

Using WALS

Prospects of quantitative approaches for linguistic typology

Michael Cysouw



MAX-PLANCK-GESELLSCHAFT

Max Planck Institute
for Evolutionary Anthropology



WALS – World Atlas of Language Structures

WALS



The World Atlas of Language Structures

edited by M. Haspelmath, M. S. Dryer, D. Gil, B. Comrie

The Interactive Reference Tool

developed by Hans-Jörg Bibiko

OXFORD
UNIVERSITY PRESS

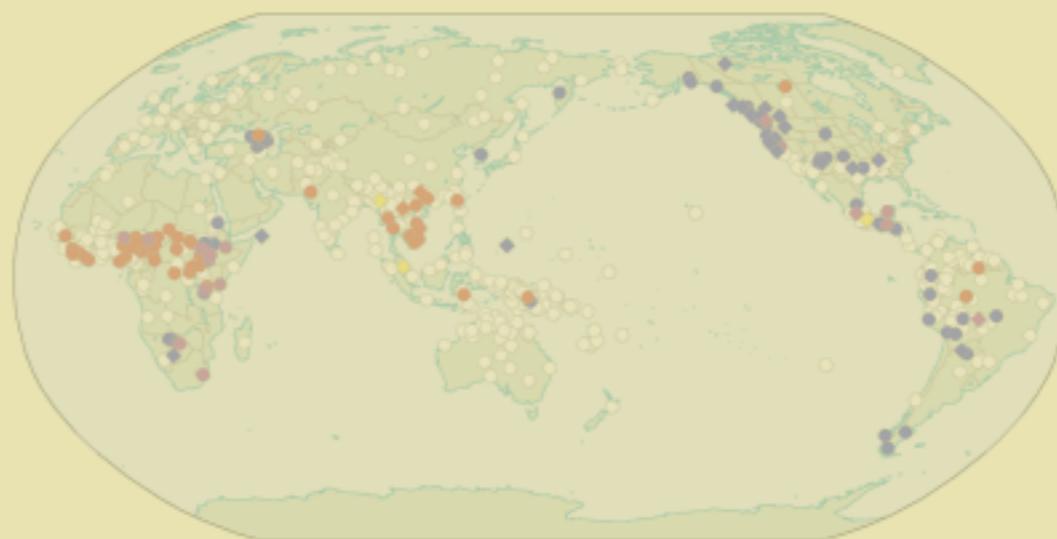
FEATURE VIEWER

LANGUAGE VIEWER

COMPOSER

ABOUT THIS PROGRAM

GUIDED TOUR



Wow !

WALS is just great

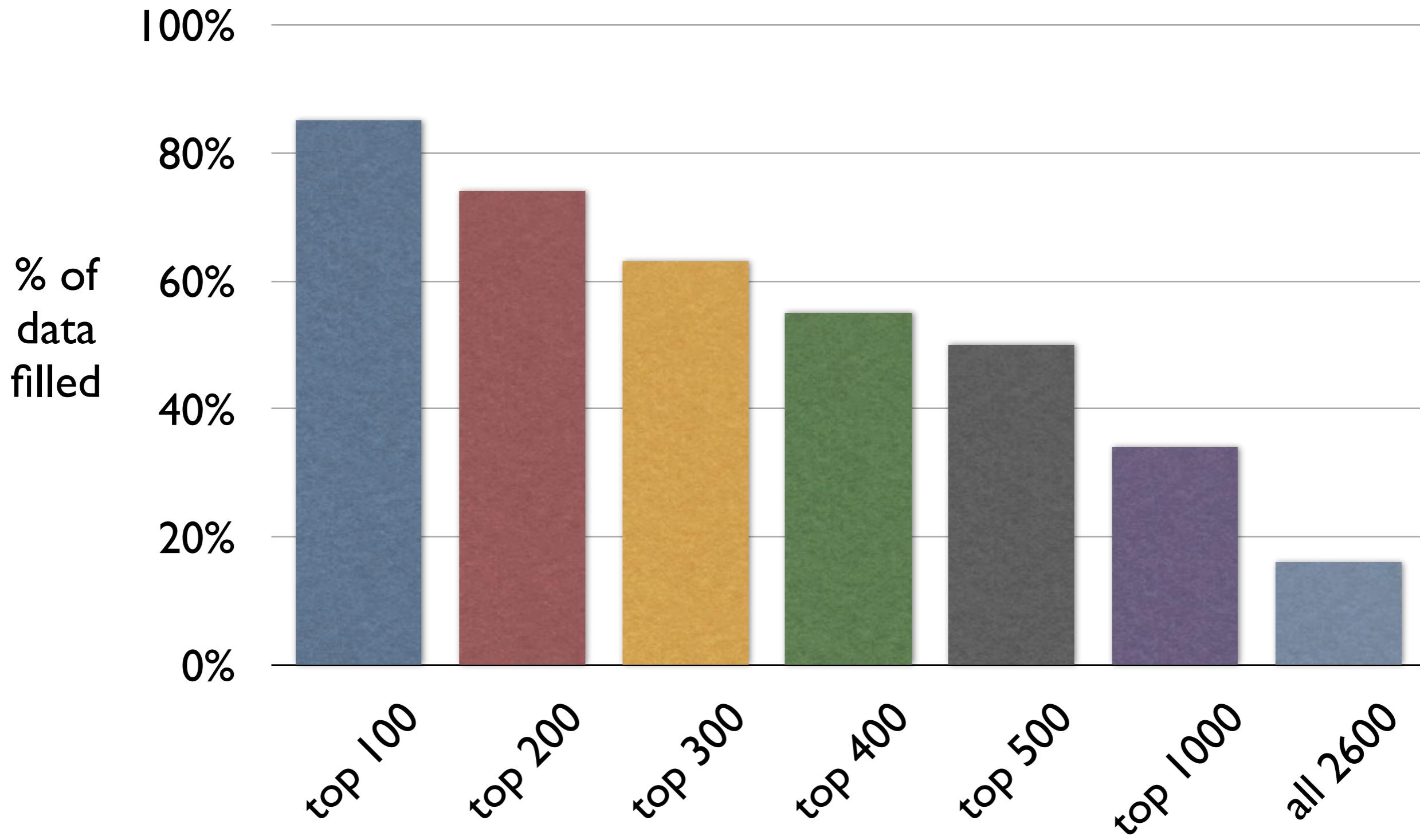
- Linguistic data is very expensive ...
(data estimated 6,000,000 EUR)
- but now we have so much data
- from different linguistic domains
- and so well organized !

Problems

So much data ...

- 2,600 languages
- 140 characteristics
- almost 60,000 datapoints
- wait a minute: $60,000/2,600*140 = 0.165$
- the datatable is only 16.5 % filled !

Choosing the best languages



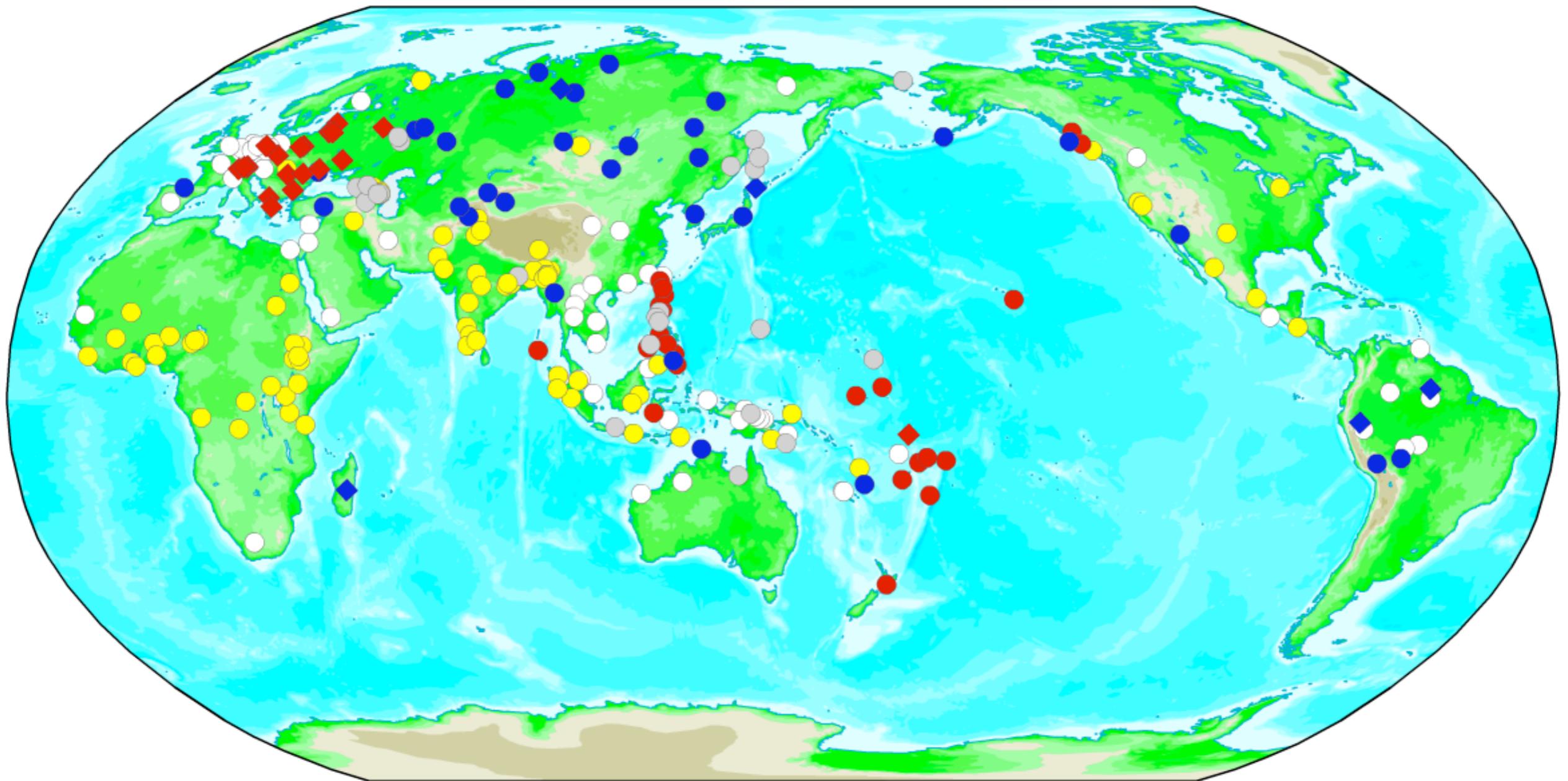
Reliability

- Latvian was checked (by B. Wälchli)
- 109 coding point in WALS
- 2 ‘technical’ errors (= 1.8 %)
- 5 ‘interpretative’ errors (= 4.6 %)

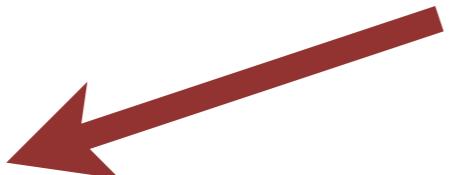
Coding Problems

- **Types consisting of dissimilar languages**
 - Independent features combined in one map
 - Definitional dependencies between maps
 - No relative similarities available

Map 54: Distributive numerals



- 1. No distributive numerals [62]
- 2. Marked by reduplication [84]
- 3. Marked by prefix [23]
- 4. Marked by suffix [32]
- ◆ 5. Marked by preceding word [21]
- ◆ 6. Marked by following word [5]
- 7. Marked by mixed or other strategies [23]



Solution:
recode all as different types

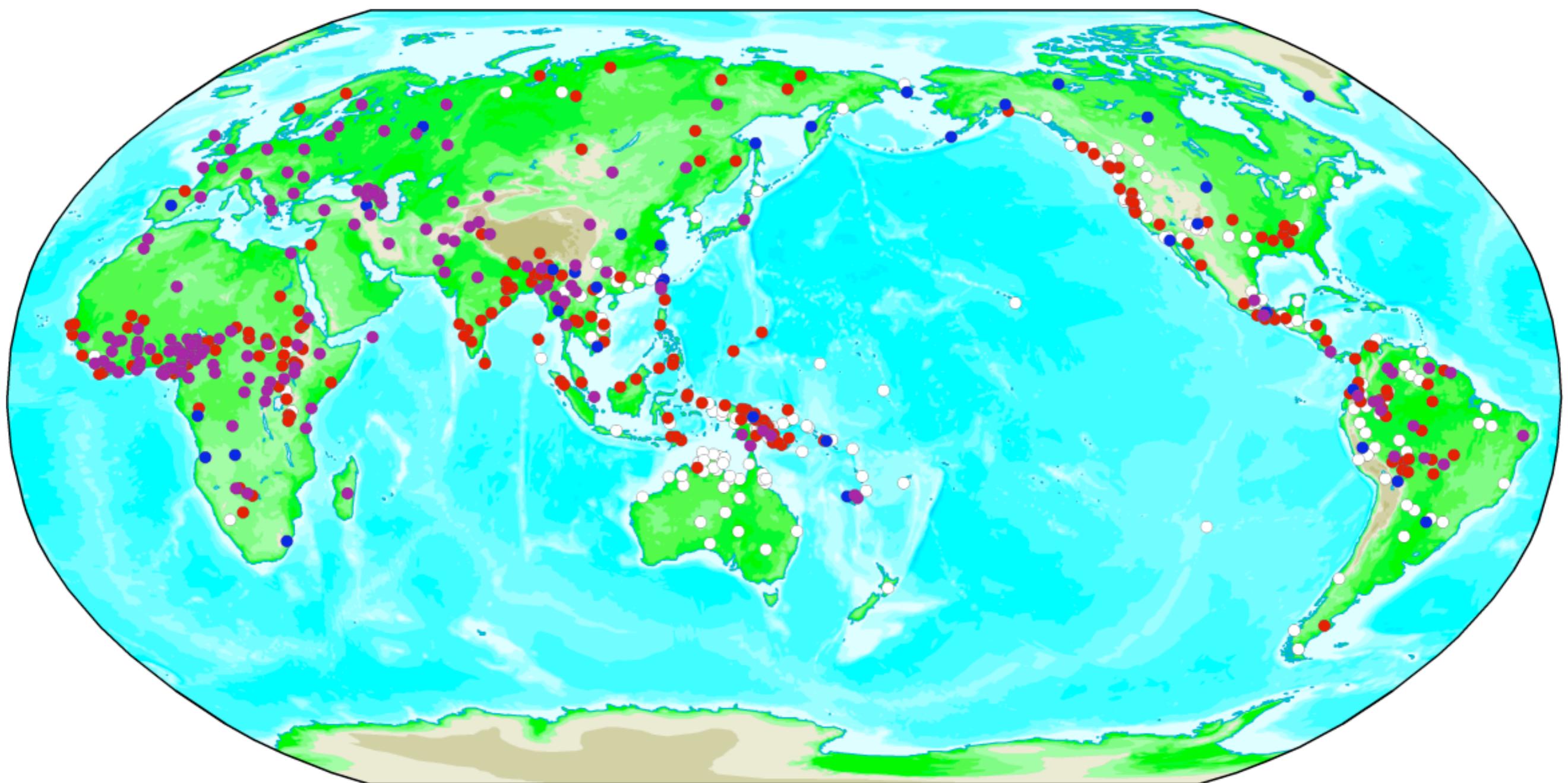
Coding Problems

- **Types consisting of dissimilar languages**
 - Independent features combined in one map
 - Definitional dependencies between maps
 - No relative similarities available

Coding Problems

- Types consisting of dissimilar languages
- **Independent features combined in one map**
- Definitional dependencies between maps
- No relative similarities available

Map 4: Voicing in Plosives and Fricatives



- 1. No voicing contrast [181]
- 2. In plosives alone [189]
- 3. In fricatives alone [38]
- 4. In both plosives and fricatives [158]

