



# Getting the nitty-gritty details into big data

*Michael Cysouw*  
Phillips-University Marburg

# How big is big ?

- “Big” does not necessarily mean “googol-size” ( $= 10^{100}$ ), or “google-ngram-size” ( $\approx 10^{10}$ )
- **“Big” often is much simpler, just: “too big to grasp in one human brain”**
- Humanity-scholars often spend decades to be able to grasp an awful amount of data in one human brain!

# The importance of the nitty-gritty stuff

- By using large datasets, many problems can be automatically solved
- Often, 70~80% precision is easily reached, with each further percentage being difficult
- **Yet, many insights in the humanities depend upon unusual, strange, rare cases**
- It is difficult to get beyond the “we-knew-that-already” state when using big data

# Ideal world ...

- **Let the big data do the easy part automatically**
- **So humans can focus on the interesting stuff!**
- **Observation:**
  - ▶ papers on big-data automatic analyses often report results only in the aggregate
  - ▶ Yet, it is normally the concrete decisions in individual cases that are more revealing!

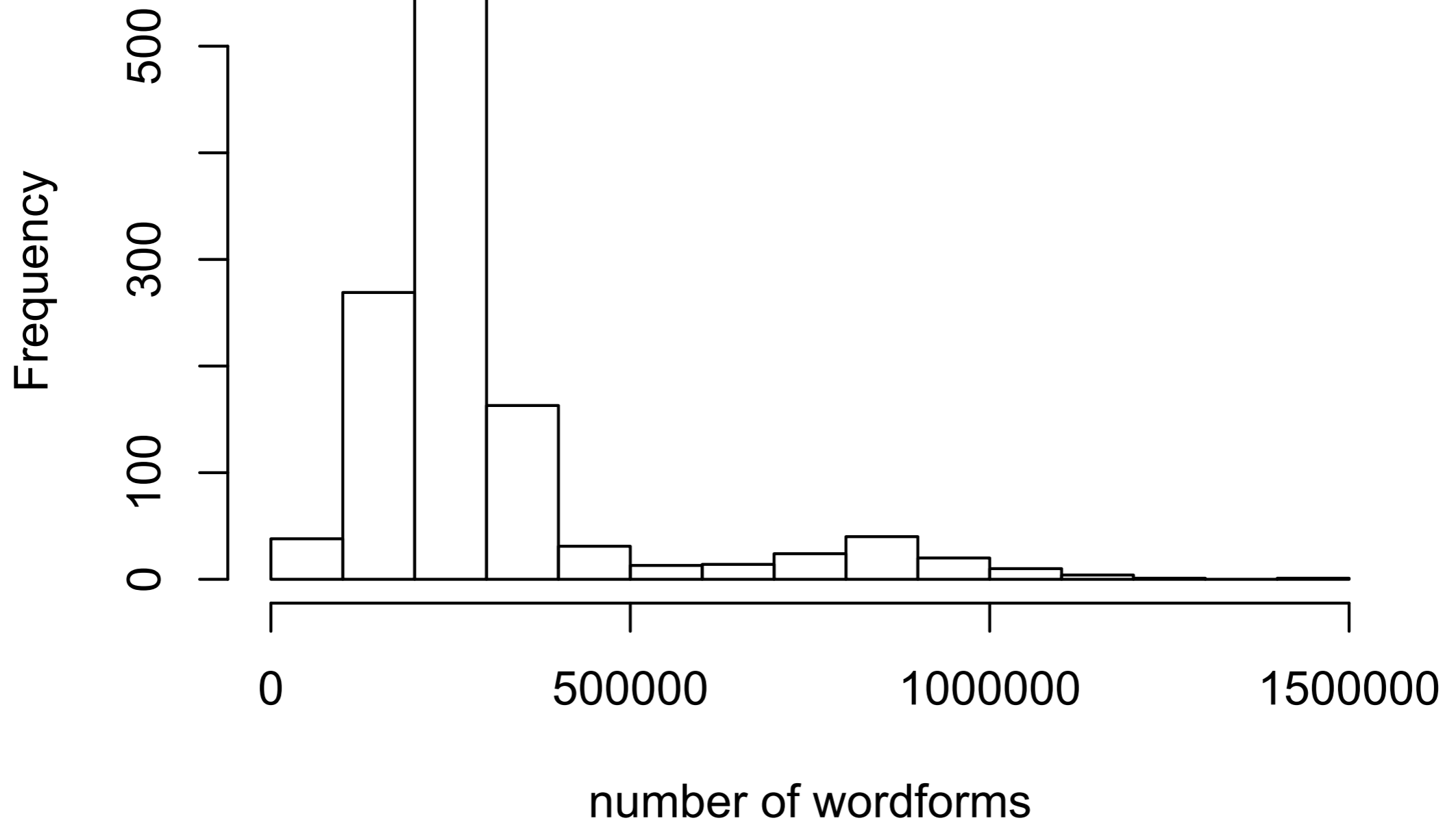




# 1. Bible corpus

paralleltext.info

- Bible translations in unified and cleaned format, tokenised for punctuation
- **Currently 1172 translation in 907 different ISO 639-3 codes (“languages”)**
- On average “only” **300 K** wordforms per translation (including punctuation...)
- Total corpus **350 M wordforms!**



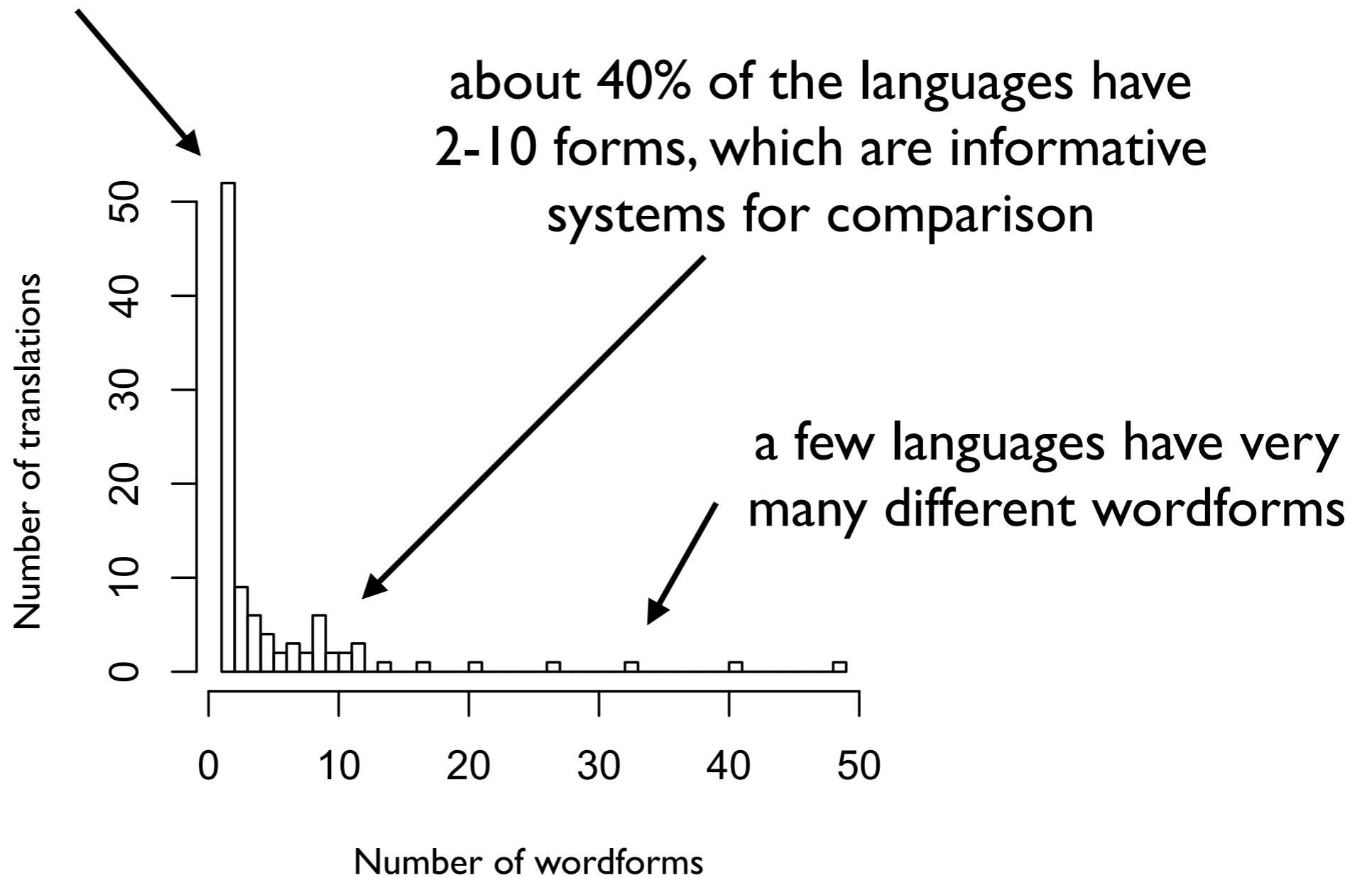
# Word Alignment

- GIZA, berkeley aligner, fast align ...
- One problem: with 1000 translations, there are **500.000 pairs**, which even optimistically calculated will take about a year to align everything
- What really is interesting for comparative linguistics is a **multiple-alignment**, and getting from pairs to multiple is not trivial either
- From a manual perspective, the quality is low

# 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 case!

more than half of the languages  
have just one wordform



# Angaataha

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

- jerusaremthanda
- jerusaremthandaahapt
- jerusaremthandi
- jerusaremthandaahiyai
- jerusaremthandaahit
- jerusaremthandaahapthiyaunt
- jerusaremthandaahiyaisangi
- jerusaremthandaahapthiya
- jerusaremthandaahitraapt
- jerusaremthandaahithit
- jerusaremthandaahithe
- jerusaremthandamt
- jerusaremthanda
- jerusaremthandapt
- jerusaremthandi
- jerusaremthandaahapito
- jerusaremthandaahapaahithit
- jerusaremthandi
- jerusaremthandaahapt
- jerusaremthandaahunt
- jerusaremthandaahapunt
- jerusaremthandaahiya
- jerusaremthandamthint
- jerusaremthandaahapthiyaatihit
- jerusaremthandaahapthiyaate
- jerusaremthandaahiyaunt
- jerusaremosthiyaate

# Amharic

(ISO 639-3 amh, spoken in Ethiopia)

