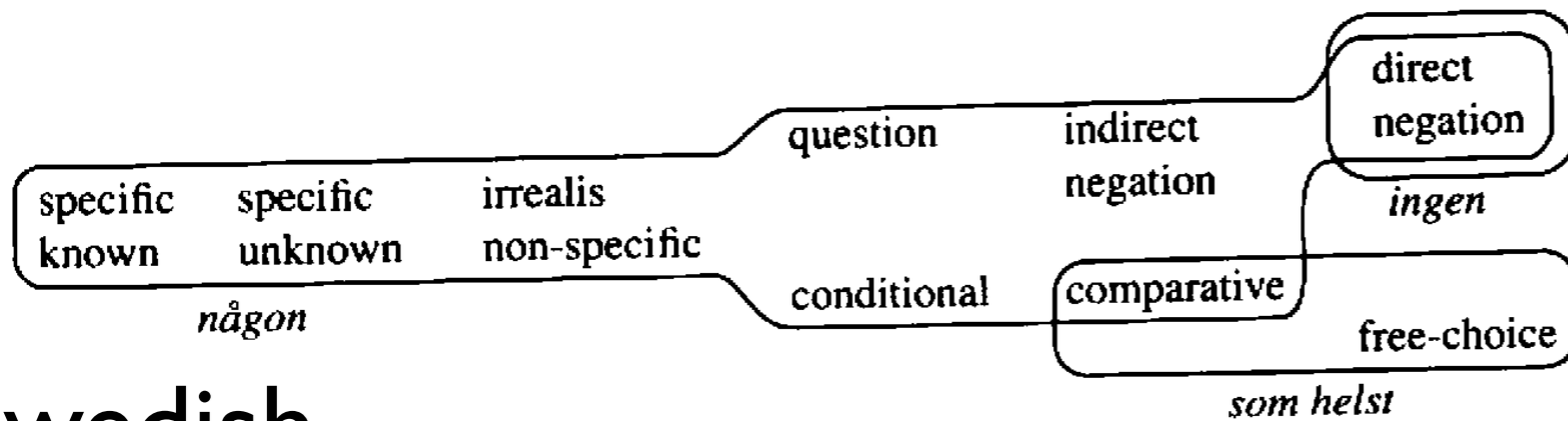
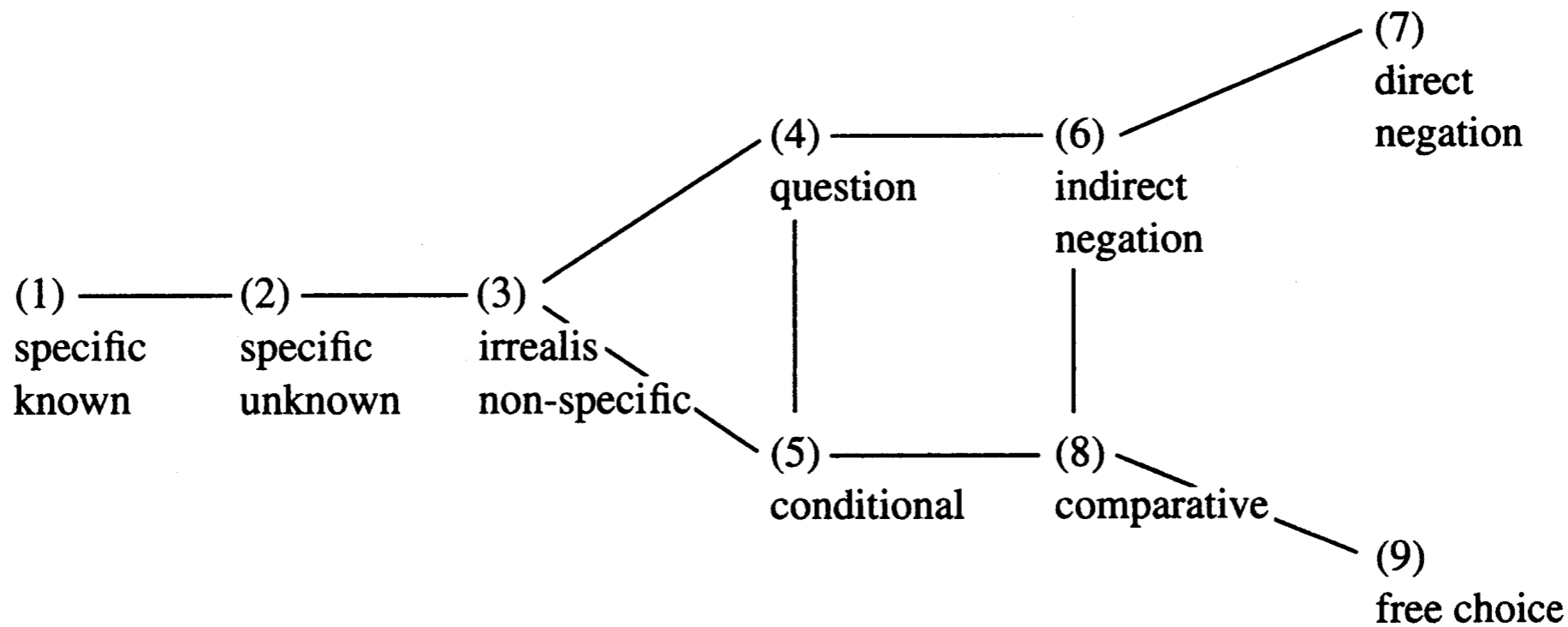