Solution:
split, and disregard the original

Coding Problems

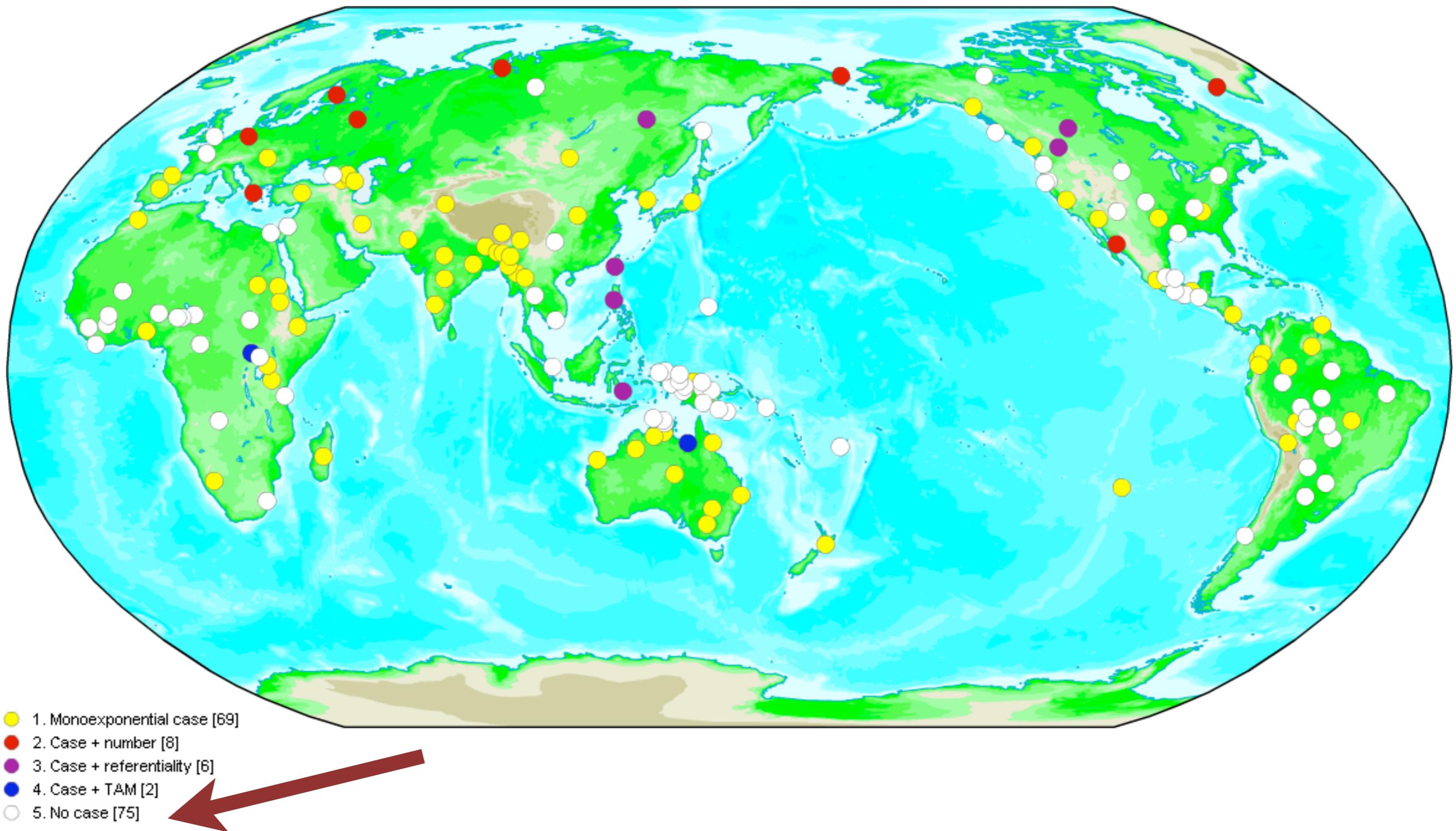
- Types consisting of dissimilar languages
- **Independent features combined in one map**
- Definitional dependencies between maps
- No relative similarities available

Coding Problems

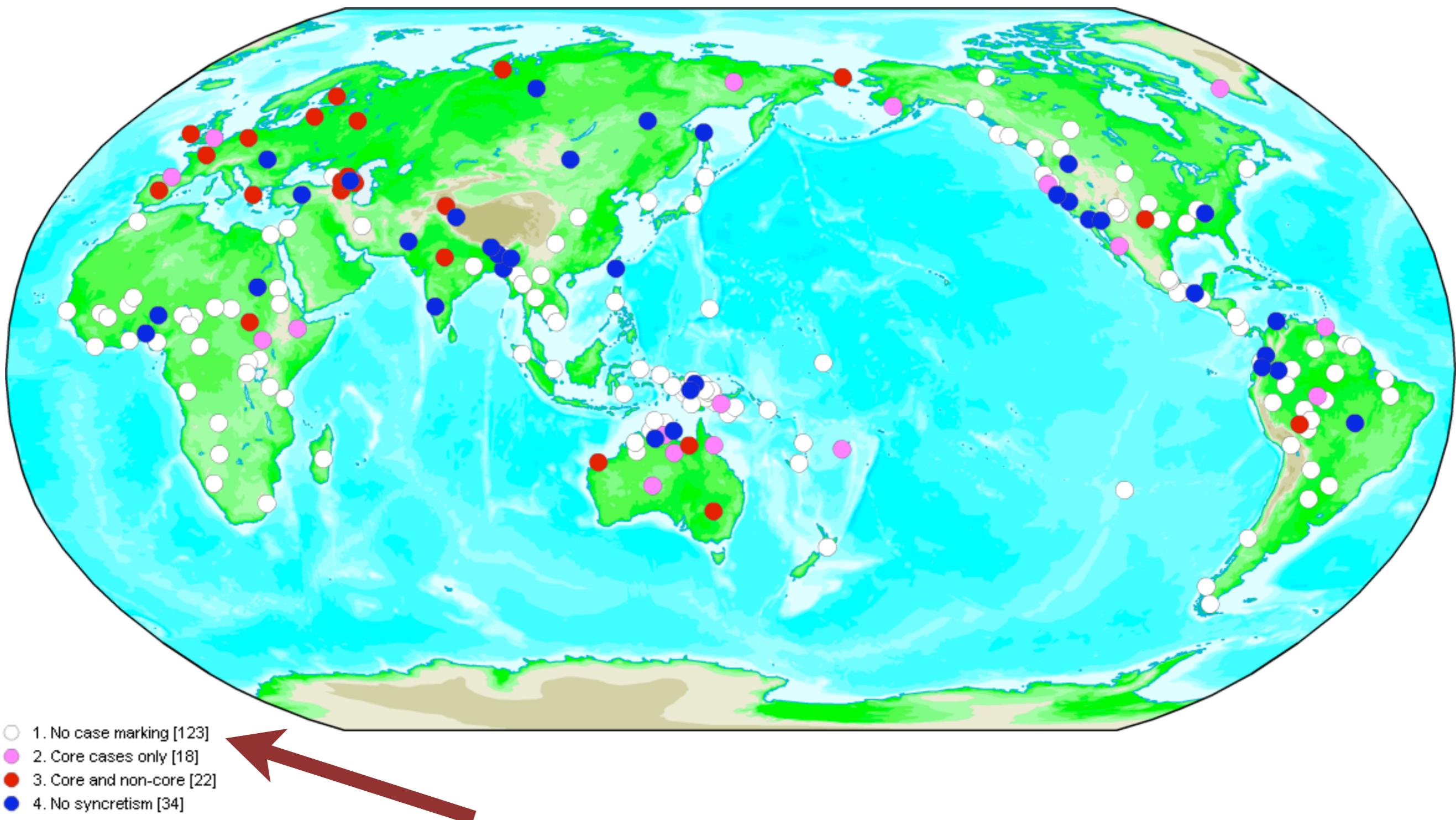
- Types consisting of dissimilar languages
- Independent features combined in one map
- **Definitional dependencies between maps**
- No relative similarities available

Problem:
(Apparently) identical values

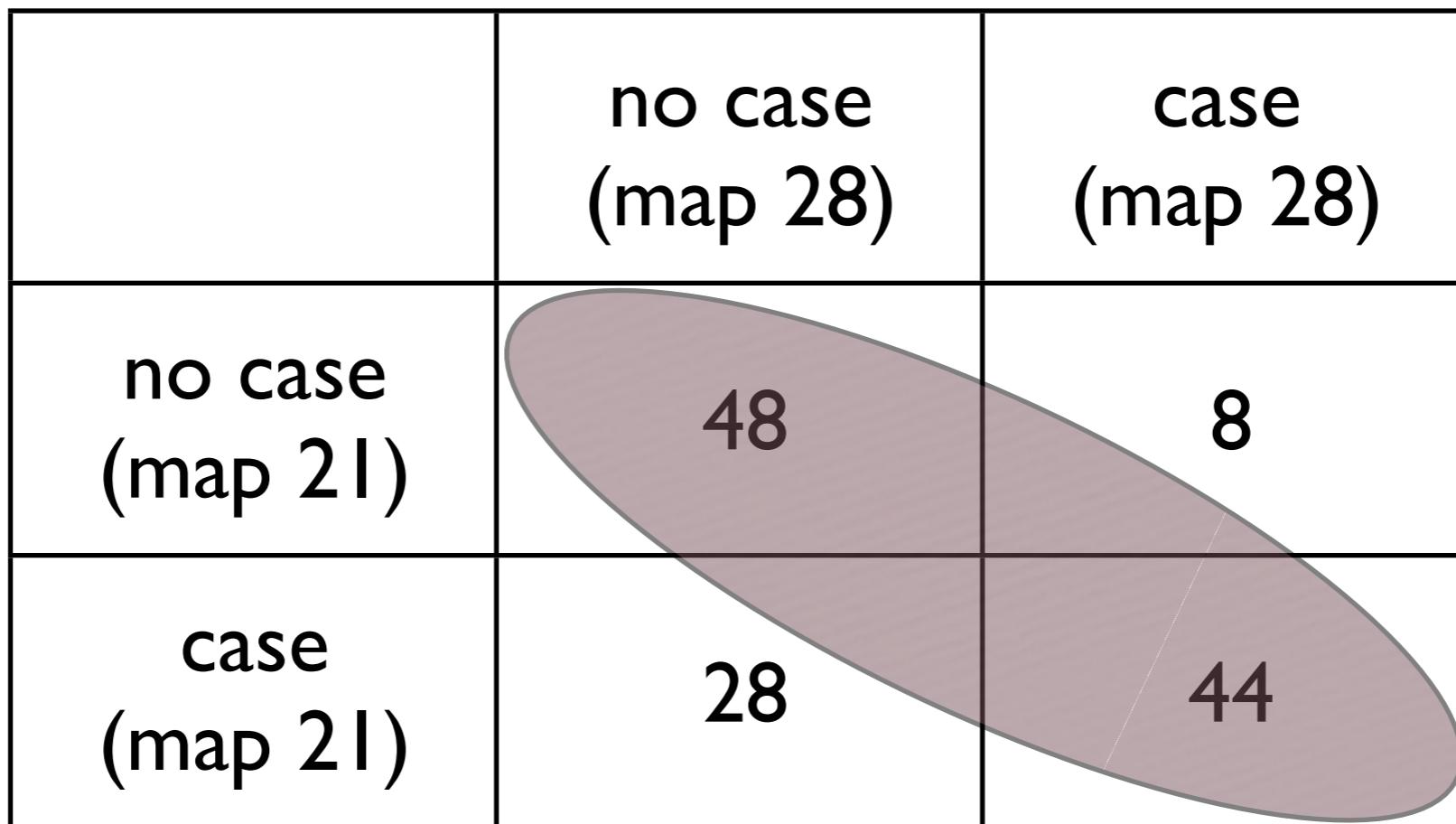
Map 2I: Exponence of Selected Inflectional Formatives