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



# Amharic

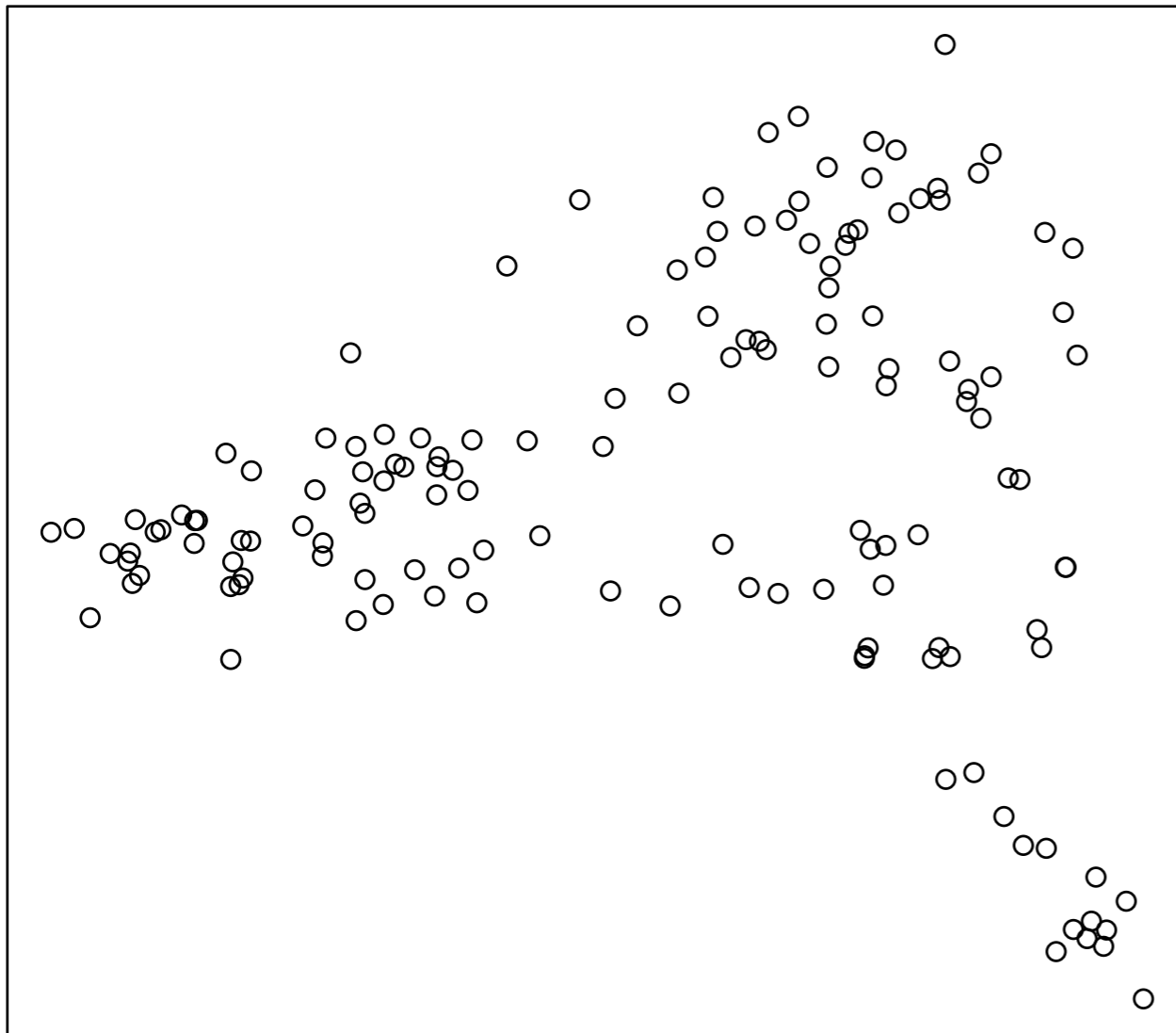
(ISO 639-3 amh, spoken in Ethiopia)

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

98 languages, in total 520 different wordforms

We selected 167 verses including *Jerusalem*

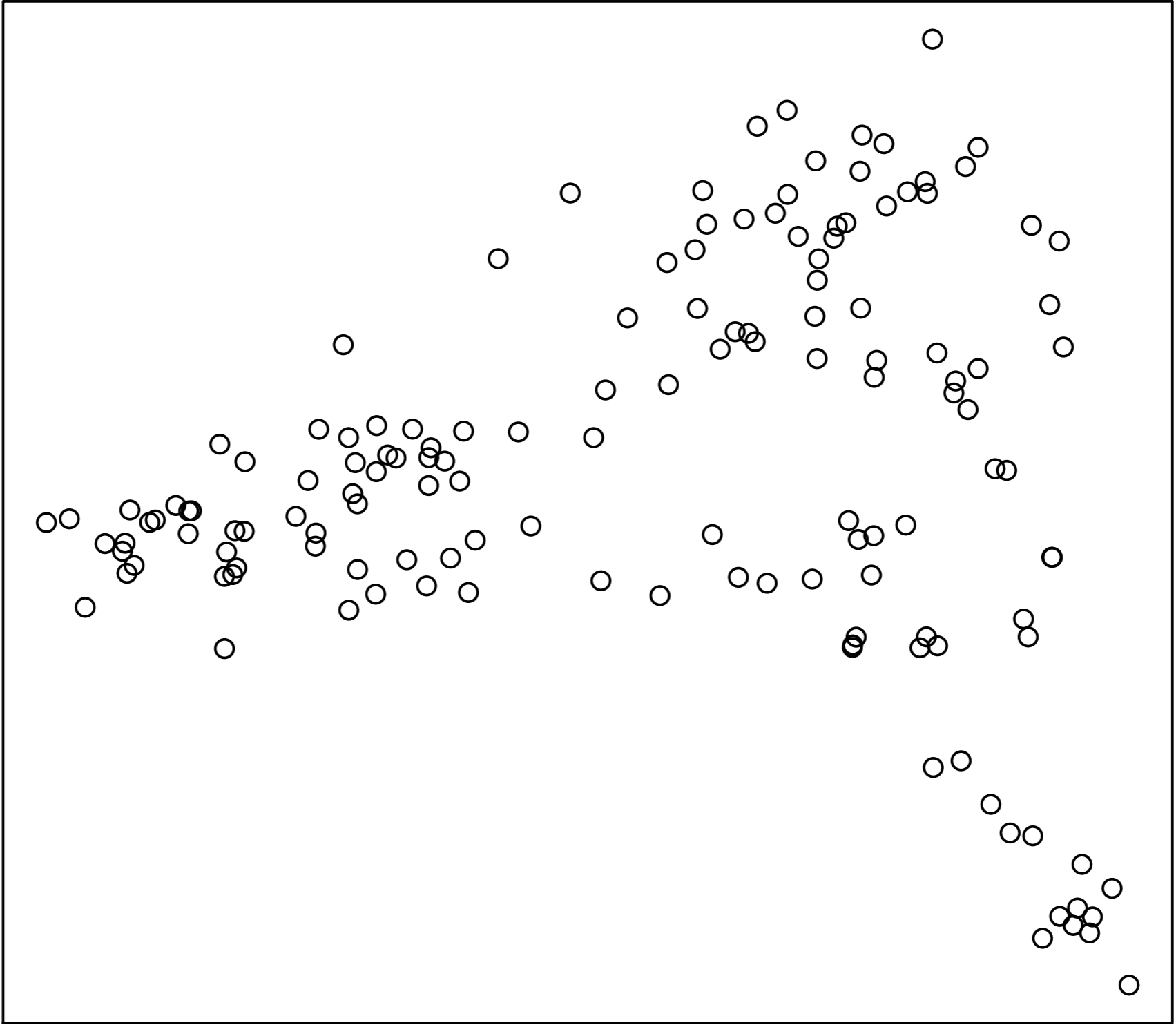
Matrix of size 520 x 167 coding the distribution of wordforms over verses



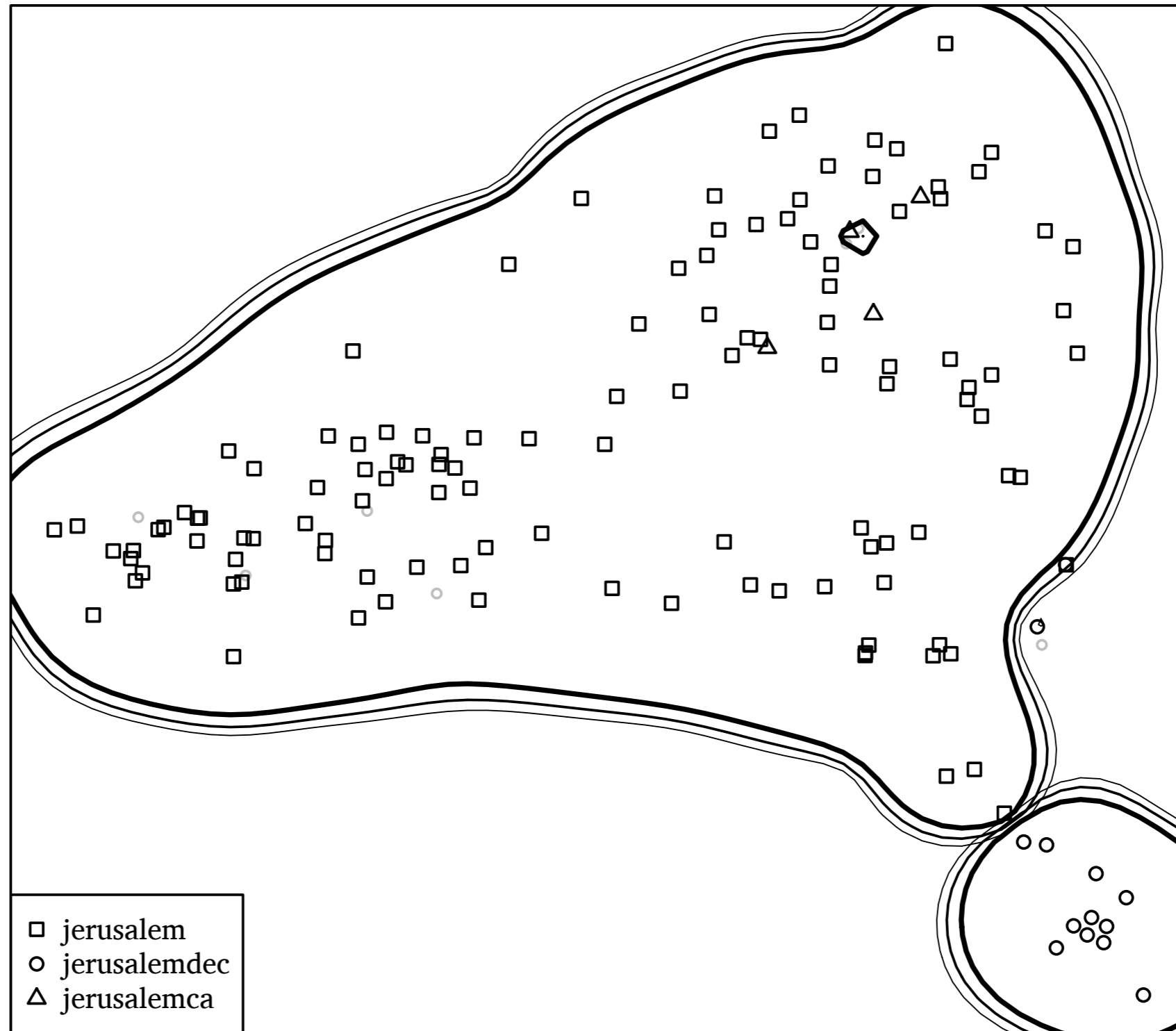
Two verses with *Jerusalem* are similar when they often share the same wordform in language after language