# swe-x-bible-folk1998.txt

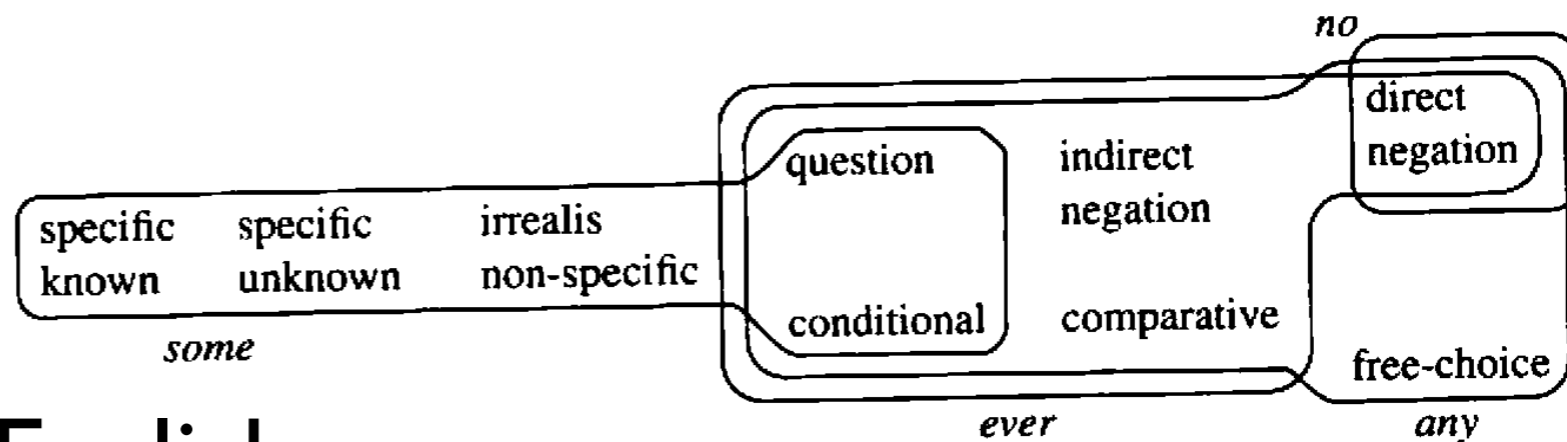


**Indefinite person  
(someone, anyone)**