Map 28: Case Syncretism

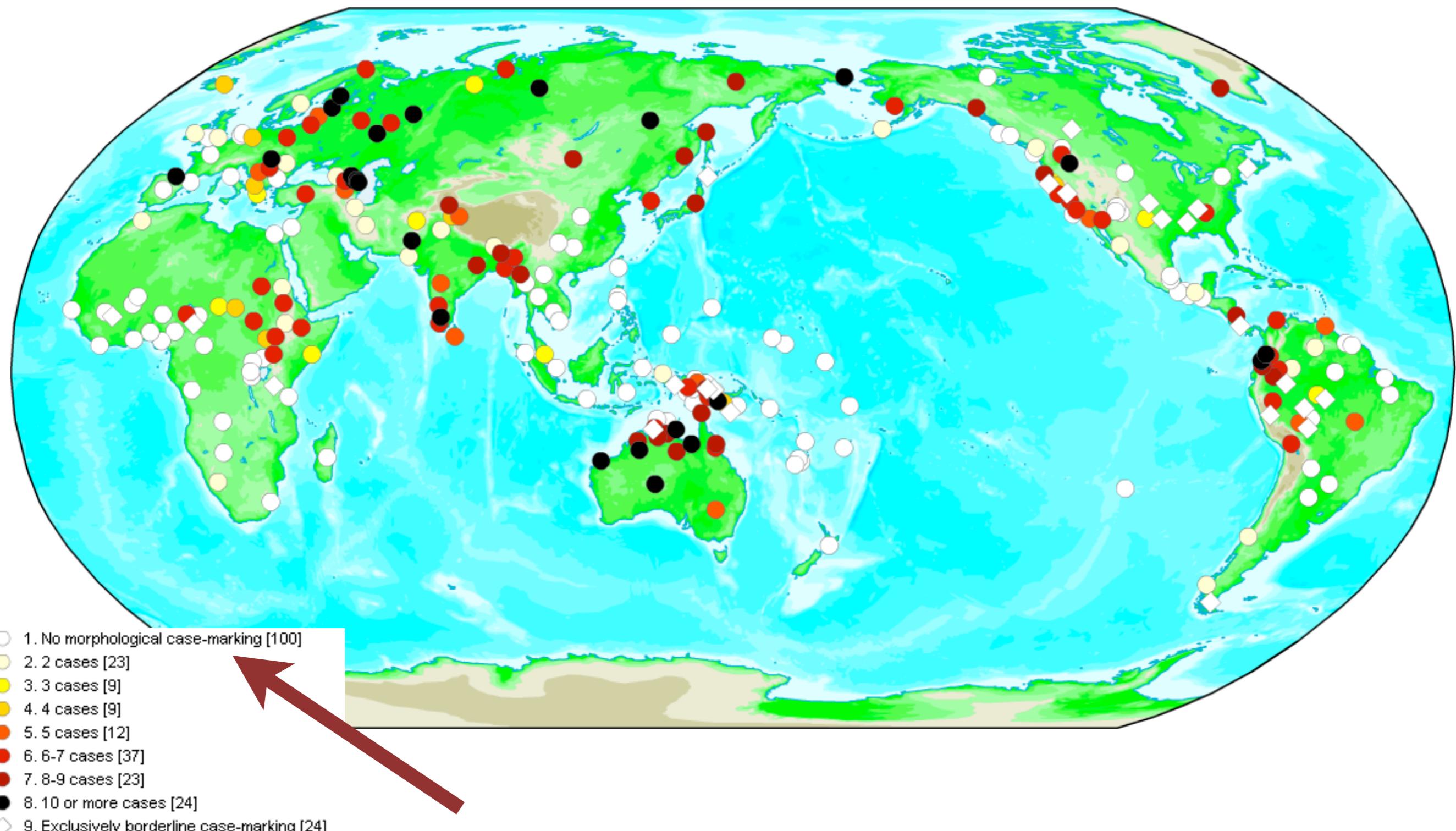


Different definitions and interpretation

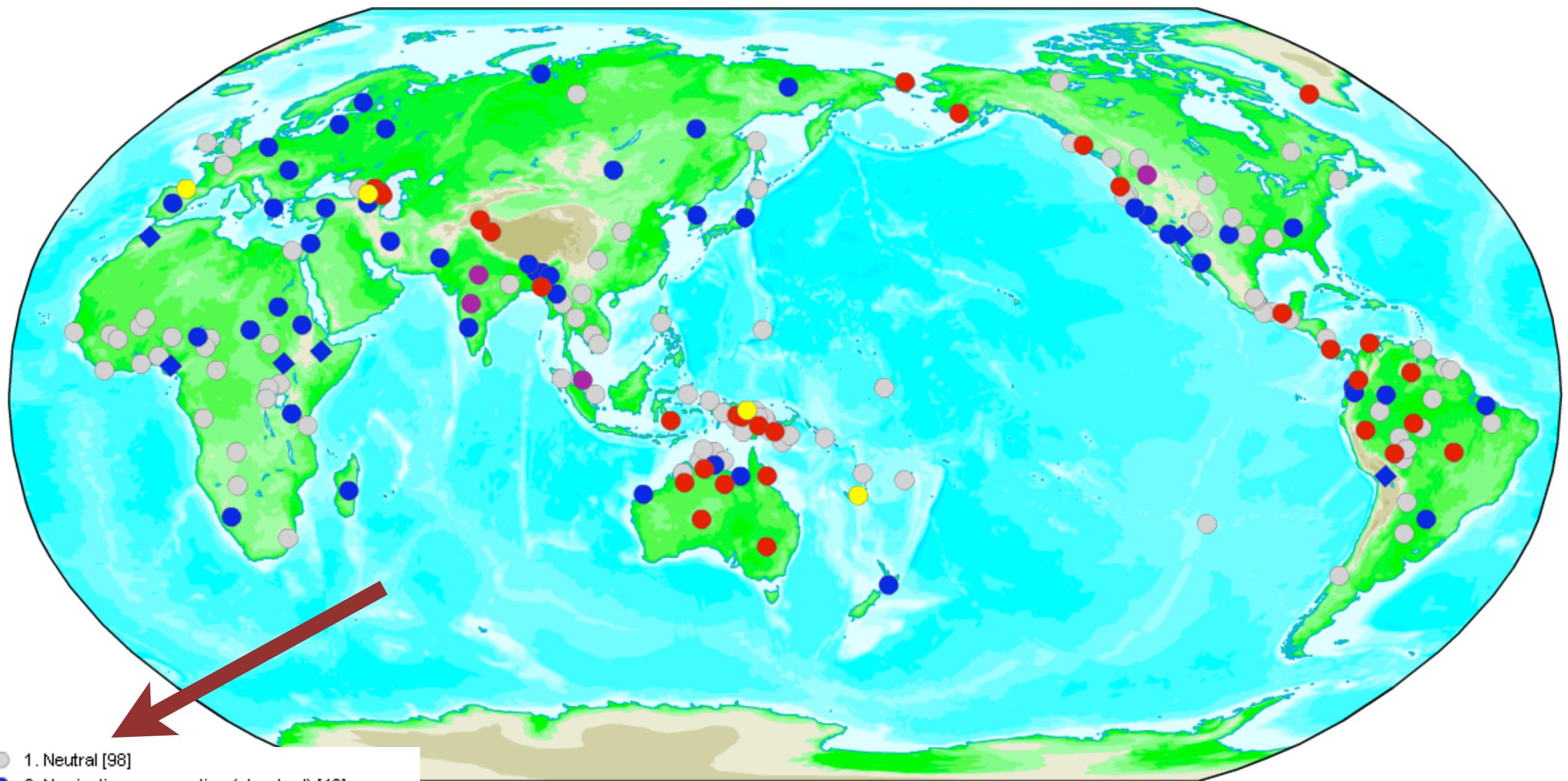


Problem:
Covert Dependencies

Map 49: Number of Cases

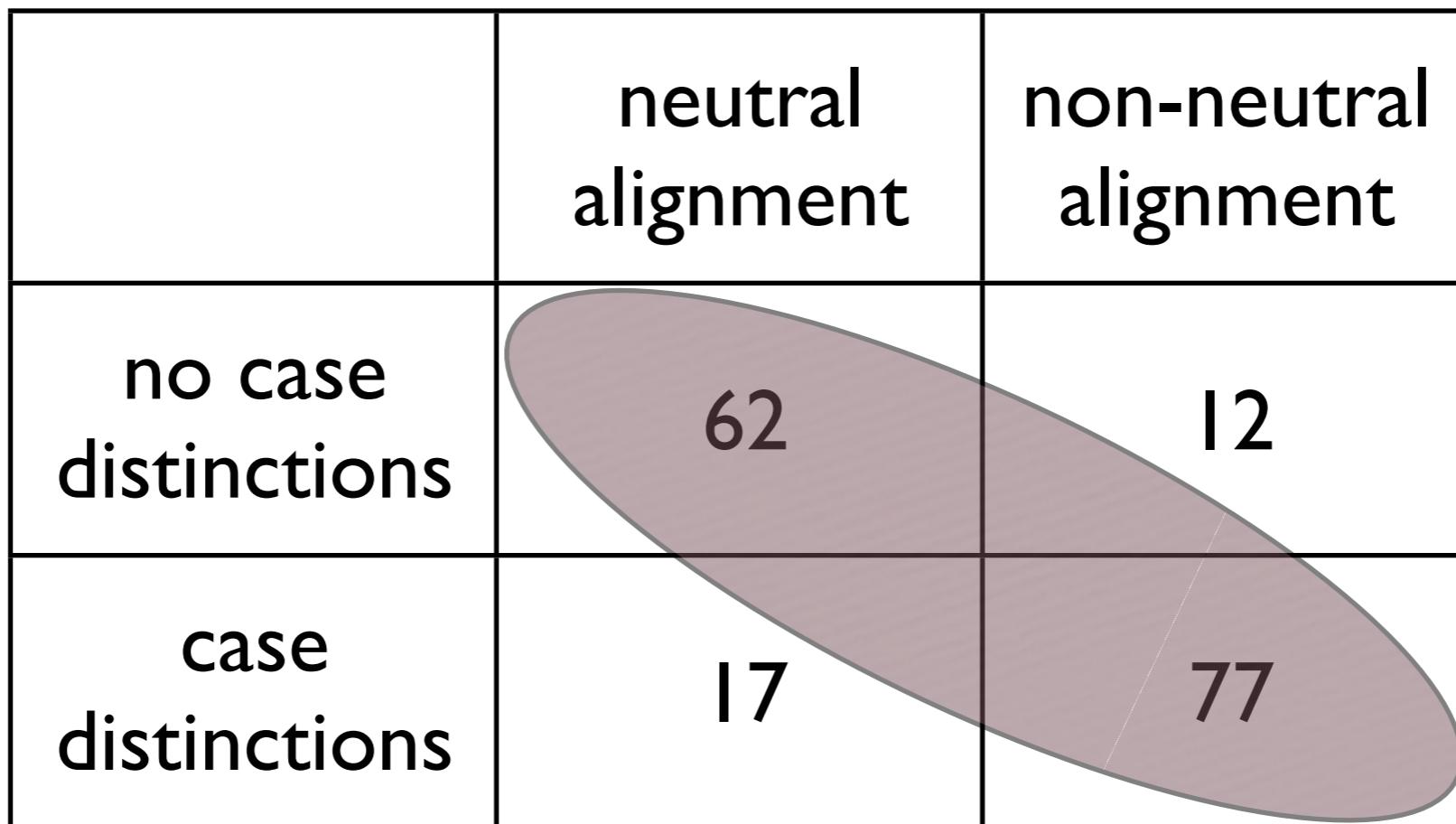


Map 98: Alignment of case marking of full noun phrases



- 1. Neutral [98]
- 2. Nominative - accusative (standard) [46]
- 3. Nominative - accusative (marked nominative) [6]
- 4. Ergative - absolute [32]
- 5. Tripartite [4]
- 6. Active-inactive [4]

Marking of full noun phrases



Solution:

**Combine features into one
feature with (very) many values**

Coding Problems

- Types consisting of dissimilar languages
- Independent features combined in one map
- **Definitional dependencies between maps**
- No relative similarities available

Coding Problems

- Types consisting of dissimilar languages
- Independent features combined in one map
- Definitional dependencies between maps
- **No relative similarities available**

Map 51: Position of Case affixes

WALS – World Atlas of Language Structures

LANGUAGE VIEWER **COMPOSER** **SHOW MAP**

select a feature

- ▶ thematically
- ▶ alphabetically
- ▶ user-defined

SHRINK LIST

search for a feature
51 **SEARCH**

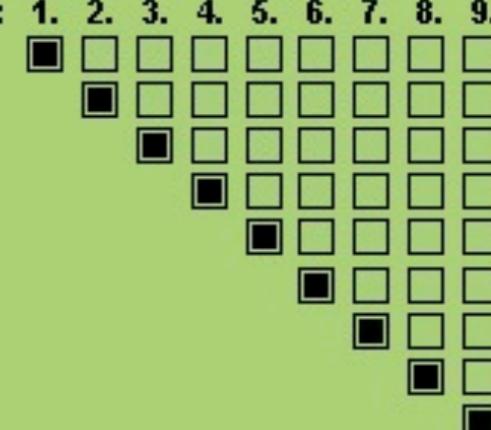
FEATURE PROFILE area: Nominal Categories

51. Position of Case Affixes
Author: Matthew S. Dryer
934 languages

symbol: include: click to list languages below [no. of lgs : of genera : of families]

	<input type="checkbox"/>	1. Case suffixes [431:174:90]	<input type="checkbox"/>						
	<input type="checkbox"/>	2. Case prefixes [35:19:14]	<input type="checkbox"/>						
	<input type="checkbox"/>	3. Case tone [4:2:1]	<input type="checkbox"/>						
	<input type="checkbox"/>	4. Case stem change [2:1:1]	<input type="checkbox"/>						
	<input type="checkbox"/>	5. Mixed morphological case [8:7:6]	<input type="checkbox"/>						
	<input type="checkbox"/>	6. Postpositional clitics [95:59:36]	<input type="checkbox"/>						
	<input type="checkbox"/>	7. Prepositional clitics [15:10:8]	<input type="checkbox"/>						
	<input type="checkbox"/>	8. Inpositional clitics [6:3:1]	<input type="checkbox"/>						
	<input type="checkbox"/>	9. No case affixes or adpositional clitics [338:145:56]	<input type="checkbox"/>						

Merge: 1. 2. 3. 4. 5. 6. 7. 8. 9.