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

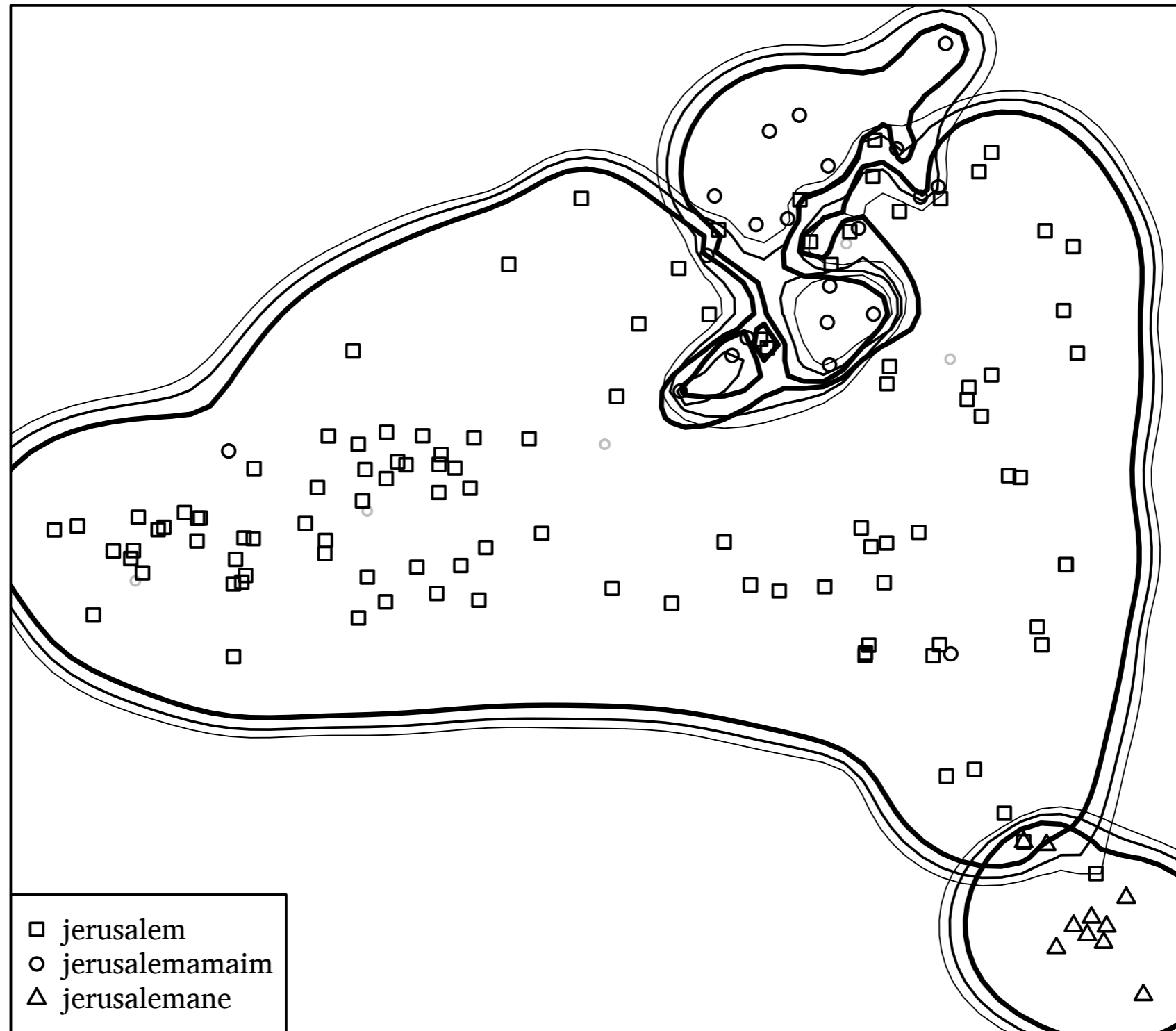


aey



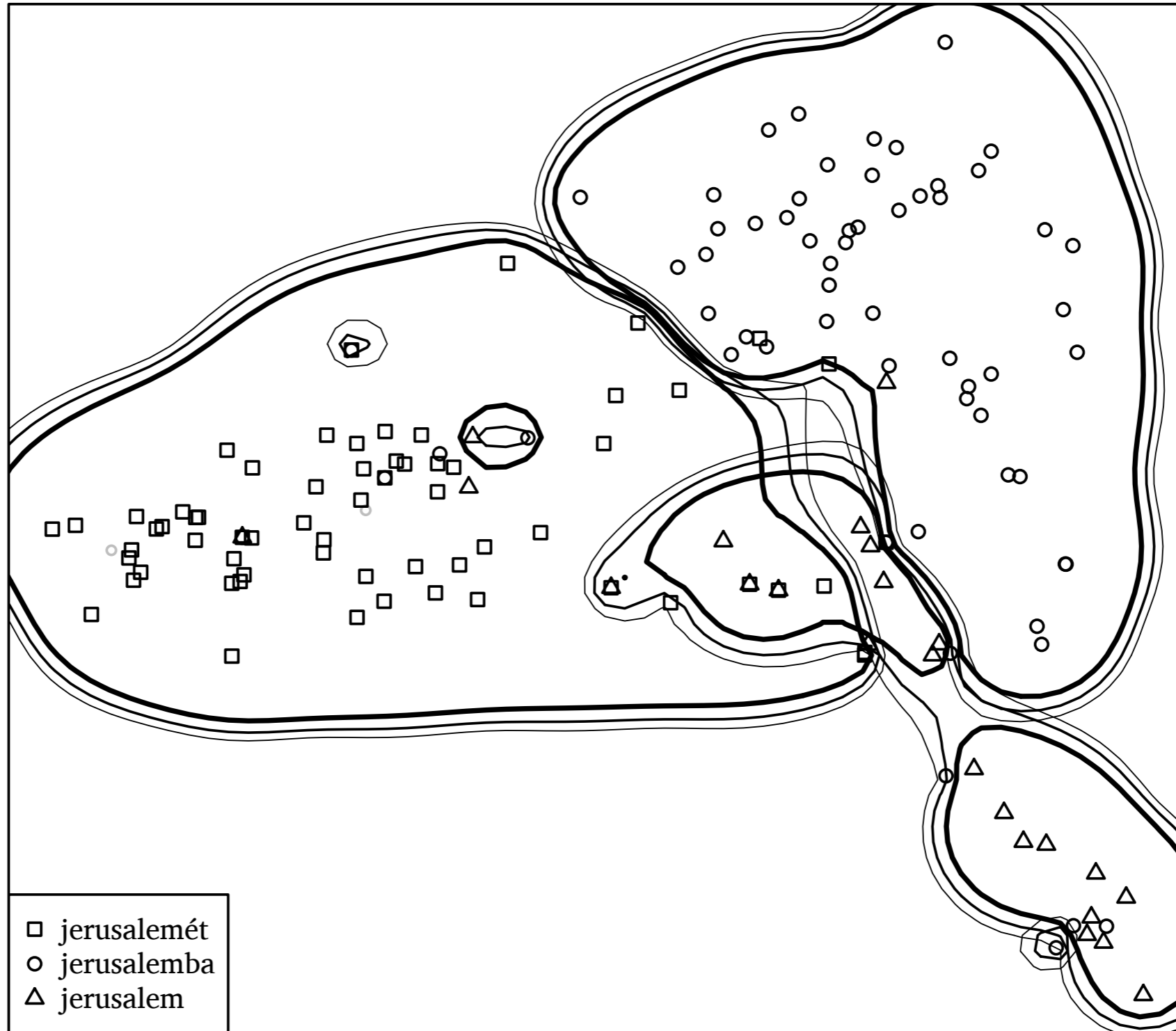
Amele (A language of Papua New Guinea)