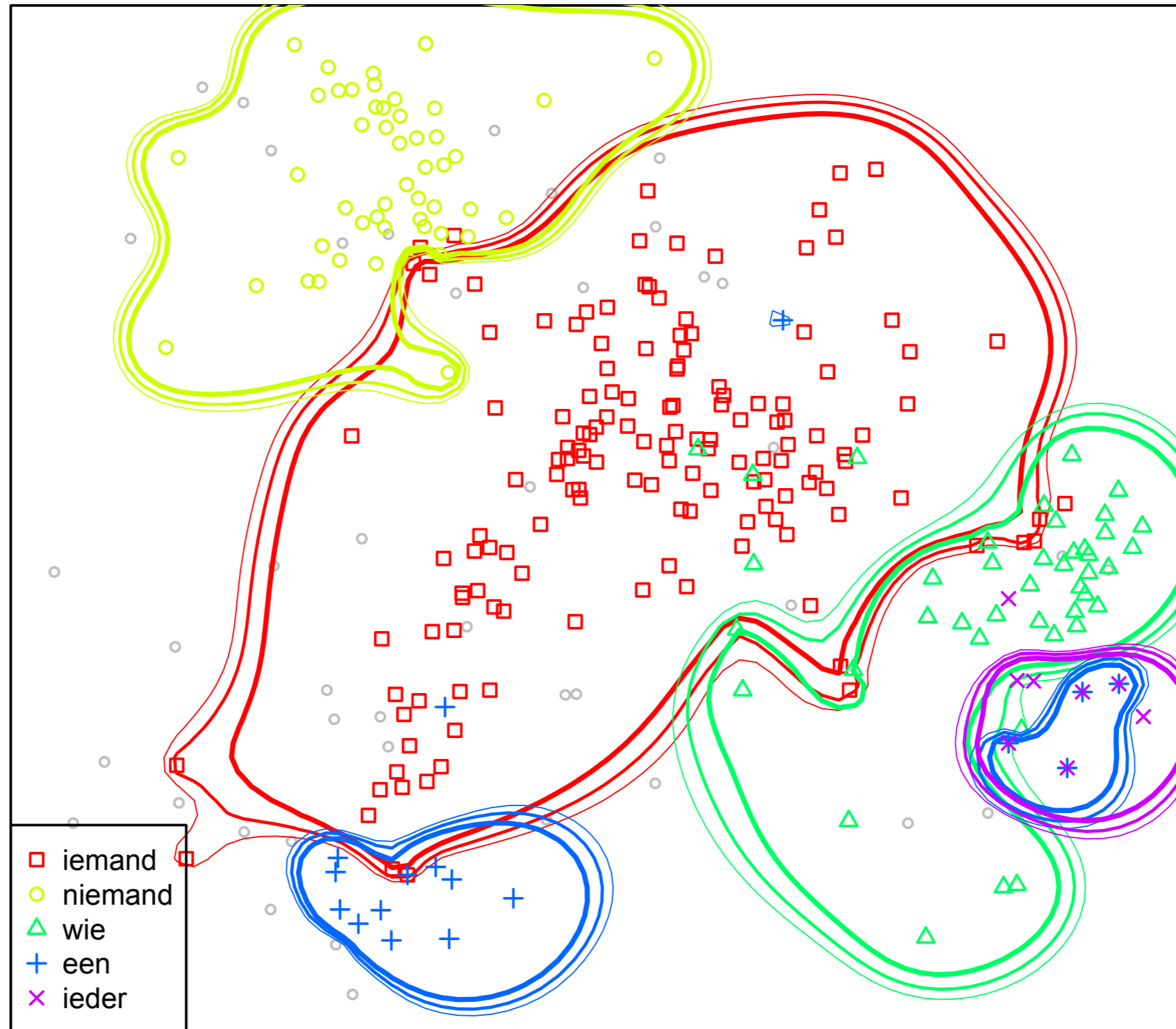


# Swedish

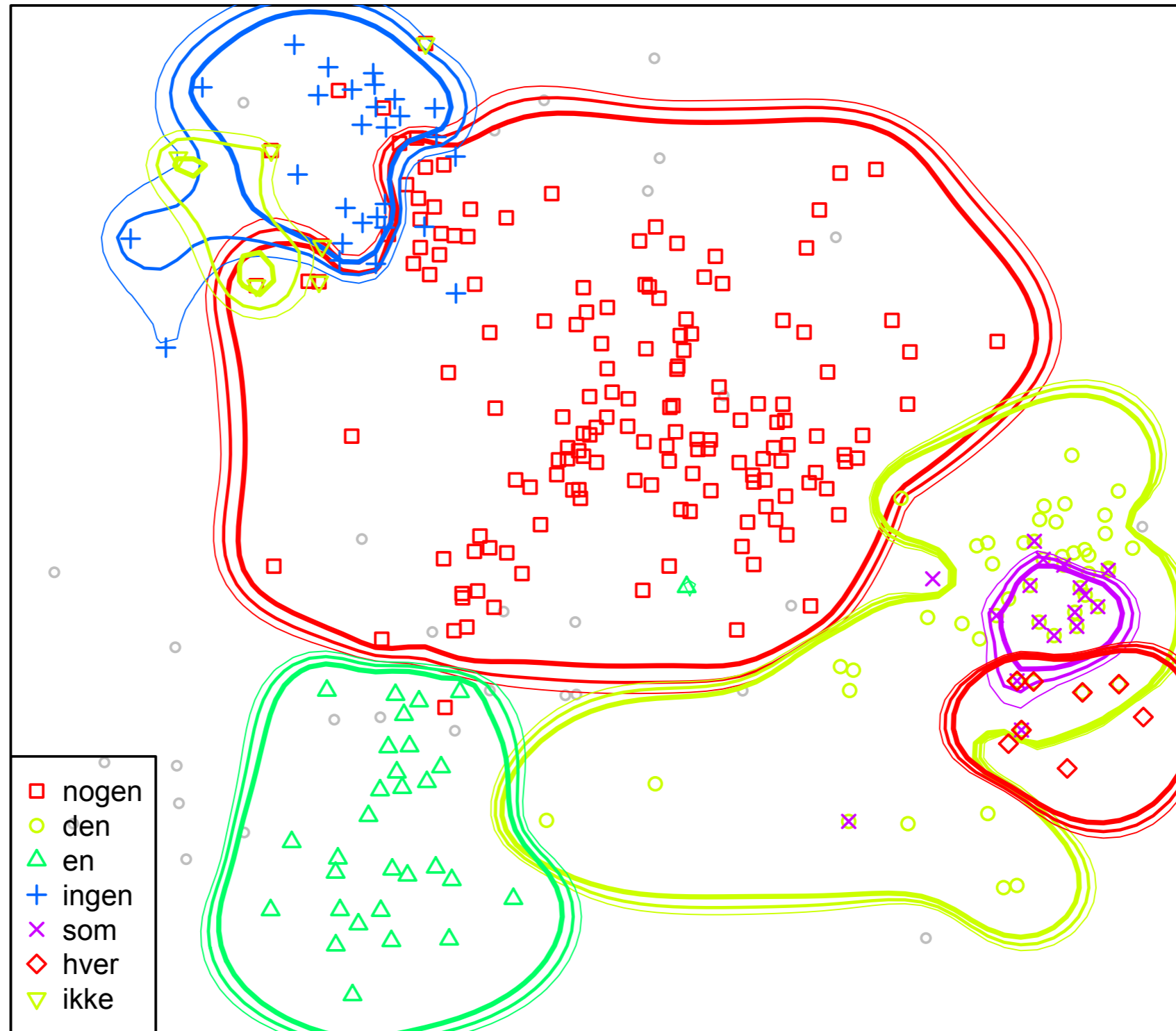


# English

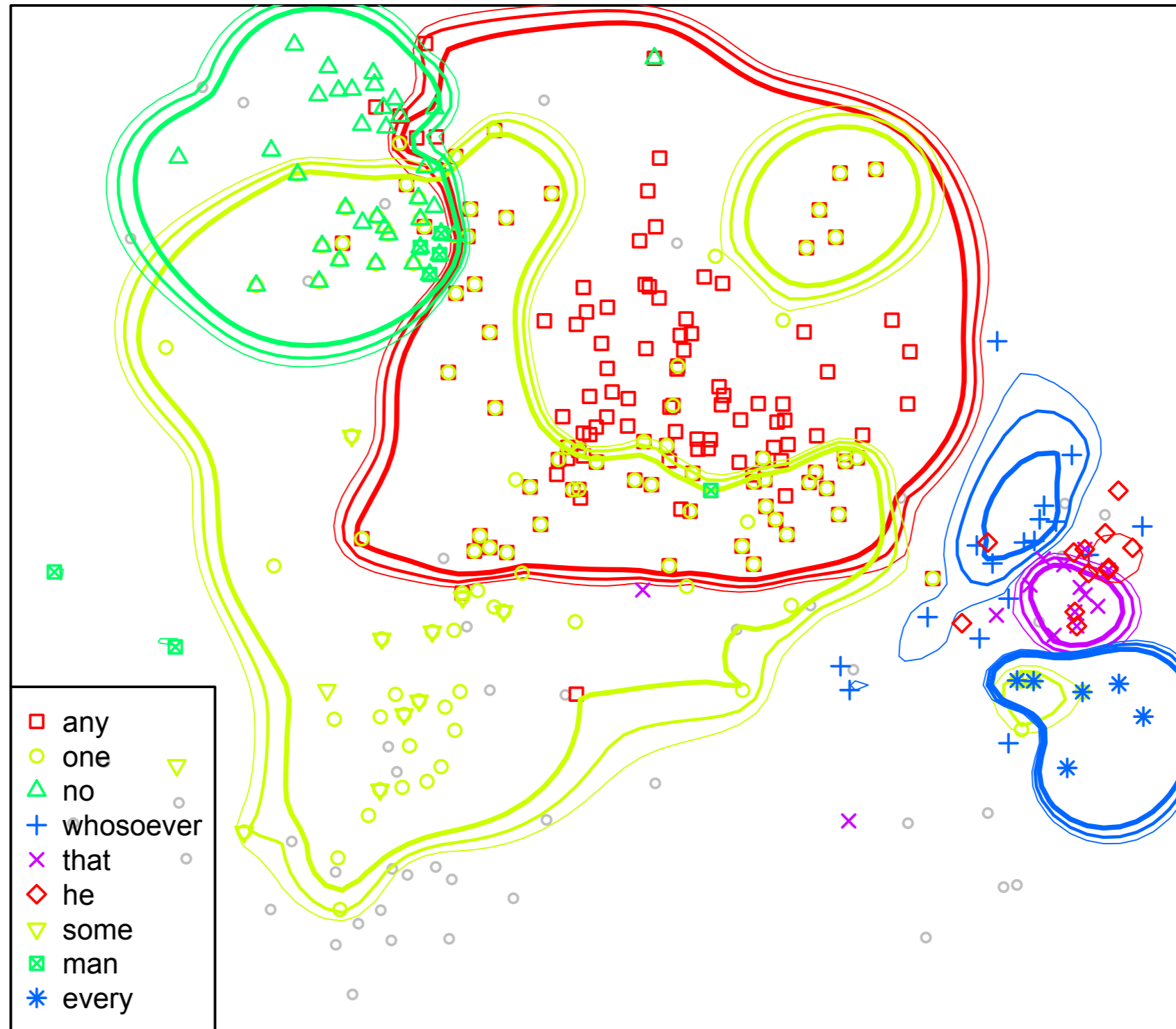
# nld-x-bible-1951.txt



# dan-x-bible-1931.txt



# eng-x-bible-darby.txt



# Conclusion

- Massively parallel texts are a goldmine for language comparison
- Much experimentation is needed to find suitable methods to enrich the data
- Collaboration welcome (using git-approach)