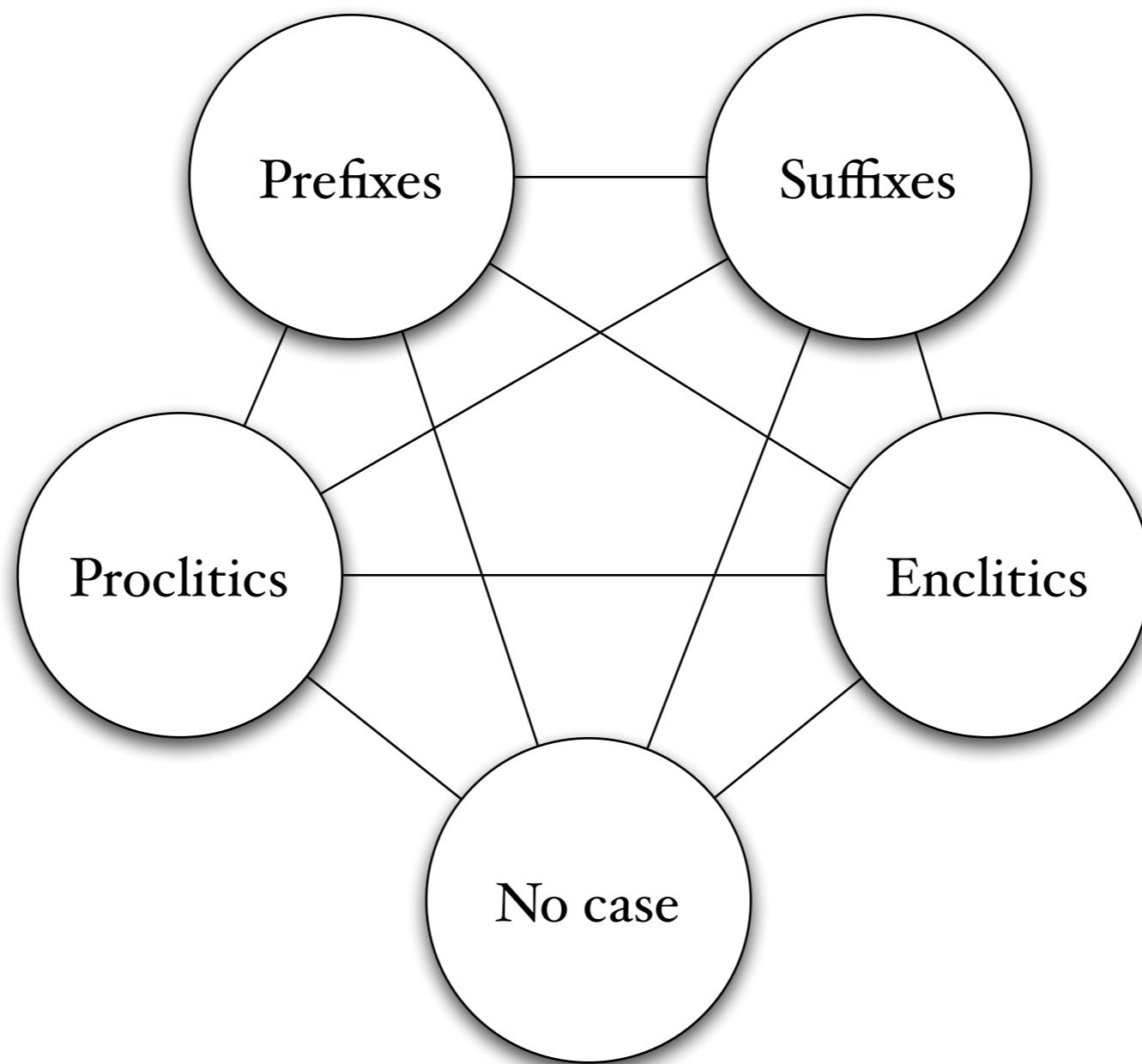


DESCRIPTION

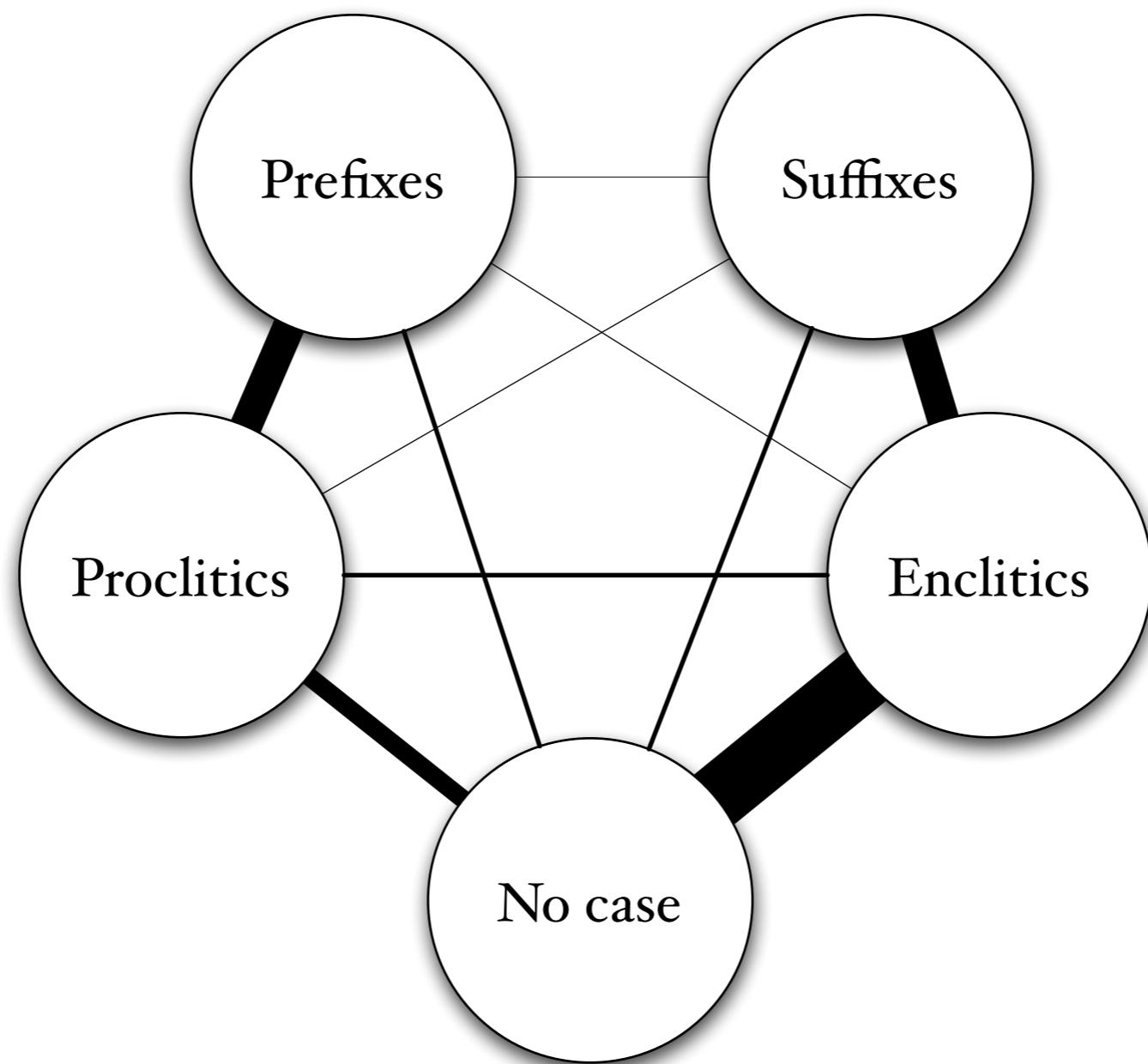
arrange the languages by
languages ▾

COPY LIST

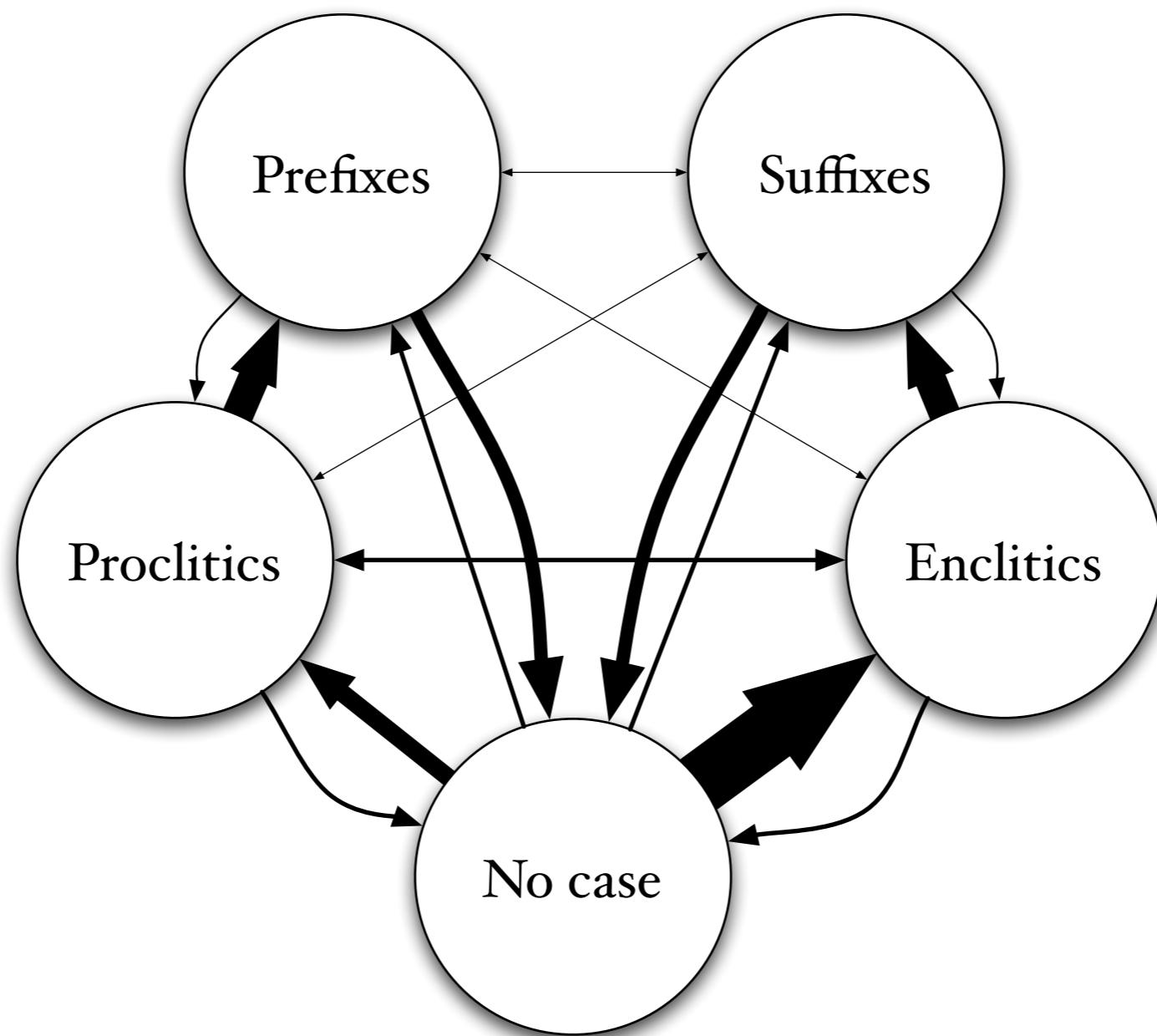
Undifferentiated typology



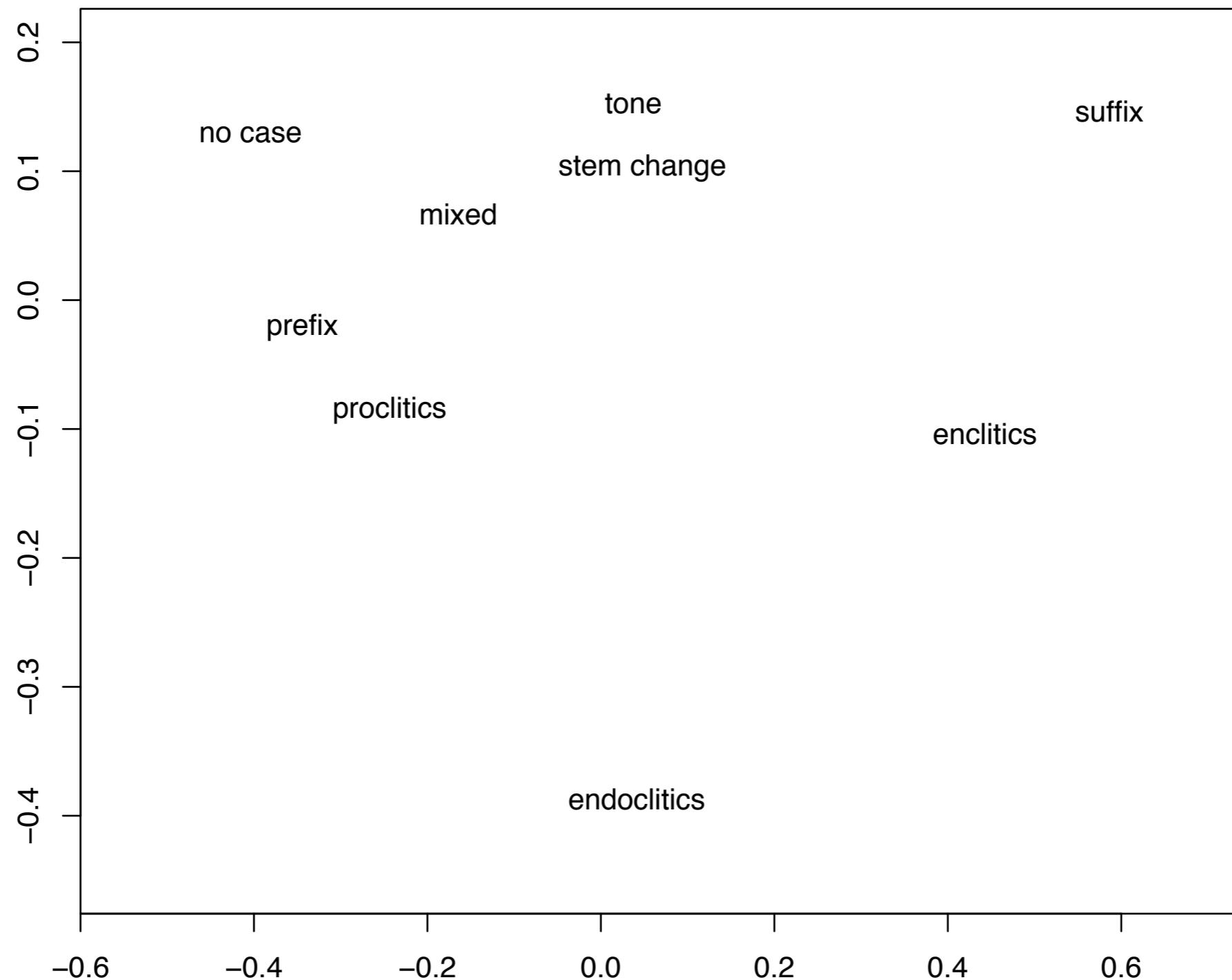
Similarities



Transitional probabilities



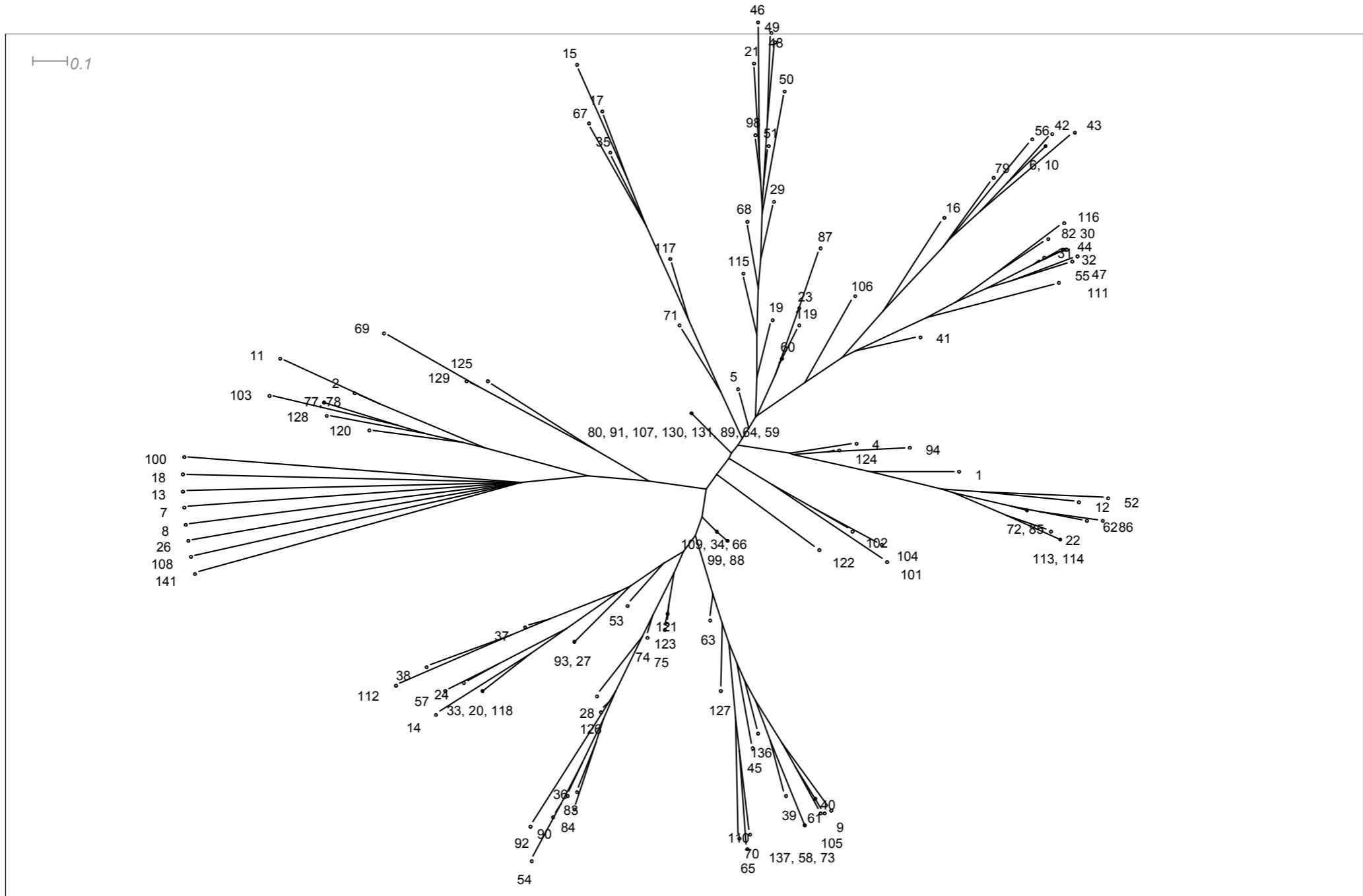
Estimate similarities from low-level genetic groups



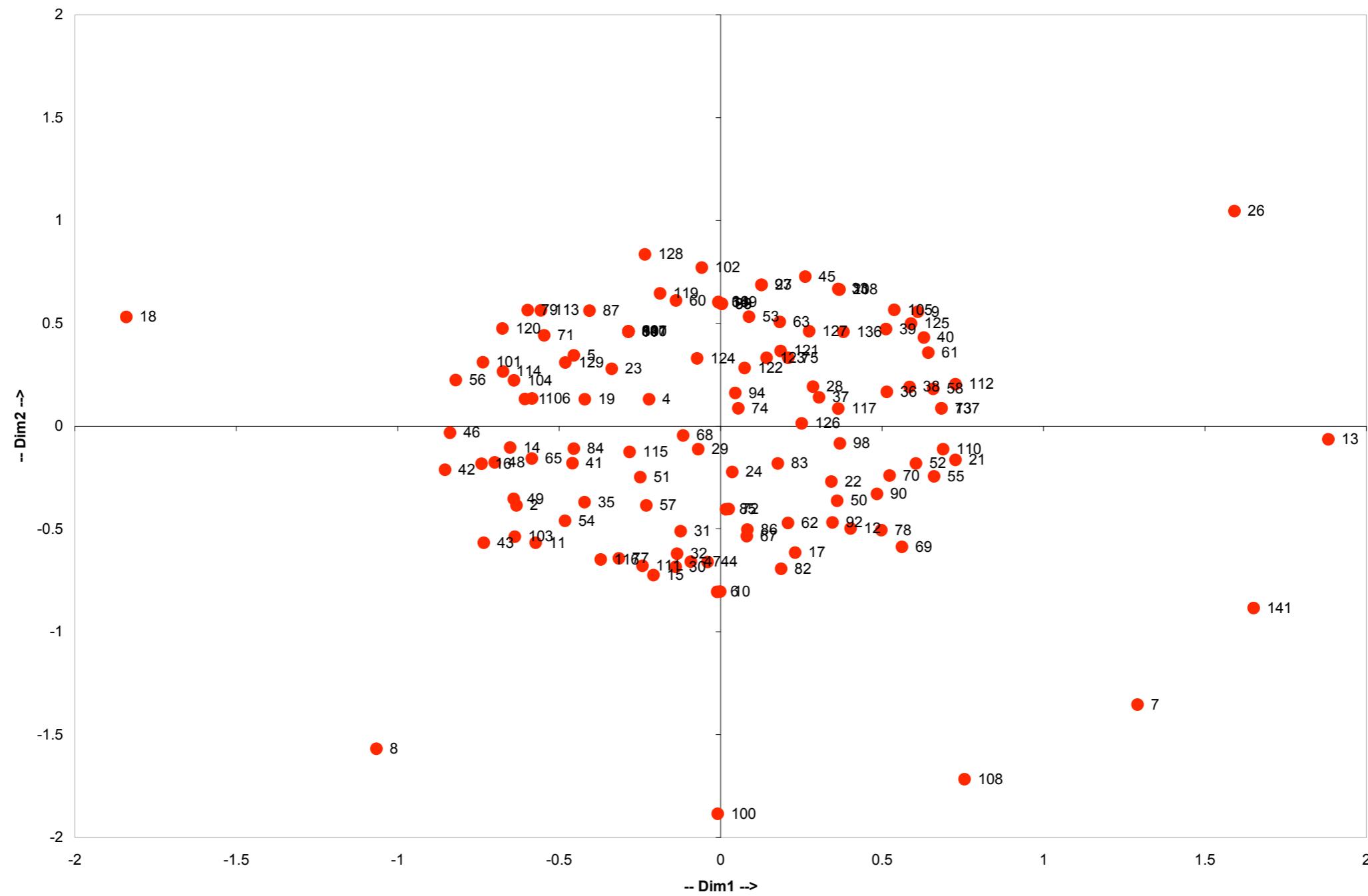
Exploration

Clustering of features

Looks promising...



... but it is an artifact

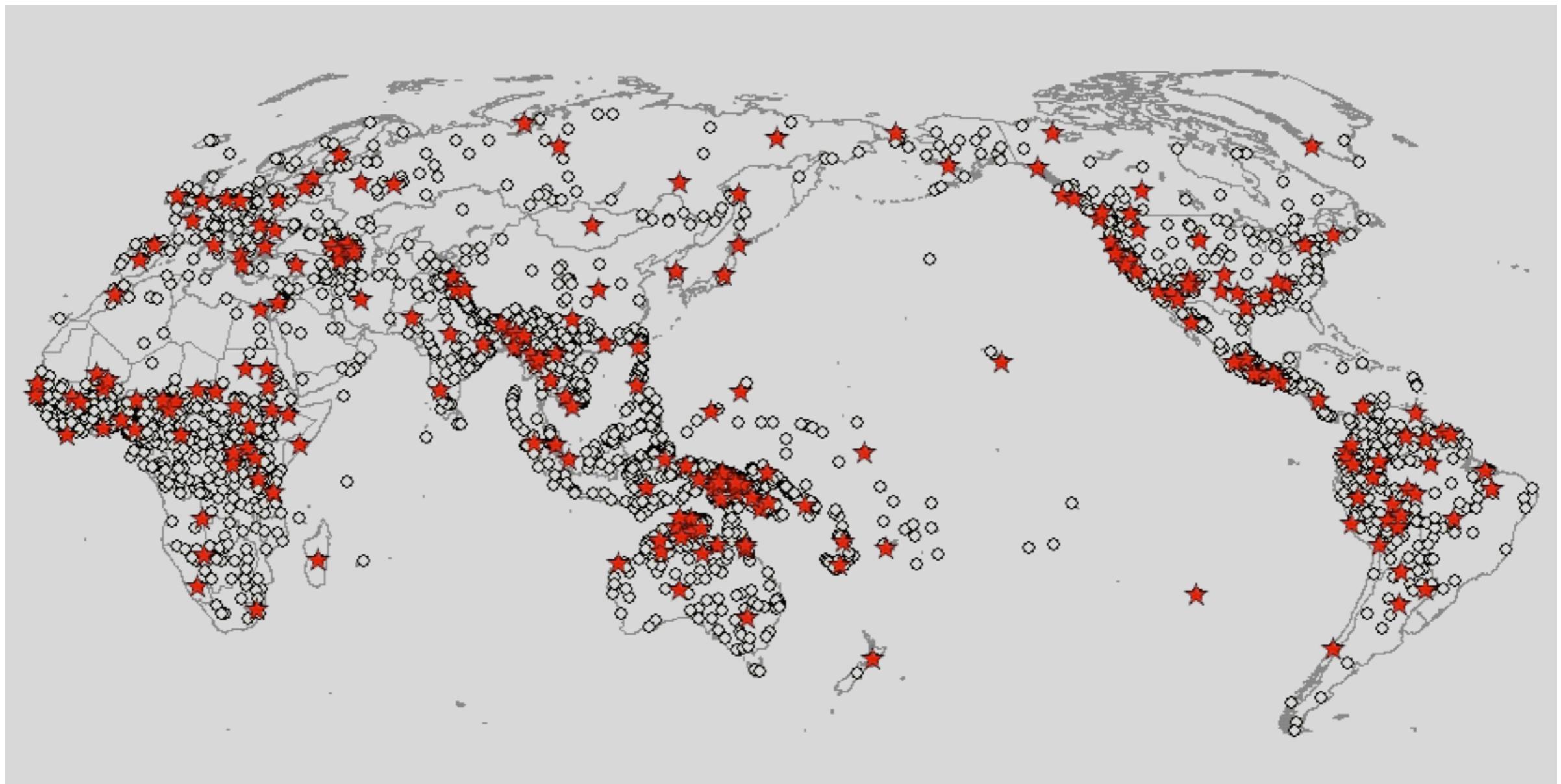


Evaluation of Sampling (by Östen Dahl)