# aai



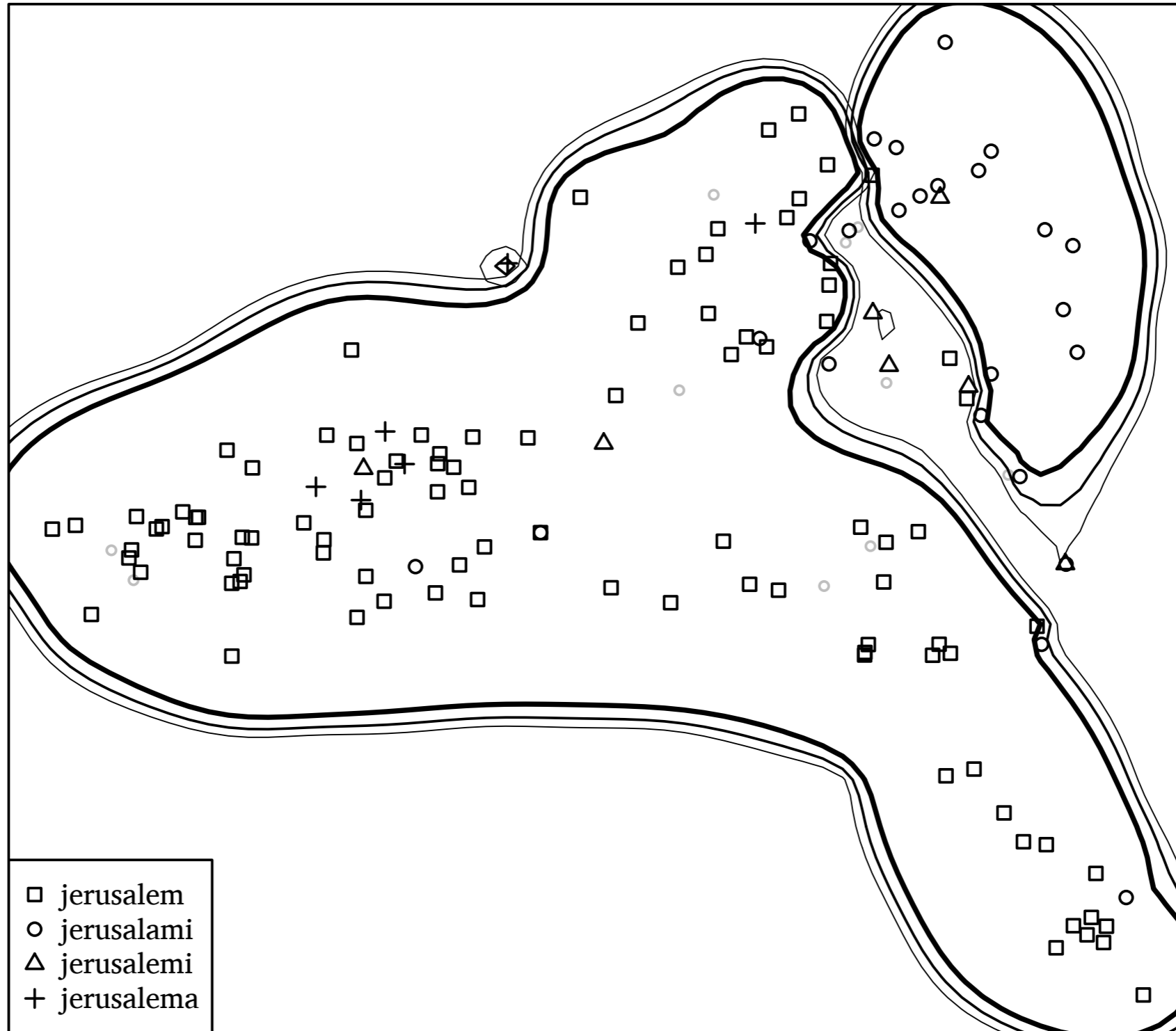
Arifama-Miniafia (a language of Papua New Guinea)

abt



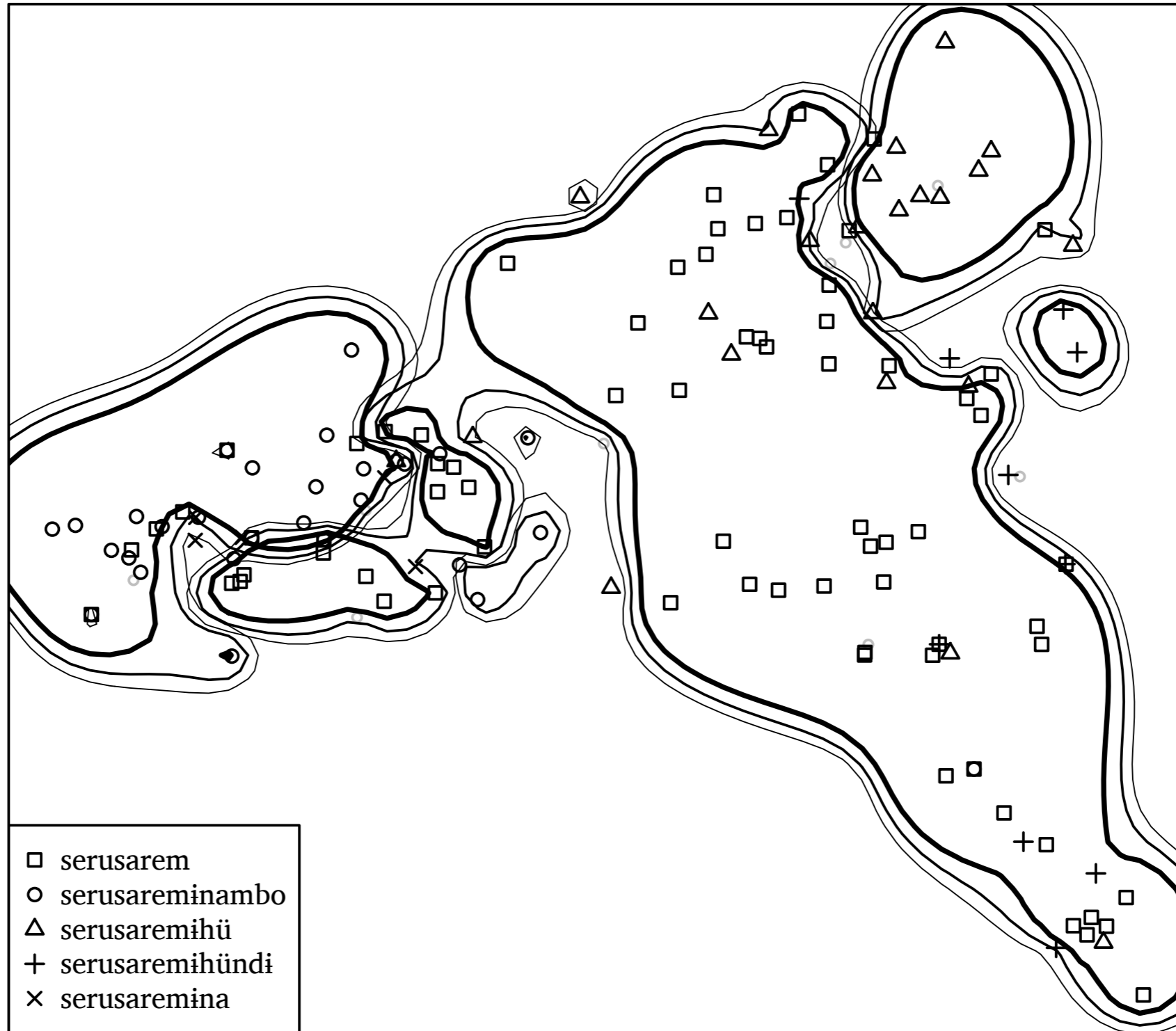
Ambulas (a language of Papua New Guinea)

aoj



Muffian (a language of Papua New Guinea)

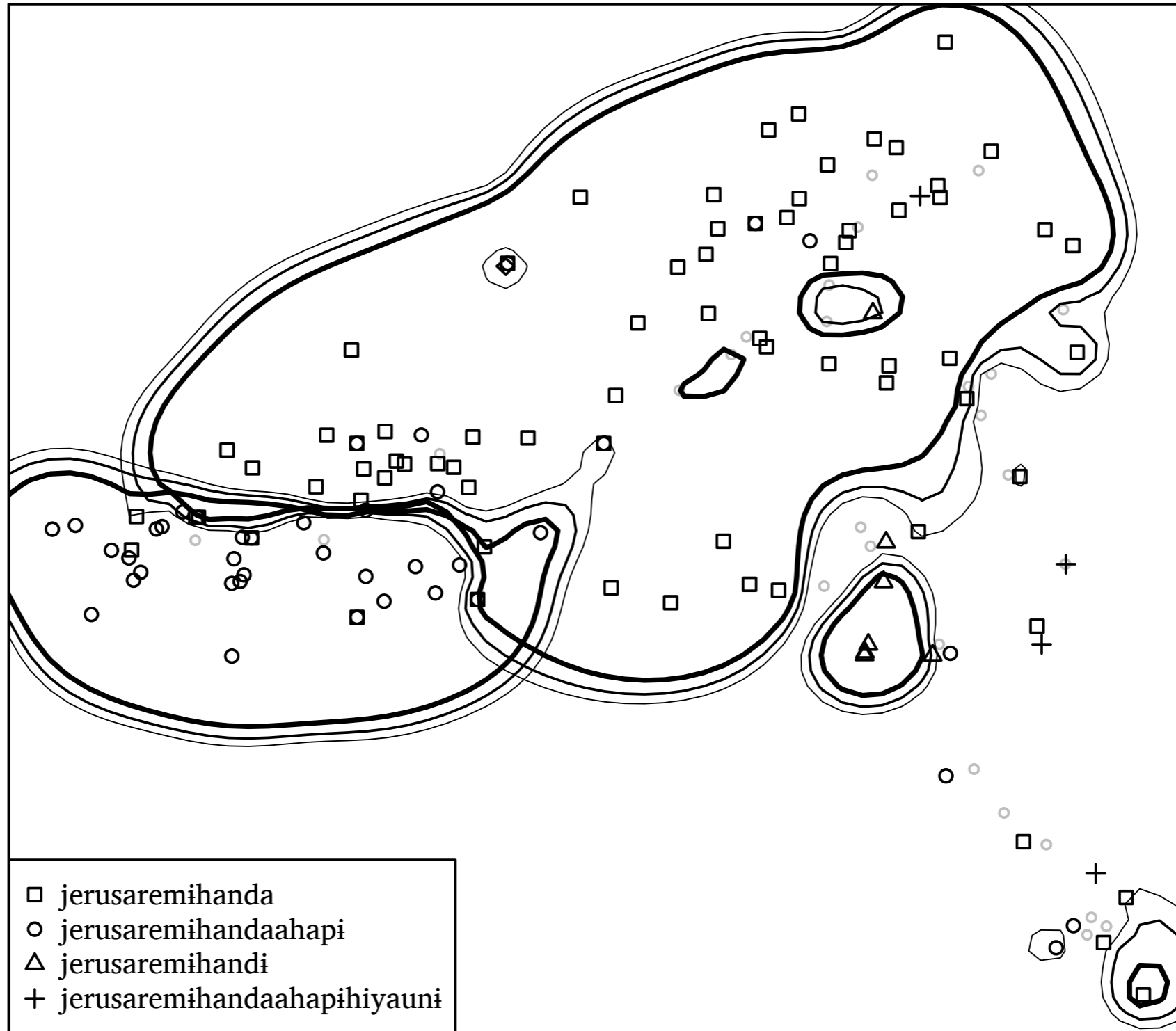
# agg



Angor (a language of Papua New Guinea)