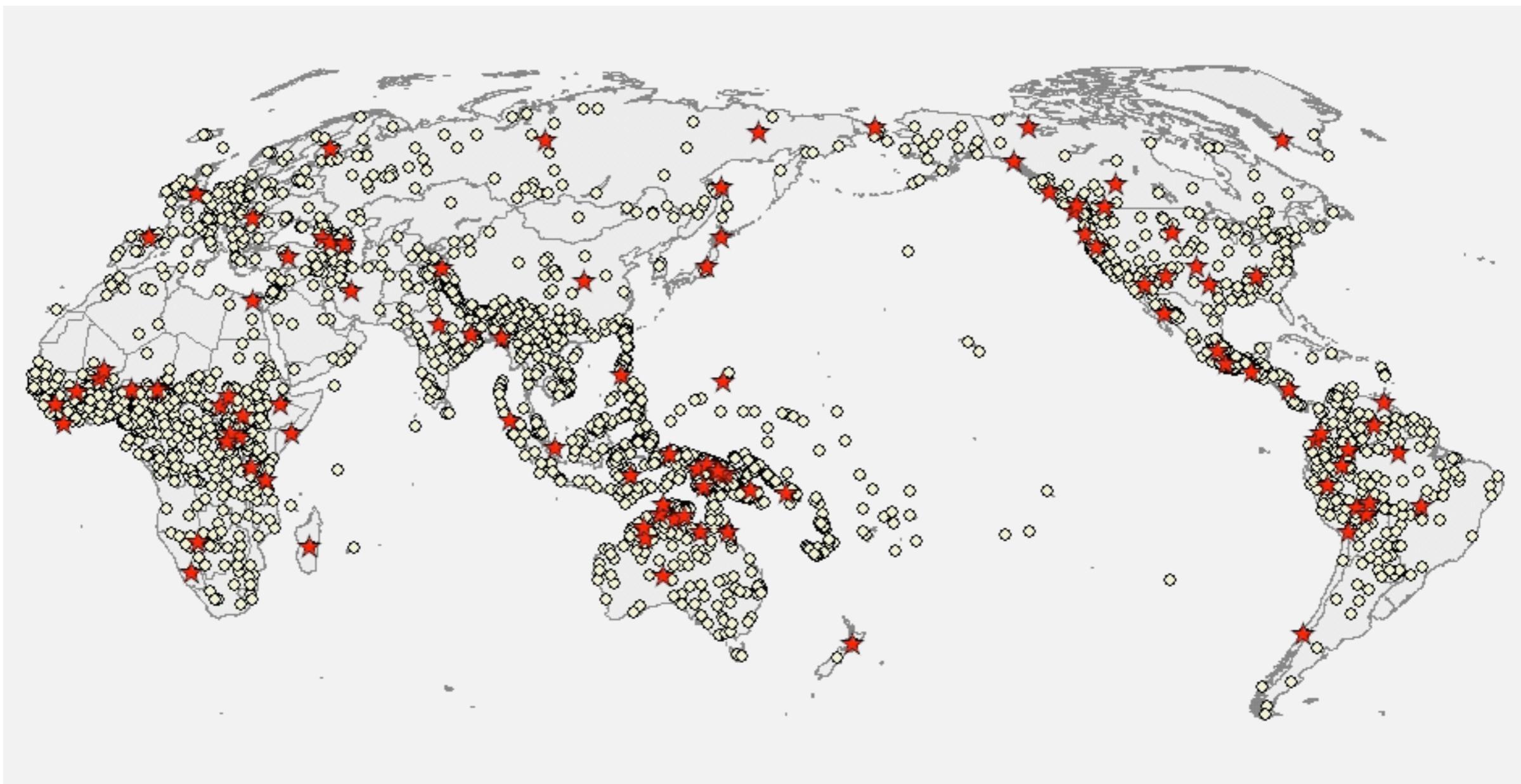
Which languages to investigate ?

- Typological samples normally take one language from each genealogical group
- But: genealogical groups are mostly based on lexical & phonological characteristics
- Is this a good basis to investigate syntax or morphology ?

Genealogical 200 language sample



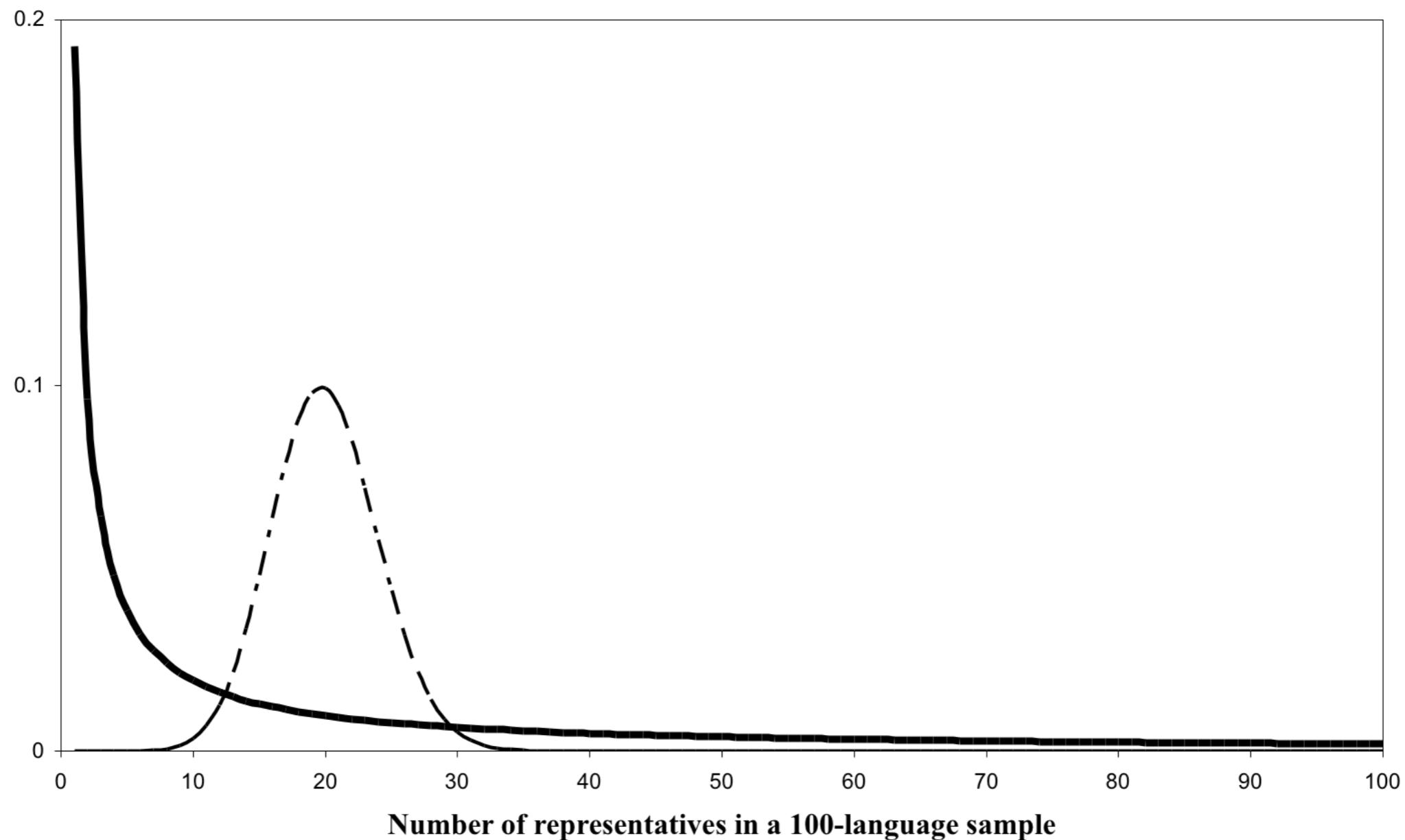
Typological 100 language sample



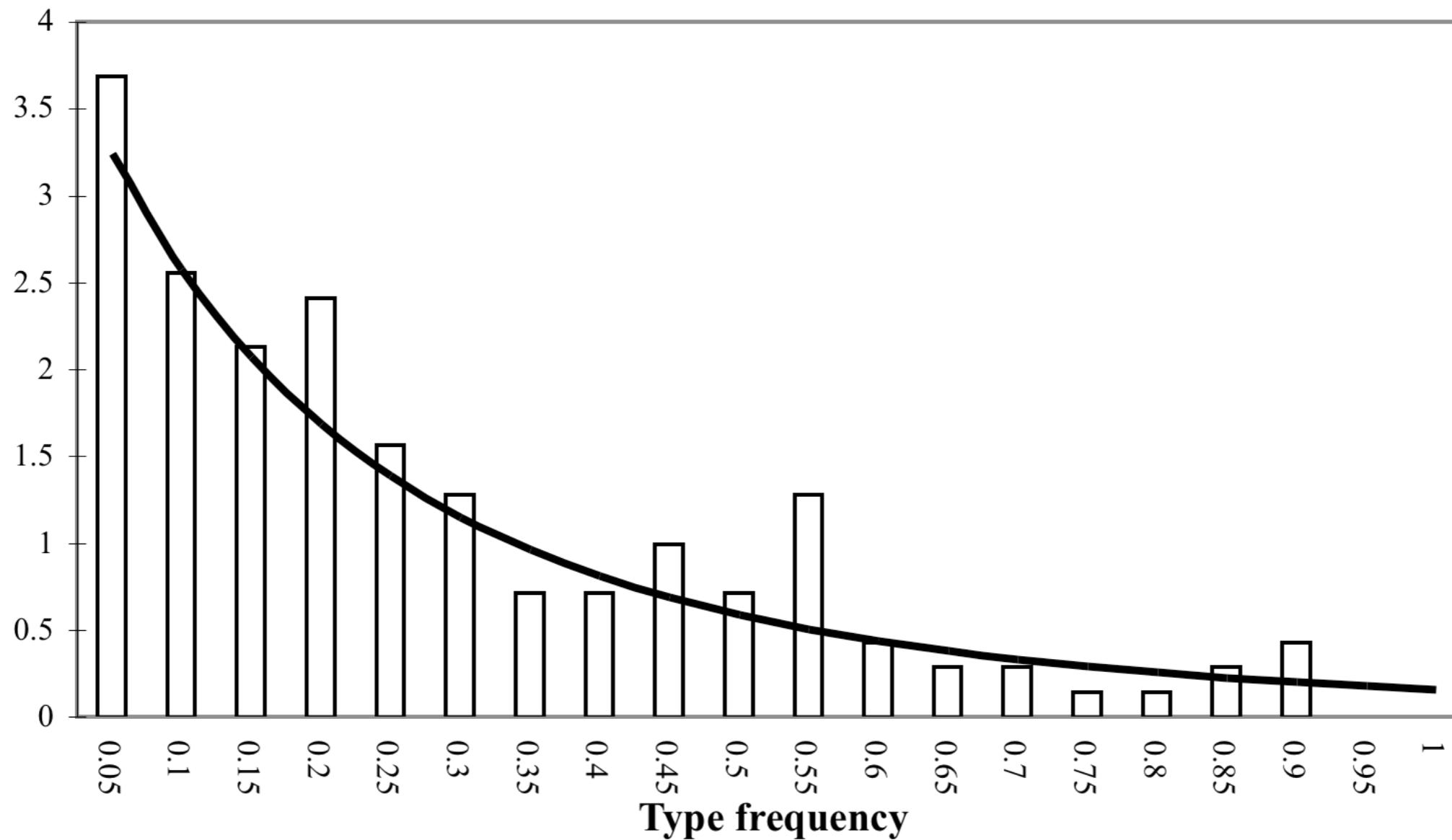
Evaluation of Expectations

(by Elena Maslova)

With five values, what do we expect to find ?



WALS lets us find out what we should expect



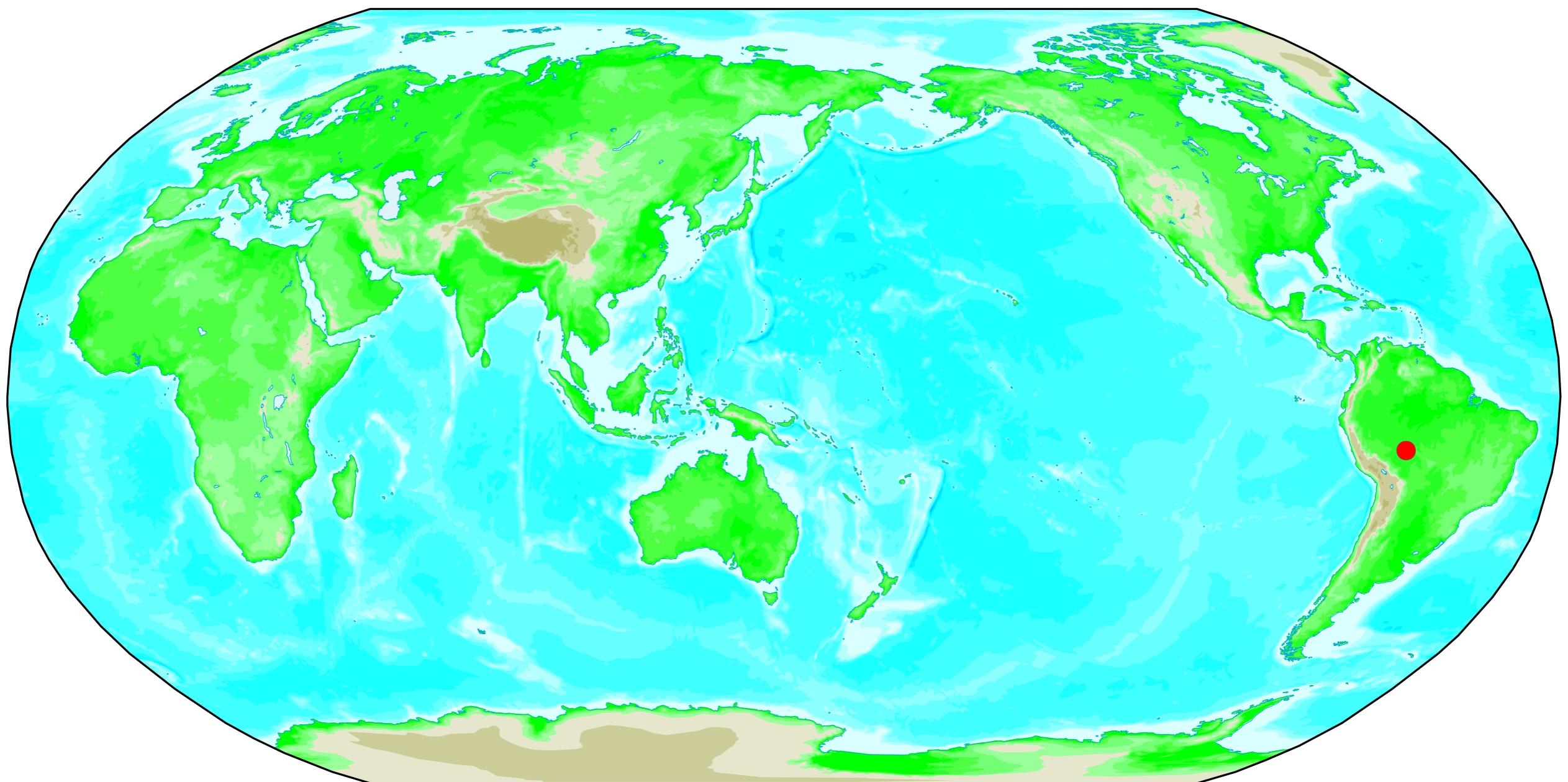
Investigating rare characteristics

And the winners are:

In the category:

‘Most Unusual
Individual Language’