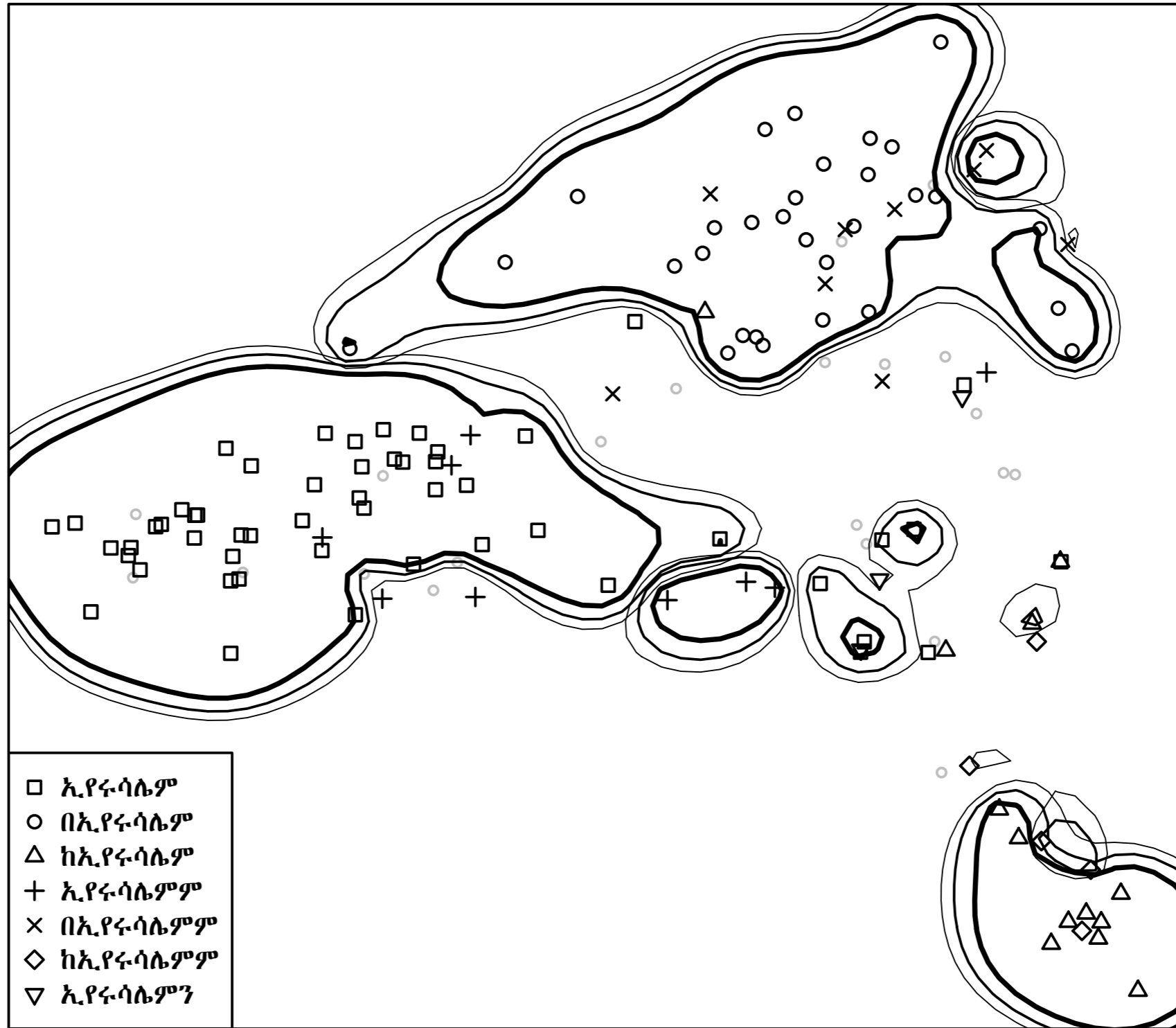


agm



Angaatiha (a language of Papua New Guinea)

amh

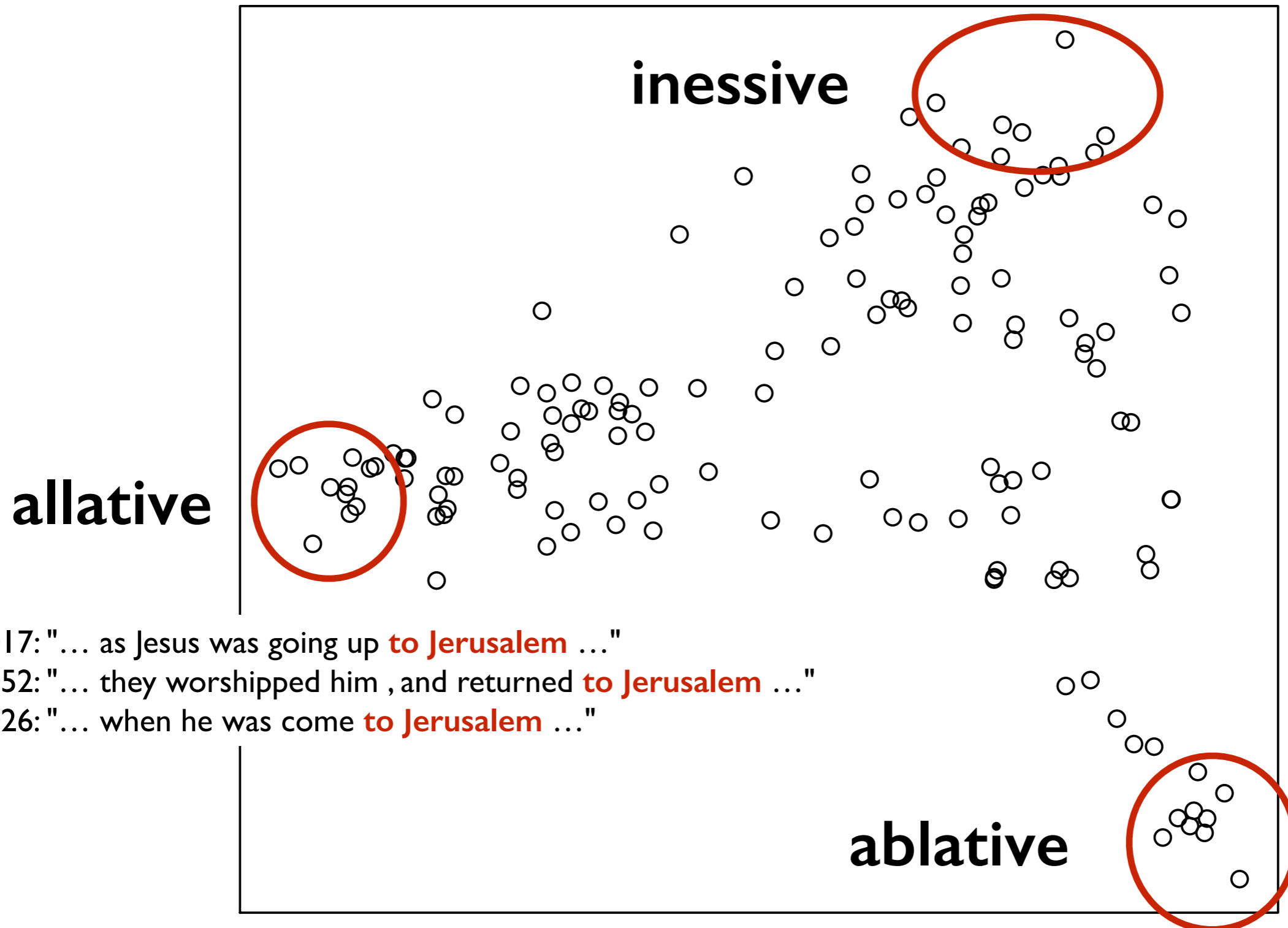


Amharic (a language of Ethiopia)

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 ..."



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 ..."

# 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, but manual decisions are needed

# 2. Dictionaries for historical comparison

- Comparative-historical method
  - ▶ use complete knowledge of languages
- Concept-list method (“Swadesh”)
  - ▶ use small dataset per language (typically ~100 words) and process automatically
- Challenge: use large datasets automatically
  - ▶ <http://quanthistling.info/data/>

# Jivaroan

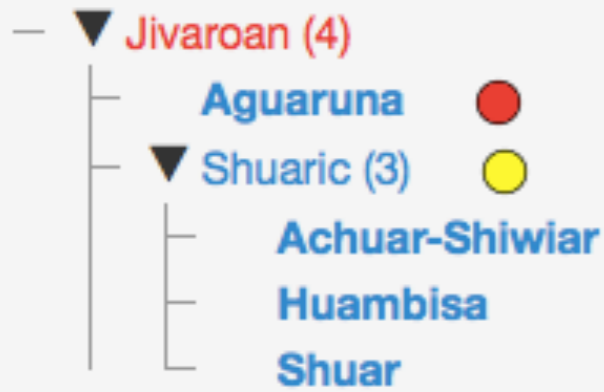
[github.com/cysouw/jivaro](https://github.com/cysouw/jivaro)

- Experiment in collaboratively making an etymological dictionary
- 4 languages, 7 dictionaries, 46 K wordforms

# Family: Jivaroan

## Classification

[open Jivaroan](#) [expand all](#) [collapse all](#)



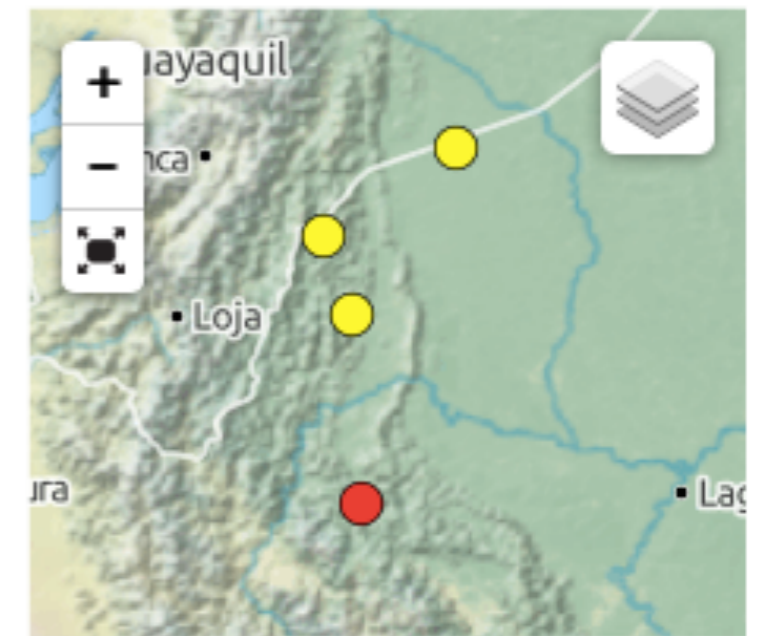
### Family membership references

- [Adelaar, Willem F. H. and Muysken, Pieter C. 2004](#)

### Subclassification references

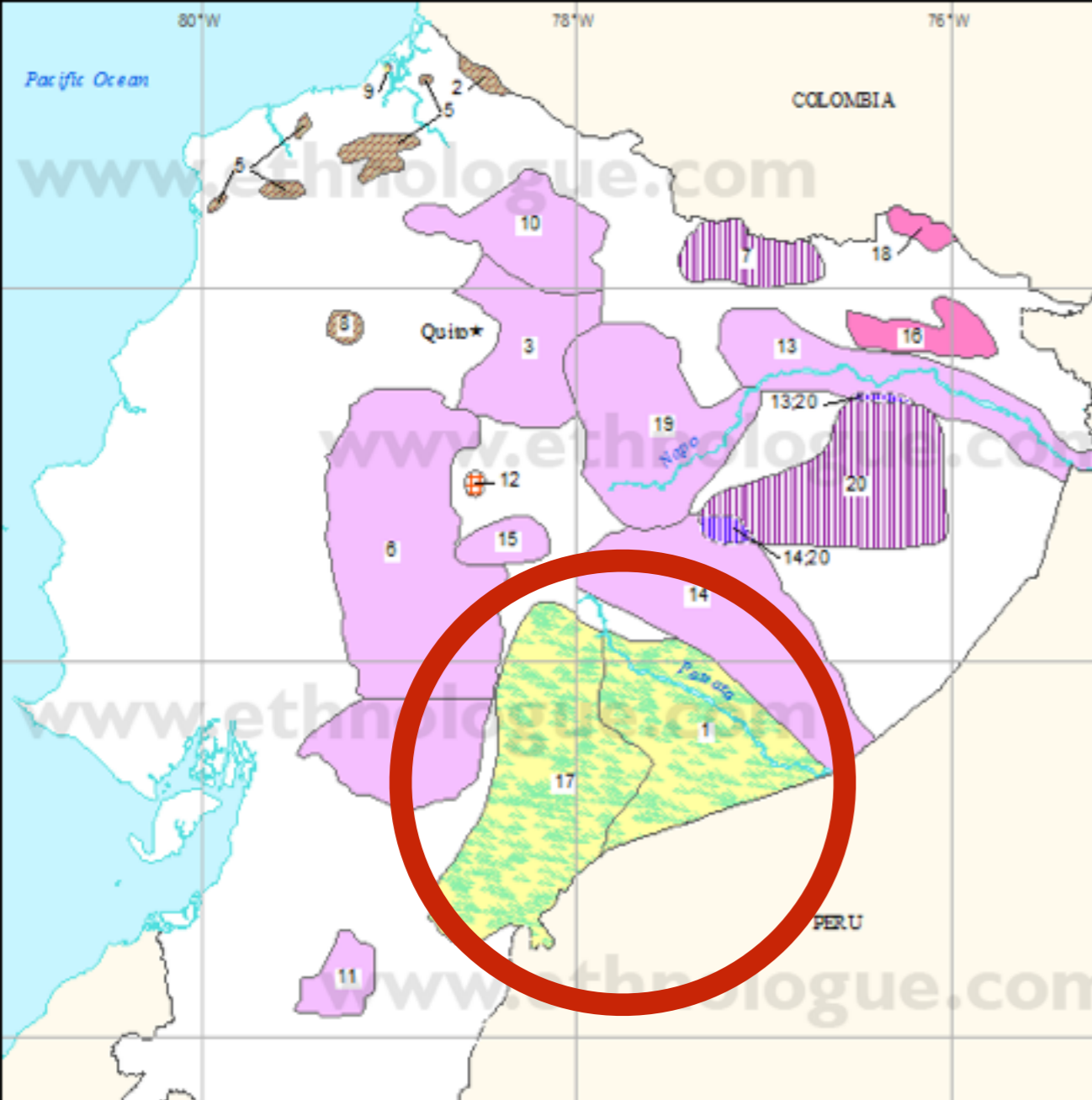
- [Overall, Simon 2007](#)

## Map



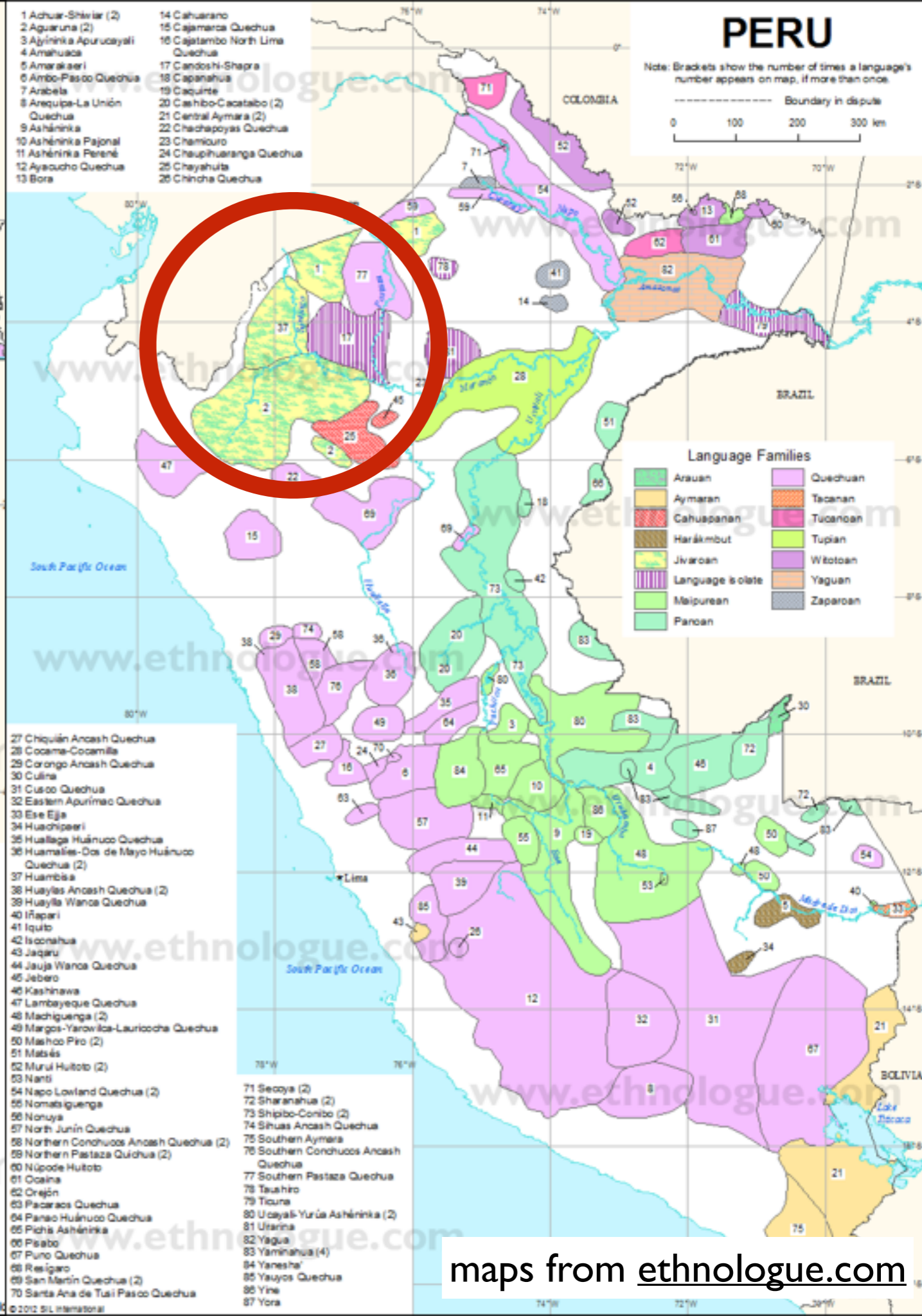
## Links





## ECUADOR

- Achuar-Shiwiar
- Awa-Cuaikuer
- Calderón Highland Quichua
- Cañar Highland Quichua
- Chachi (2)
- Chimborazo Highland Quichua
- Cofán
- Colorado
- Epena
- Imbabura Highland Quichua
- Loja Highland Quichua
- Media Lengua
- Napo Lowland Quichua (2)
- Northern Pastaza Quichua (2)
- Salasaca Highland Quichua
- Shuar
- Siona
- Tena Lowland Quichua
- Wao-rani (3)



- Achuar-Shiwiar (2)
- Aguaruna (2)
- Ajinka Apurucayali
- Amahuaca
- Amarakaeri
- Ambo-Pasto Quechua
- Arabela
- Arequipa-La Unión Quechua
- Asháninka
- Asháninka Pajonal
- Asháninka Perené
- Ayacucho Quechua
- Bora
- Cahuarano
- Cajamarca Quechua
- Cajatambo North Lima Quechua
- Candoshi-Shapra
- Capashua
- Caquinte
- Cashibo-Cacataibo (2)
- Central Aymara (2)
- Chachapoyas Quechua
- Chemuro
- Chaphuaranga Quechua
- Chayahua
- Chincha Quechua
- Chiquián Ancash Quechua
- Cocoma-Cocamilla
- Corongo Ancash Quechua
- Culina
- Cusco Quechua
- Eastern Apurímac Quechua
- Ese Eja
- Huachipero
- Huallaga Huánuco Quechua
- Huamán - Dos de Mayo Huánuco Quechua (2)
- Humbia
- Huaylas Ancash Quechua (2)
- Huaylla Wanca Quechua
- Irapari
- Iquito
- Isonahua
- Jaquiru
- Jauja Wanca Quechua
- Jebero
- Kashinawa
- Lambayeque Quechua
- Machiguenga (2)
- Margos-Yarowilca-Lauricocha Quechua
- Mashco Piro (2)
- Matsigena
- Muril Huilto (2)
- Nanti
- Secoya (2)
- Sharanahua (2)
- Shipibo-Conibo (2)
- Sihuas Ancash Quechua
- Southern Aymara
- Southern Conchucos Ancash Quechua
- Southern Pastaza Quechua
- Taushiro
- Ticuna
- Ucayali-Yurúa Asháninka (2)
- Urarina
- Yagua
- Yaminahua (4)
- Yanesha'
- Yauyos Quechua
- Yine
- Yora

maps from [ethnologue.com](http://ethnologue.com)



# Using the github-approach

- Keep everything in one big file
  - ▶ [github.com/cysouw/jivaro](https://github.com/cysouw/jivaro)
- Automatic grouping in ‘cognate sets’
  - ▶ using quick-and-dirty sparse matrix algebra from [github.com/cysouw/qlcMatrix](https://github.com/cysouw/qlcMatrix)
- Automatic alignment
  - ▶ using LingPy from [github.com/lingpy](https://github.com/lingpy)
- Manual corrections using Alignment Editor
  - ▶ [github.com/digitallinguists/msa-editor](https://github.com/digitallinguists/msa-editor)



&lt;data.tsv&gt;