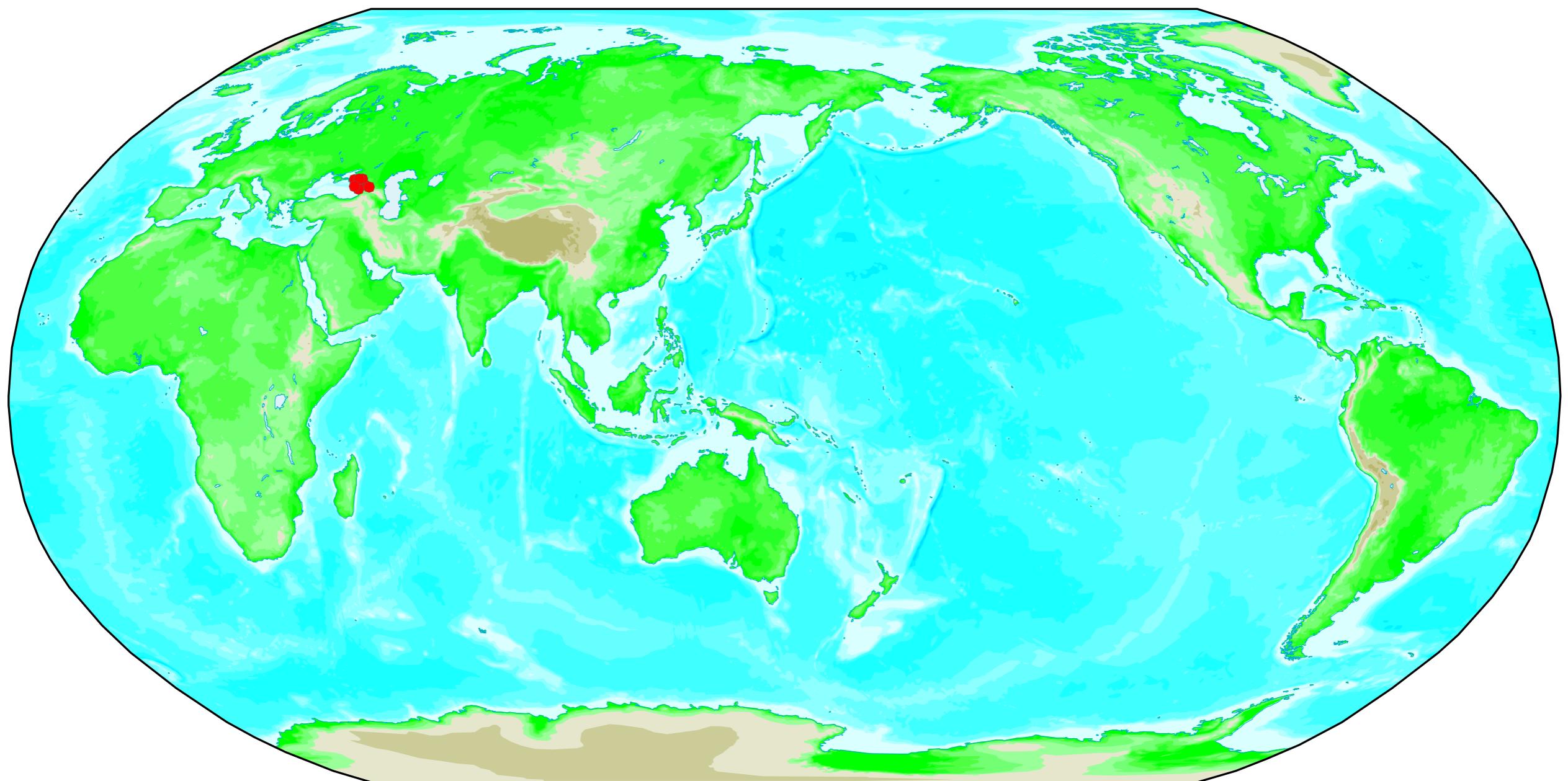
Wari'



In the category:

‘Most Unusual
Genealogical Group’

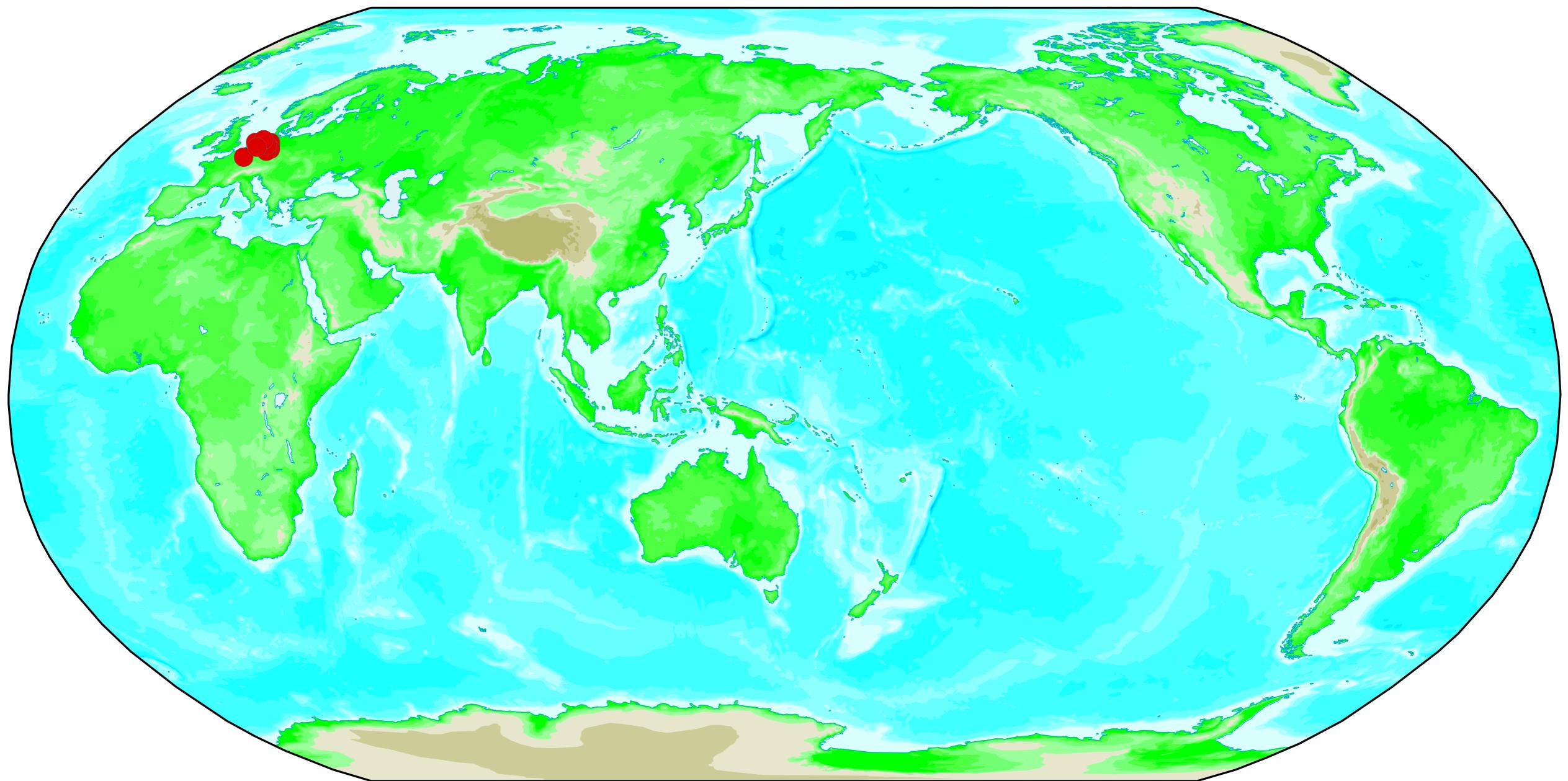
Northwest Caucasian



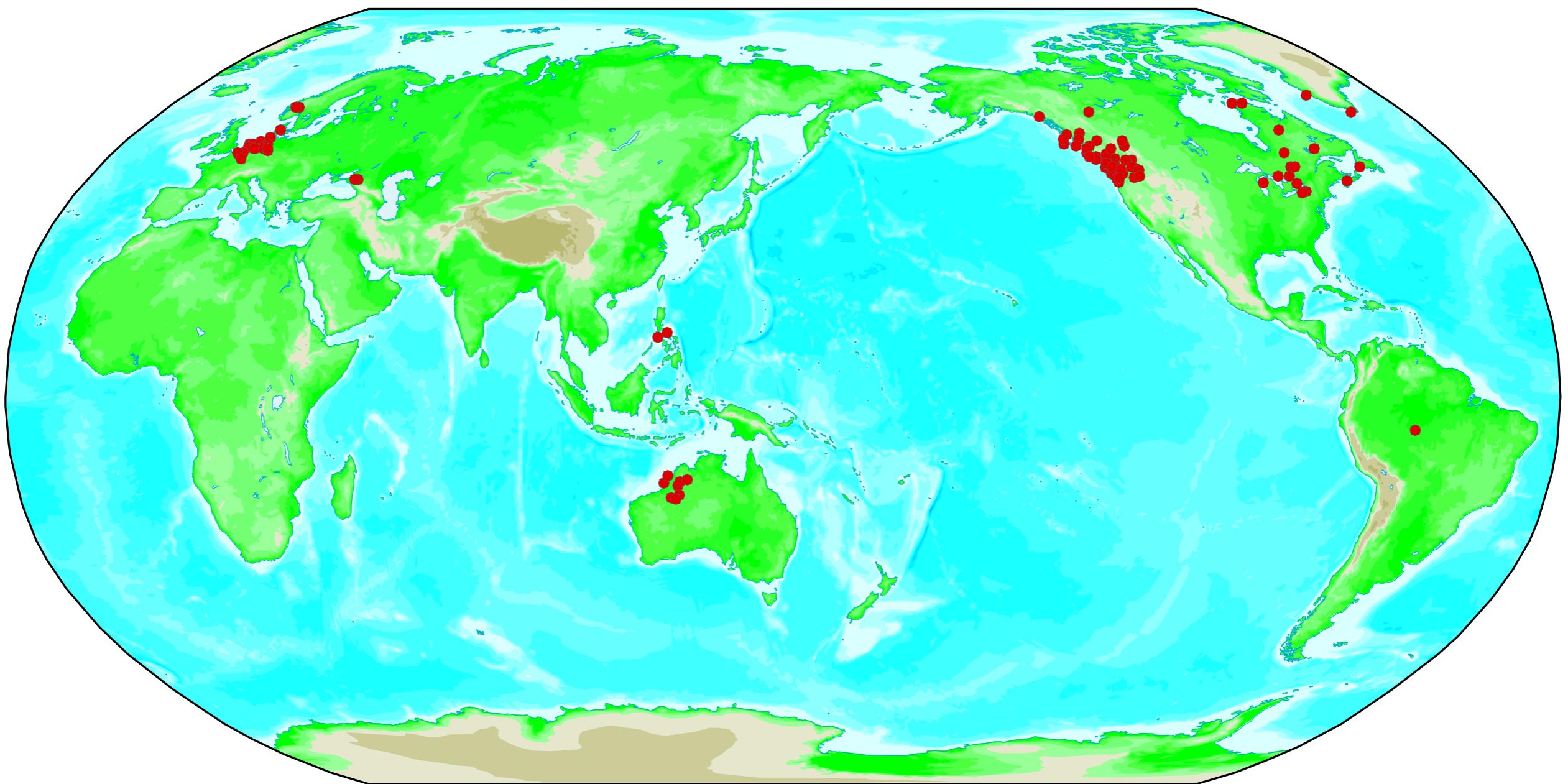
In the category:

‘Most Unusual
Geographical Area’