Showing 1 - 16 of 16 entries

START

ID	LANGUAGE	SOURCE	ETYMONID	ALIGNMENT	TRANSLATION
53	Achuar-Shiwiar	fastmowitz2008	21	a ch í - m k a - - t i n -	agarrarse
62	Achuar-Shiwiar	fastmowitz2008	21	a ch i - - - u - - - - -	agarrador; uno que agarra
8460	Huambisa	jakway2008	21	a ch i - - - - - a m u	agarrar; capturar; sacar chonta
8464	Huambisa	jakway2008	21	a ch i k m - a - - - - u	agarrado
8468	Huambisa	jakway2008	21	a ch i - n - a - - t - - -	agarrar; capturar; sacar chonta
14812	Aguaruna	larson1966	21	a ch í - t - - - - - - -	agarrar; prender; clavar
17642	Aguaruna	mori1981	21	a ch í - t - - - - - - -	agarrar
29172	Aguaruna	wipiodeicat1996	21	a ch i - m - á - - t - - -	agarrar; prender
35619	Aguaruna	yagkug1998	21	a ch í - - - - - a m u	clavar; hincar un objeto; sacar ...
35637	Aguaruna	yagkug1998	21	a ch i - m - á - - t - - -	agarrarse en algún objeto
35638	Aguaruna	yagkug1998	21	a ch í - m - a - u - - - -	agarrarse en algún objeto
35639	Aguaruna	yagkug1998	21	a ch i - m k á - - t a s a	agarrarse en algún objeto
35644	Aguaruna	yagkug1998	21	a ch i - n - í - - - a m u	agarrarse mutuamente; darse l...
35645	Aguaruna	yagkug1998	21	a ch i - n - í k - t a s a	agarrarse mutuamente; darse l...
35646	Aguaruna	yagkug1998	21	a ch i - n - í - - t - - -	agarrarse mutuamente; darse l...
35649	Aguaruna	yagkug1998	21	a ch í - t - - - - - - -	clavar; hincar un objeto; sacar ...



<data.tsv> Showing 1 - 16 of 16 entries START

ID	LANGUAGE	SOURCE	ETYM
53	Achuar-Shiwiar	fastmowitz2008	21
62	Achuar-Shiwiar	fastmowitz2008	21
8460	Huambisa	jakway2008	21
8464	Huambisa	jakway2008	21
8468	Huambisa	jakway2008	21
14812	Aguaruna	larson1966	21
17642	Aguaruna	mori1981	21
29172	Aguaruna	wipiodeicat1996	21
35619	Aguaruna	yagkug1998	21
35637	Aguaruna	yagkug1998	21
35638	Aguaruna	yagkug1998	21
35639	Aguaruna	yagkug1998	21
35644	Aguaruna	yagkug1998	21
35645	Aguaruna	yagkug1998	21
35646	Aguaruna	yagkug1998	21
35649	Aguaruna	yagkug1998	21

ETYMONID "21" links the following 16 entries:

Achuar-Shiwiar	a	ch	í	-	m	k	a	-	-	t	i	n	-
Achuar-Shiwiar	a	ch	í	-	-	-	u	-	-	-	-	-	-
Huambisa	a	ch	i	-	-	-	-	-	-	a	m	u	-
Huambisa	a	ch	i	k	m	-	a	-	-	-	-	u	-
Huambisa	a	ch	i	-	n	-	a	-	-	t	-	-	-
Aguaruna	a	ch	í	-	t	-	-	-	-	-	-	-	-
Aguaruna	a	ch	í	-	t	-	-	-	-	-	-	-	-
Aguaruna	a	ch	i	-	m	-	á	-	-	t	-	-	-
Aguaruna	a	ch	í	-	-	-	-	-	-	a	m	u	-
Aguaruna	a	ch	i	-	m	-	á	-	-	t	-	-	-
Aguaruna	a	ch	í	-	m	-	á	u	-	-	-	-	-
Aguaruna	a	ch	i	-	m	k	á	-	-	t	a	s	a
Aguaruna	a	ch	i	-	n	-	í	-	-	a	m	u	-
Aguaruna	a	ch	i	-	n	-	í	k	-	t	a	s	a
Aguaruna	a	ch	i	-	n	-	í	-	-	t	-	-	-
Aguaruna	a	ch	í	-	t	-	-	-	-	-	-	-	-

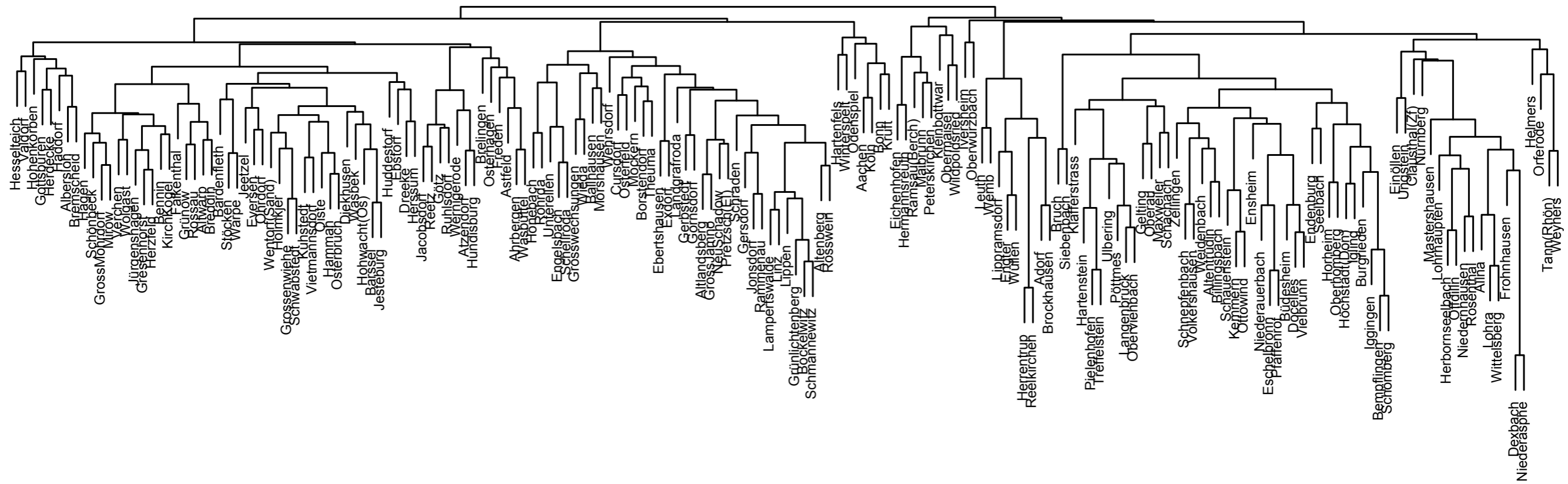
IGNORE

EDIT SUBMIT EXPORT CLOSE

# 3. German Dialect Data

[github.com/cysouw/pad](https://github.com/cysouw/pad)

- 31 K phonetically transcribed words, distributed over 186 meanings and 183 villages
- How to make sense of the details?
- Agglomerative methods are trivial
  - ▶ e.g. using sparse matrix algebra from [github.com/cysouw/qlcMatrix](https://github.com/cysouw/qlcMatrix)



“we knew that already ...”

“... except that you go wrong here and there!”

# Sound Alignment

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

A	FF	E
a:	ph	-
ɑ:	b <sup>h</sup>	ə
o→ɔ	ph	ə
ɑ:	p <sup>h</sup>	ə
ɑ	ϕh	-
ʌ	f	ɛ
a	f	-
ɑ'	f	ə'
o:	ph	-
ɒ':	p <sup>h</sup>	ə
a	f	ɛ
ʌ'	f	ə
ɔ:	p̄ϕ	-
ɒ:	p <sup>h</sup>	ə
a	f:	-
ɔ	p <sup>h</sup>	-
a	f	-
ʌ	v	ə
ɑ:	p'	-
ʏ	f:	-
ɒ:	ph	-
ɑ	f̄β	ə
ɒ':	ph	ə
...	...	...

● **Workflow:**

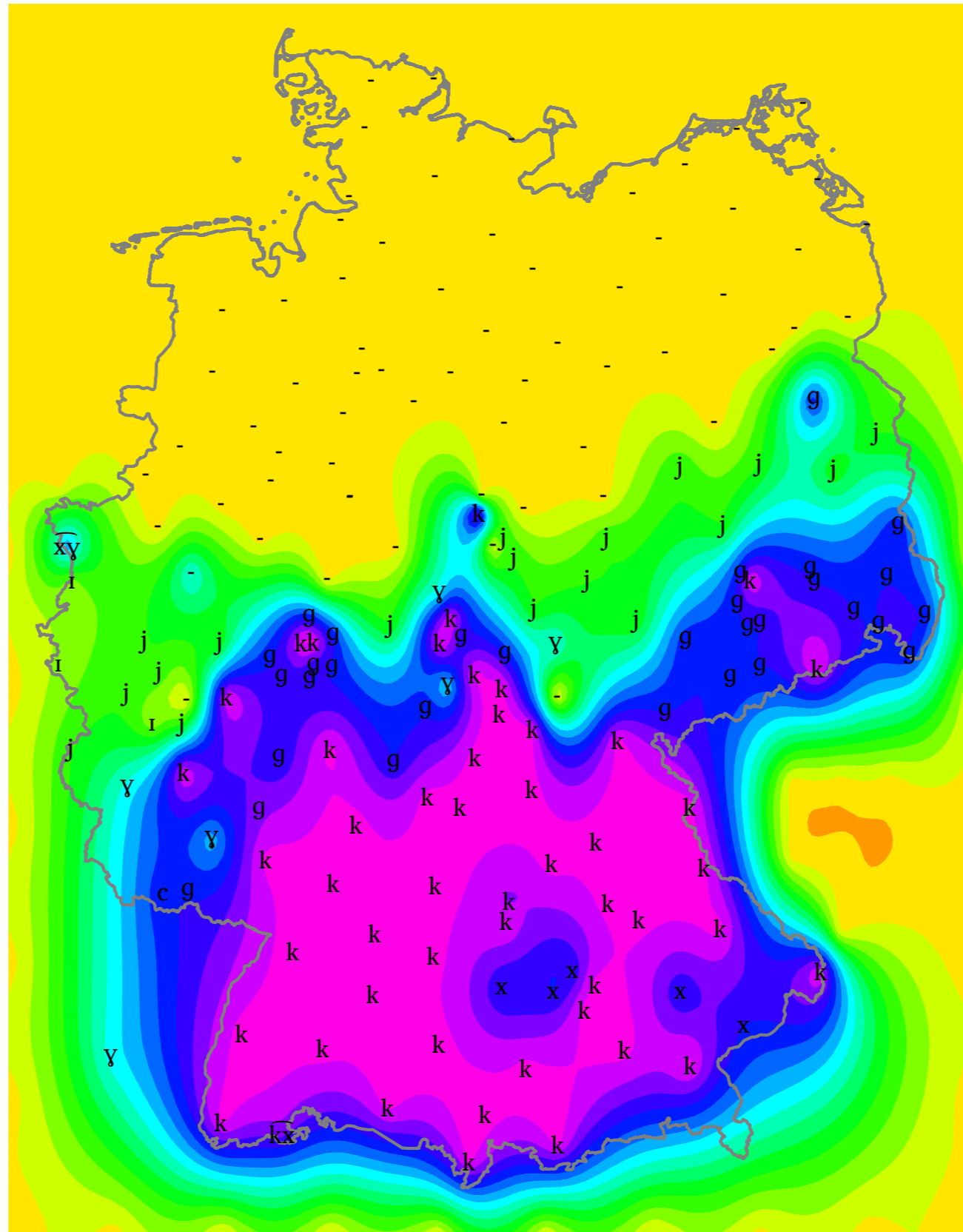
- ▶ Tokenisation of segments using **qlcTokenize** ([github.com/cysouw/qlcTokenize](https://github.com/cysouw/qlcTokenize))
- ▶ Automatic alignment using **LingPy** ([github.com/lingpy](https://github.com/lingpy))
- ▶ Manual correction using **Alignment Editor** ([github.com/digitallinguists/msa-editor](https://github.com/digitallinguists/msa-editor))
- ▶ Separation of cognates (e.g. *Samstag* vs. *Sonnabend*)
- ▶ Annotation of columns (e.g. many-to-one alignments, metathesis)

# Investigation of Detail

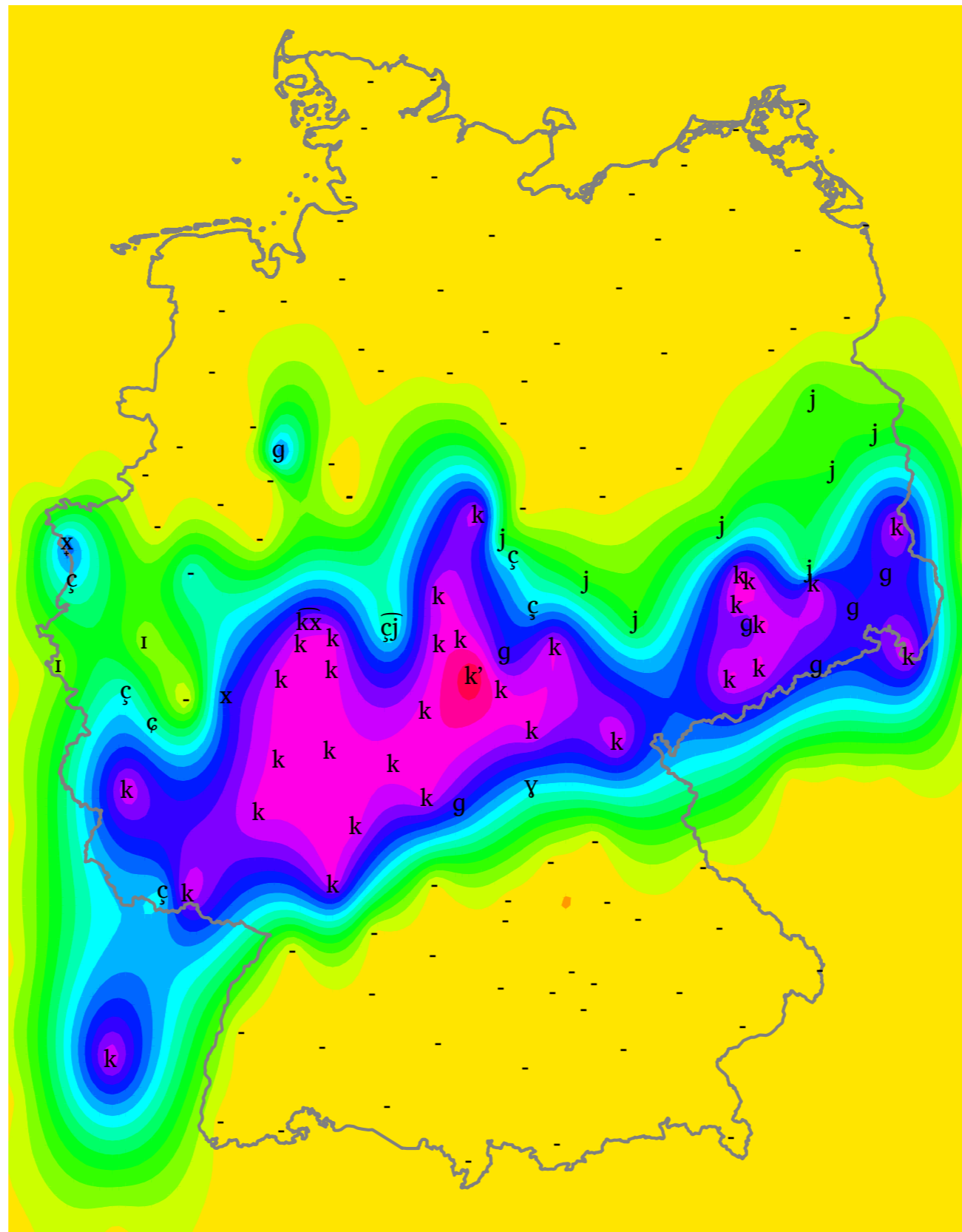
- **Consonant** [g] in German Perfect-prefix <ge>
- Ordered to ‘strength’, visualised as ‘height’
- (NULL) ʏ ɪ j̣ j̣ j̣ ʔ ʧ̣ʝ̣ ẓ ẓ ʒ̣ ʧ̣ ʏ̣ ʧ̣ ʧ̣ʝ̣ ʏ̣ ʝ̣ ʧ̣ ʏ̣  
 ʧ̣̣̣ ɡ̣̣̣ ʃ̣̣̣ ɡ̣̣̣ ʃ̣̣̣ c̣̣̣ ɡ̣̣̣ ḳ̣̣ʧ̣̣̣ ḳ̣̣j̣̣̣ x̣̣̣ c̣̣̣ ʃ̣̣̣ k<sup>h</sup> ḳ̣̣x̣̣̣ ḳ̣̣ ʀ̣̣̣  
 χ̣ ḳ ḳ̣ q̣ ḳ'



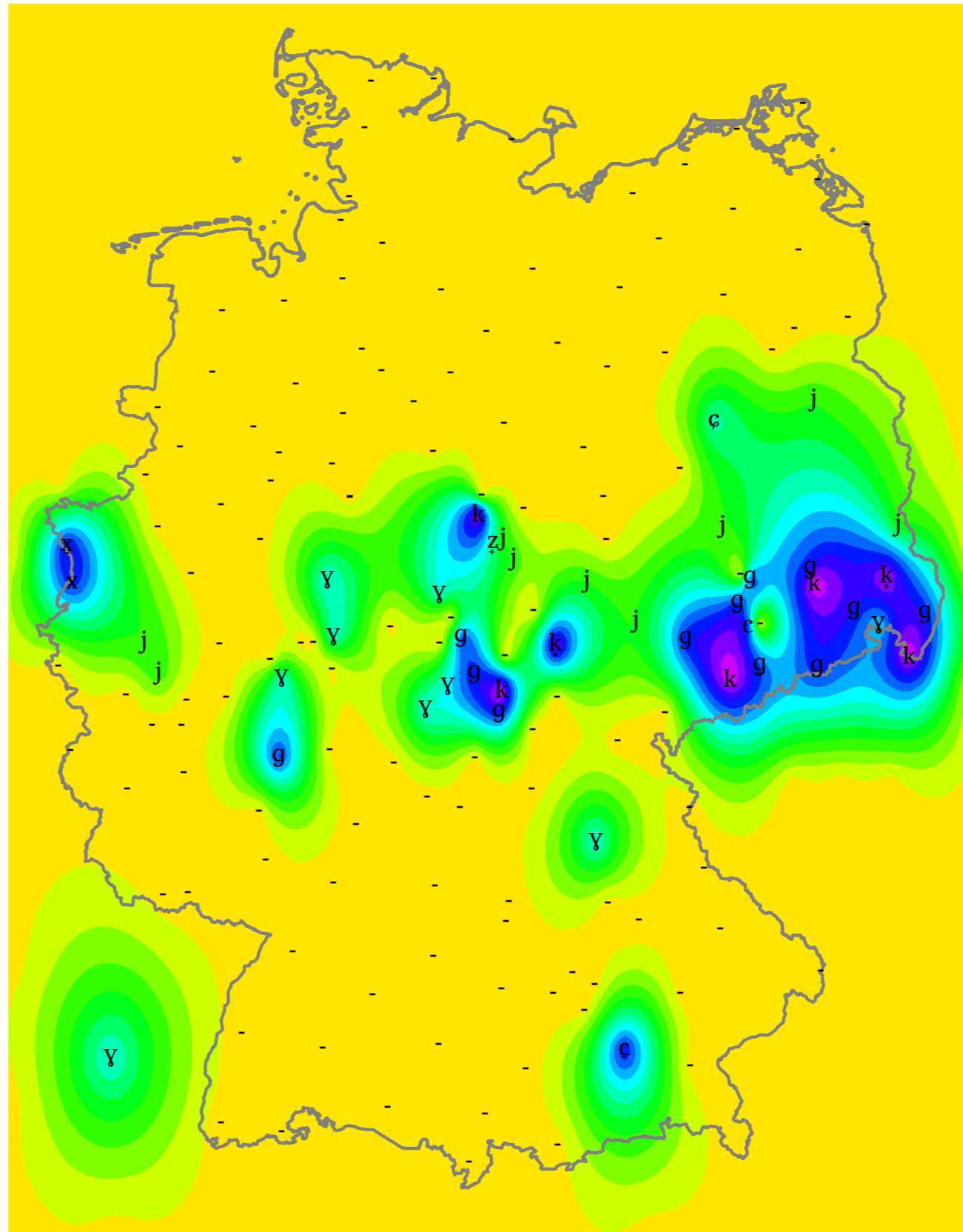
# gefahren



# gebrannt



# gekommen



# Big Data and Small Stories

- Big data are a curse and blessing
- Big (agglomerative and approximative) stories and not a step forward
- Large numbers of cases allow for much detail, but the details have to be identified first
- The detail and minutiae that can be told with big data are the real improvement!