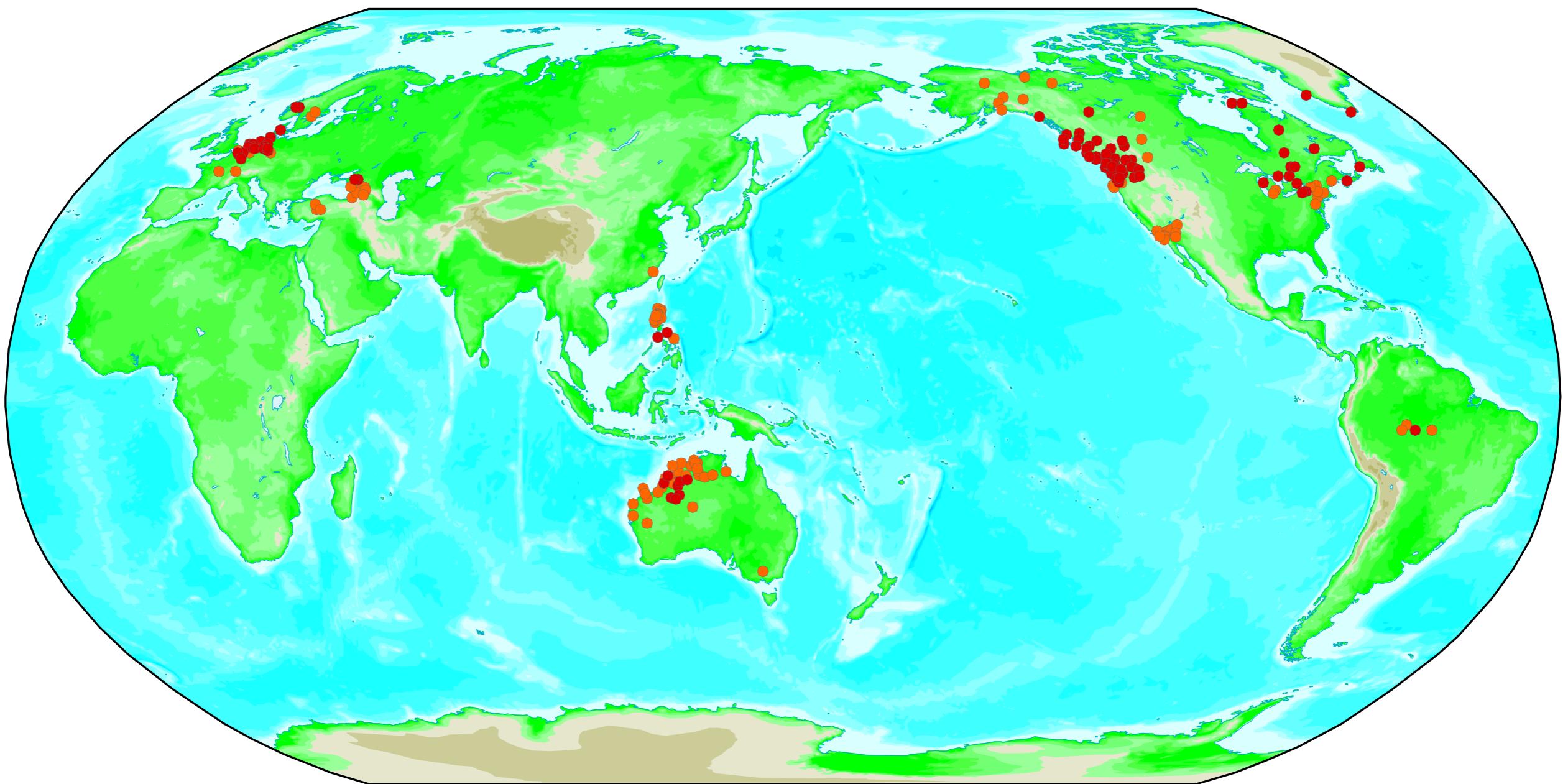
Northwest Continental Europe



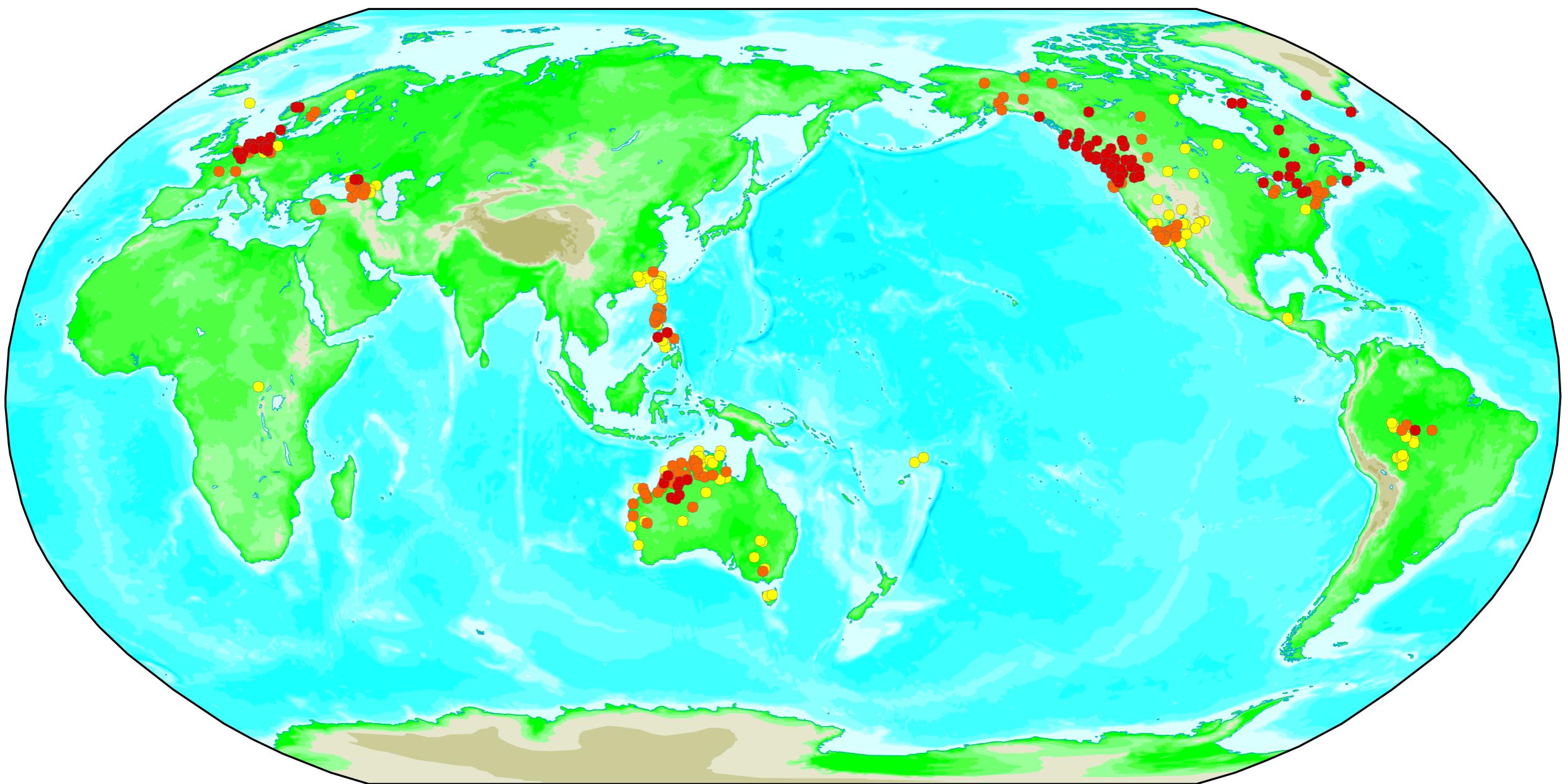
Top 100



Top 200

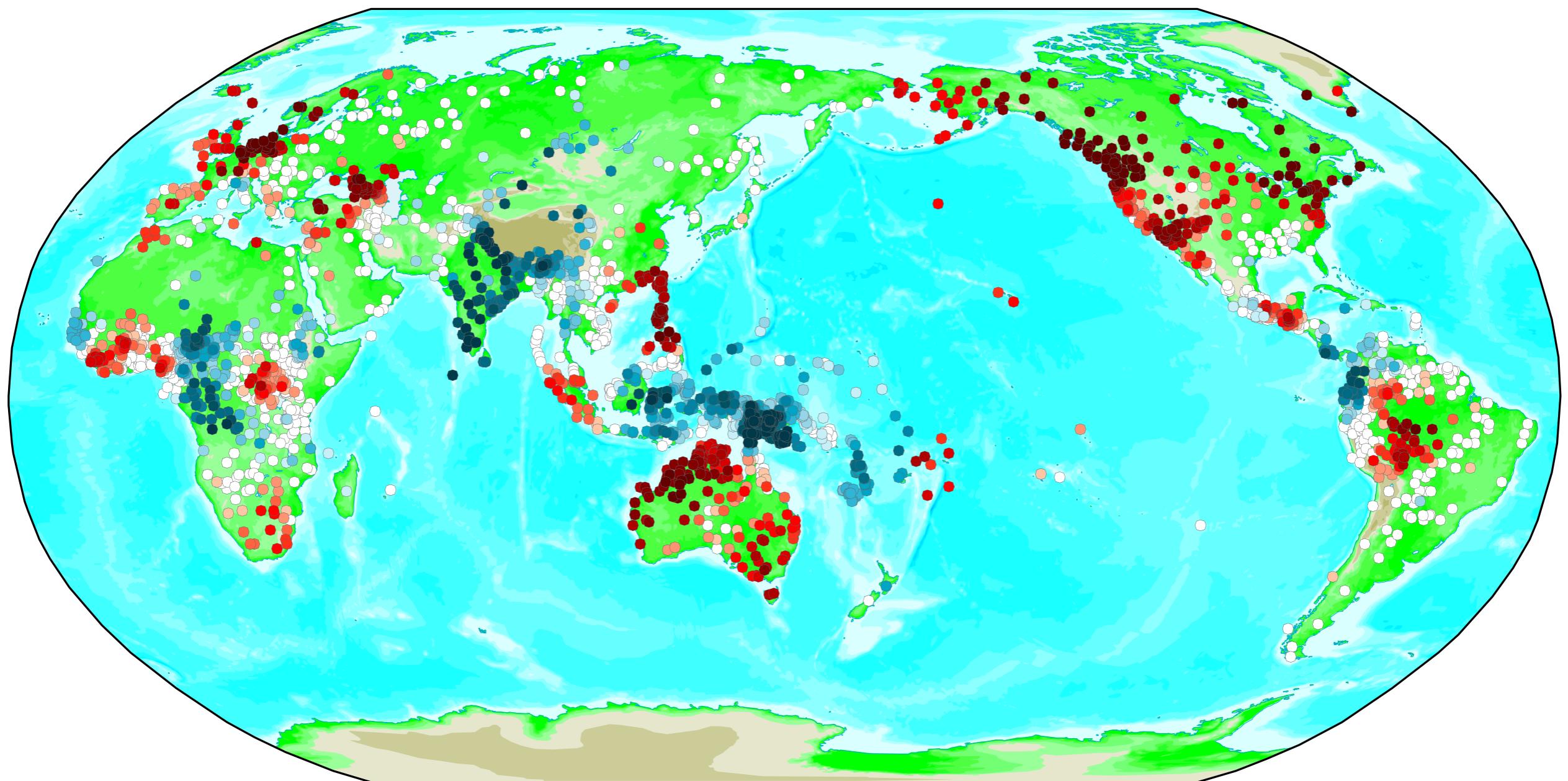


Top 300

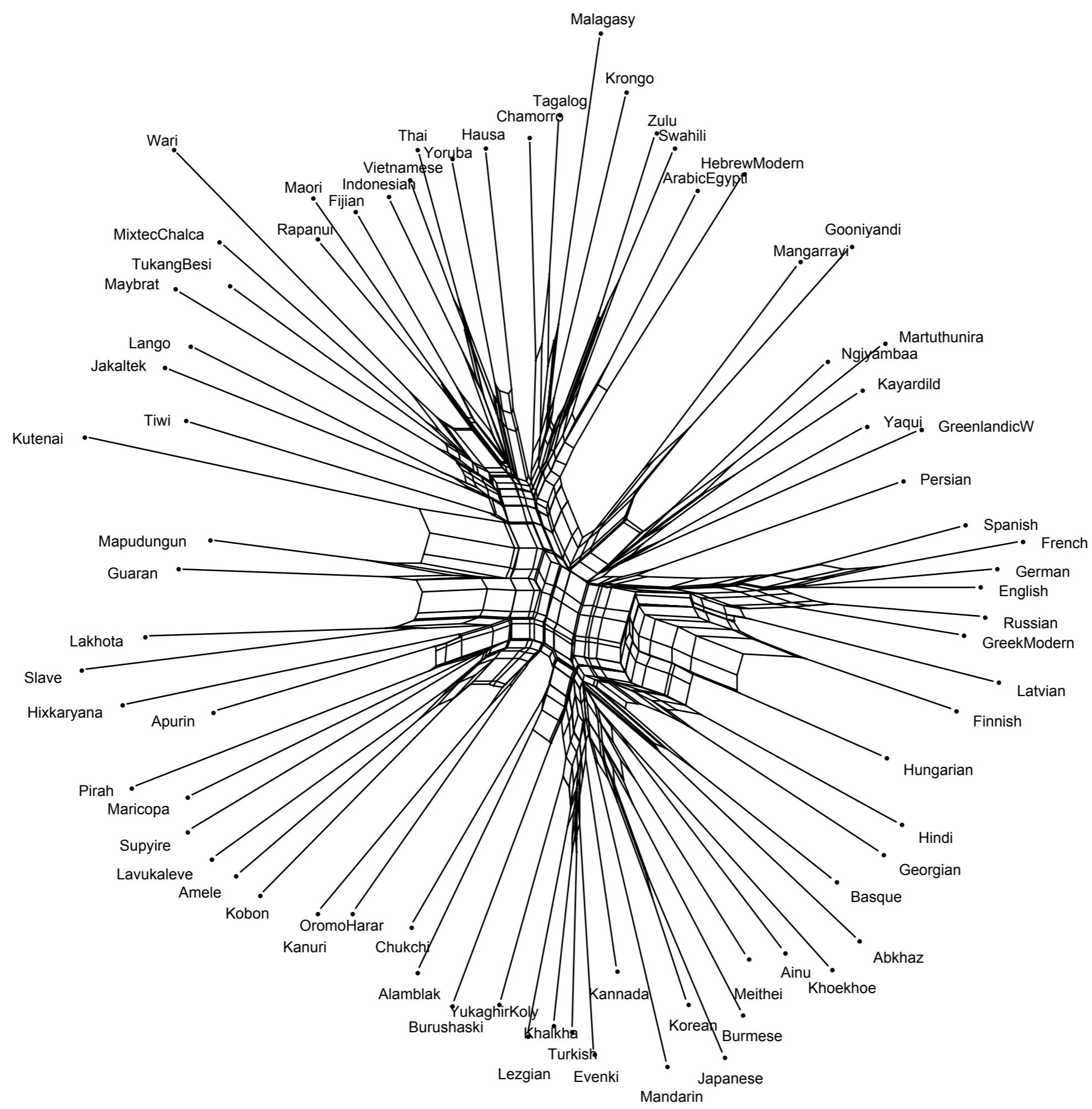


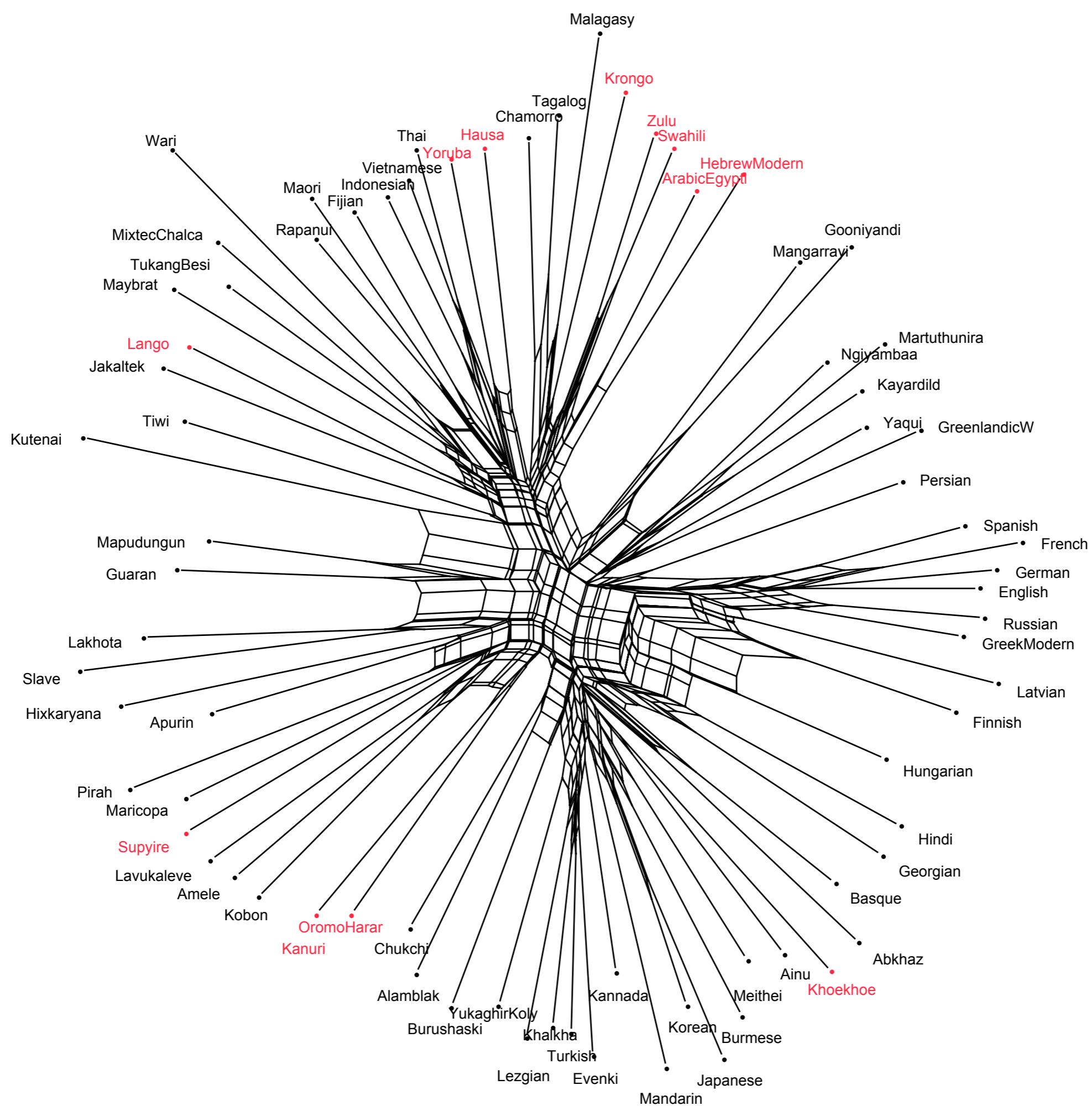
All languages

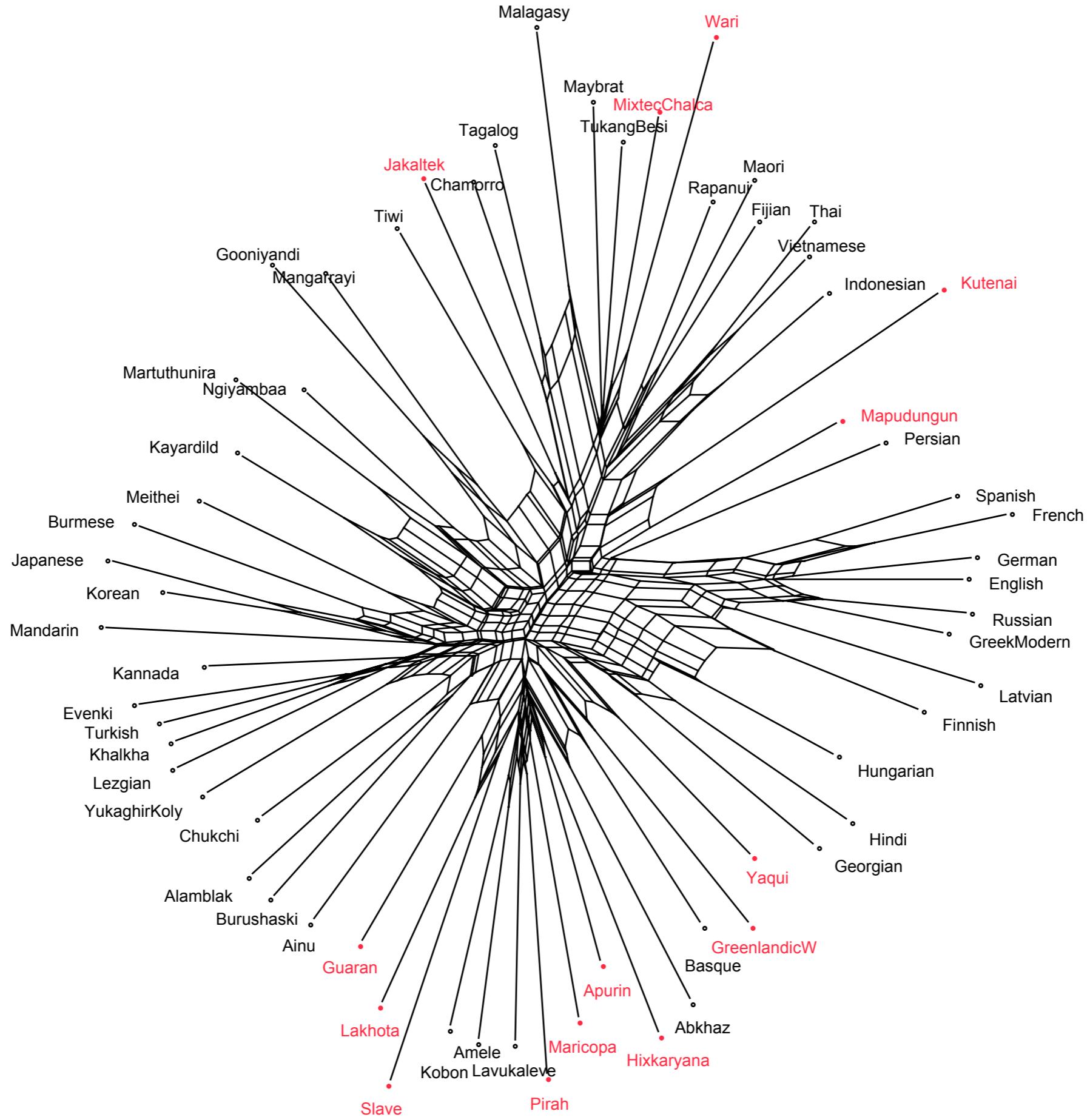
(red = rare, blue = common)

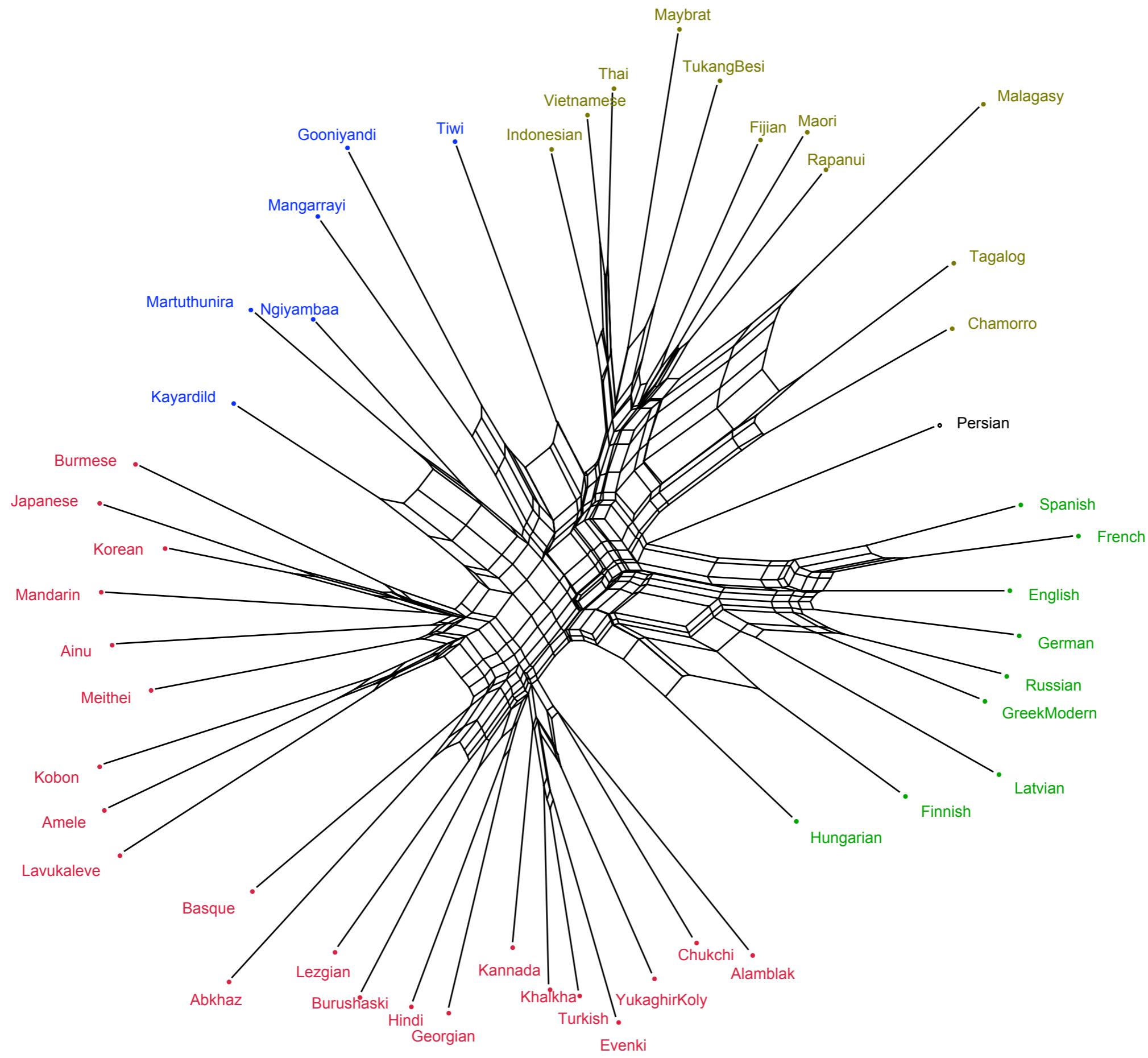


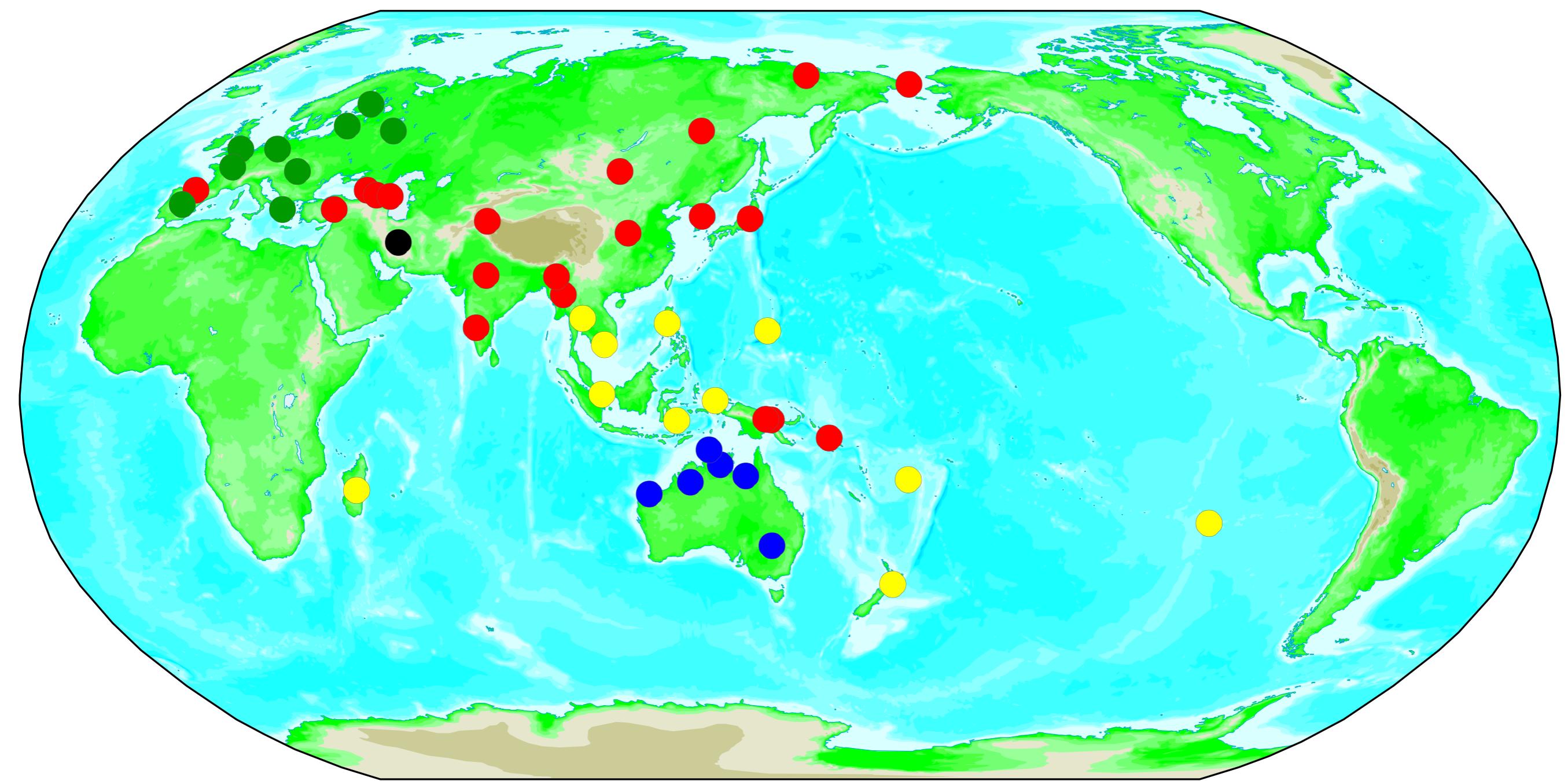
Investigating language similarity



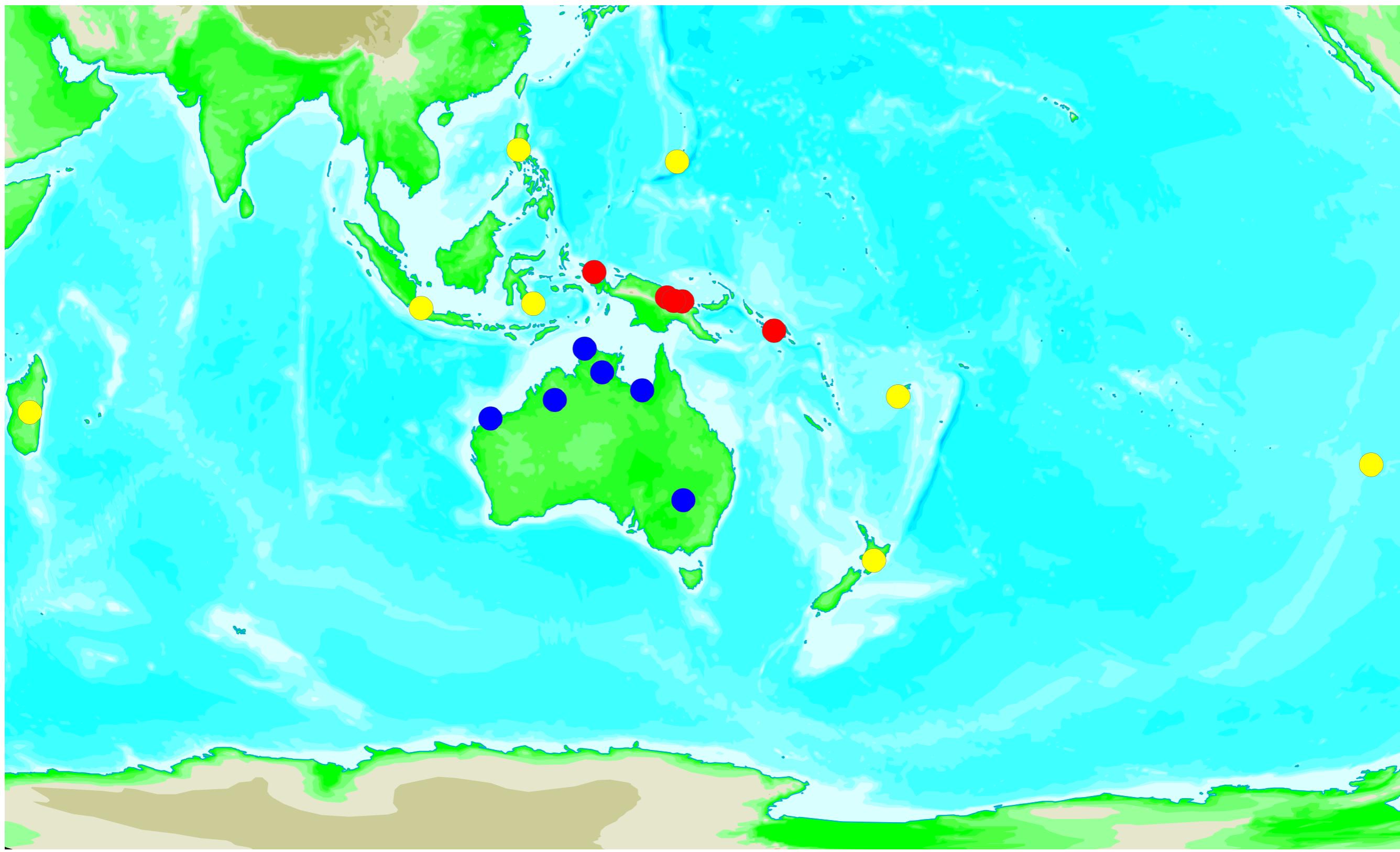








A closer look at geography: the case of Oceania

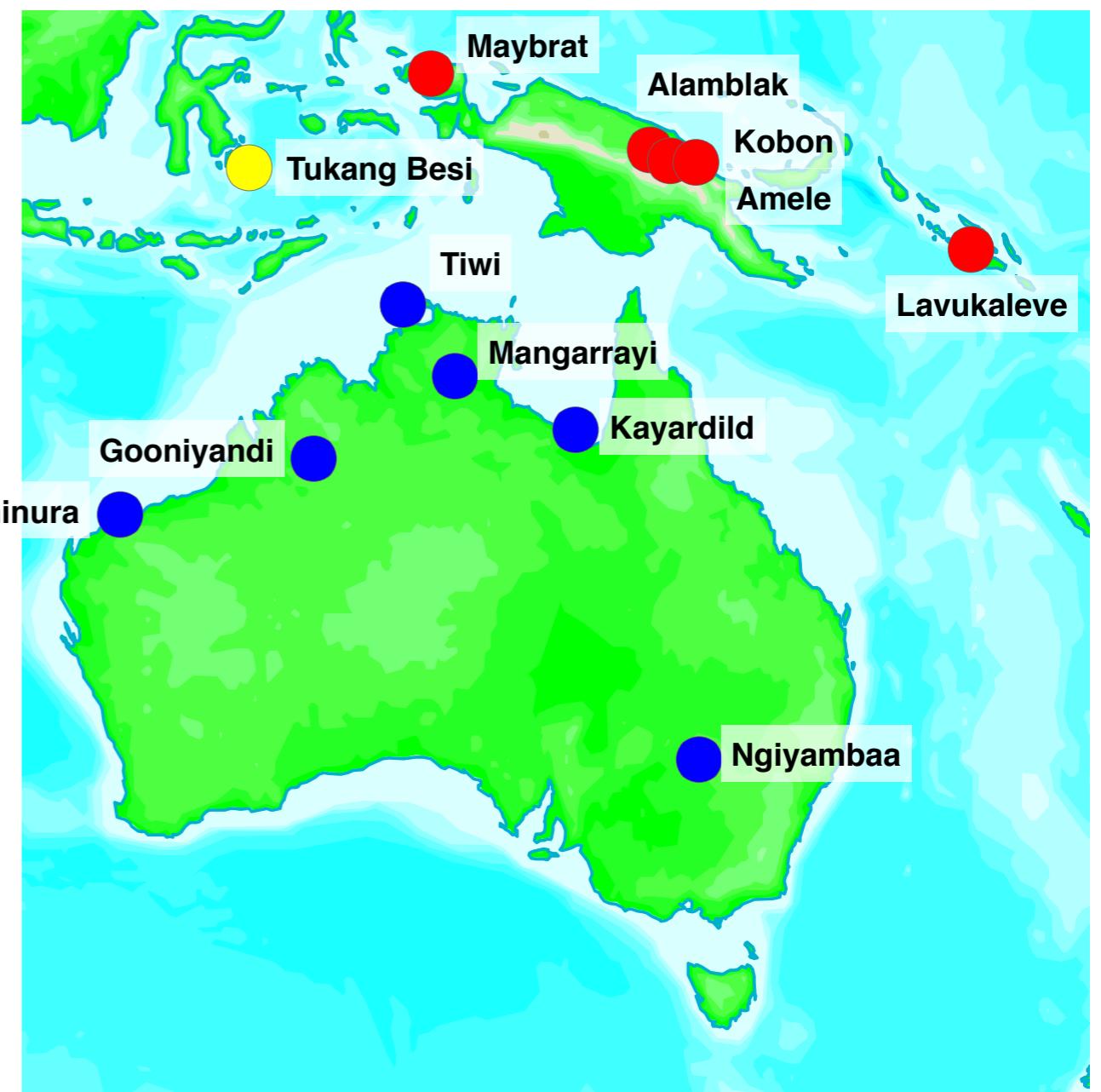
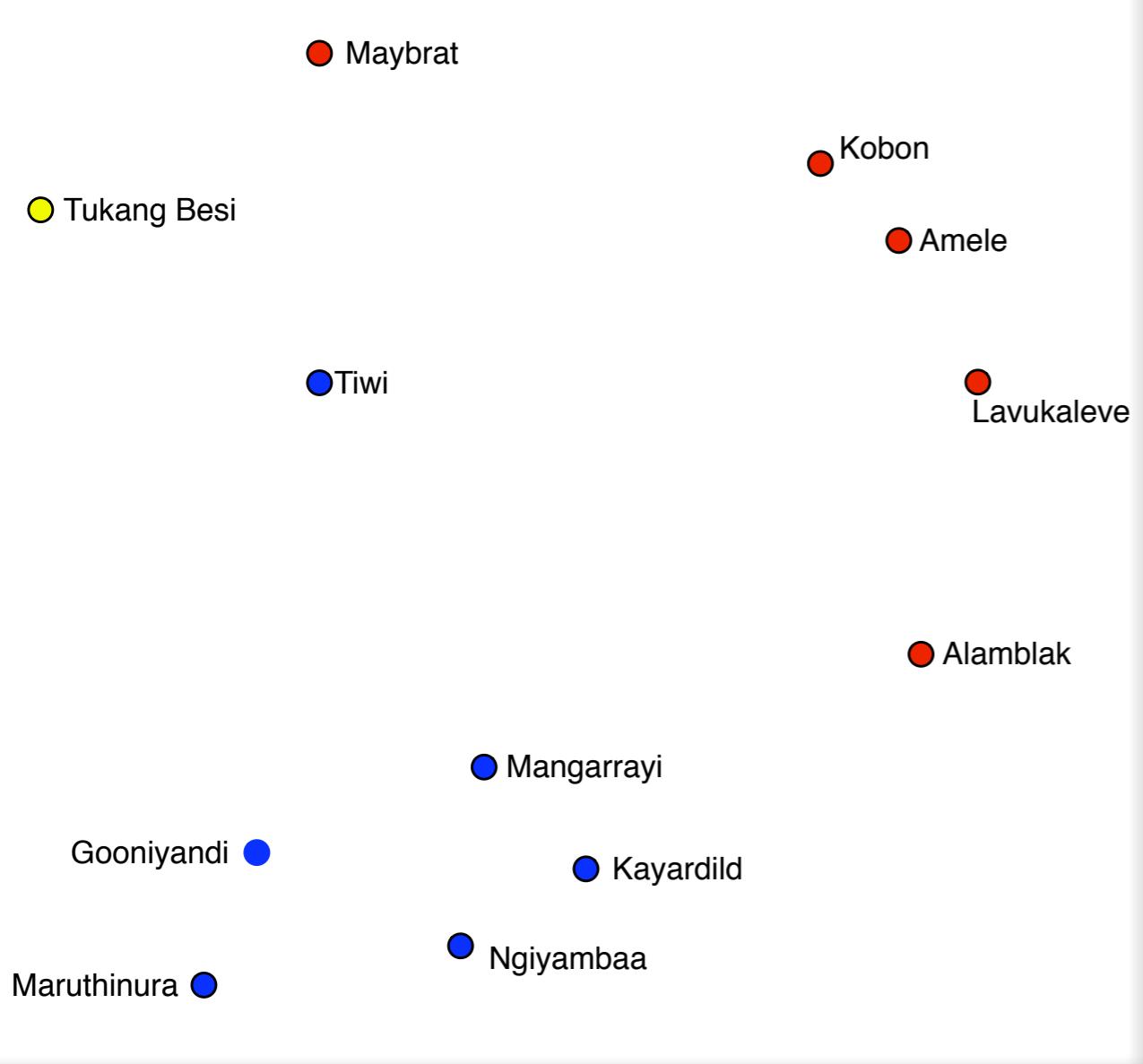


MDS of typological distances

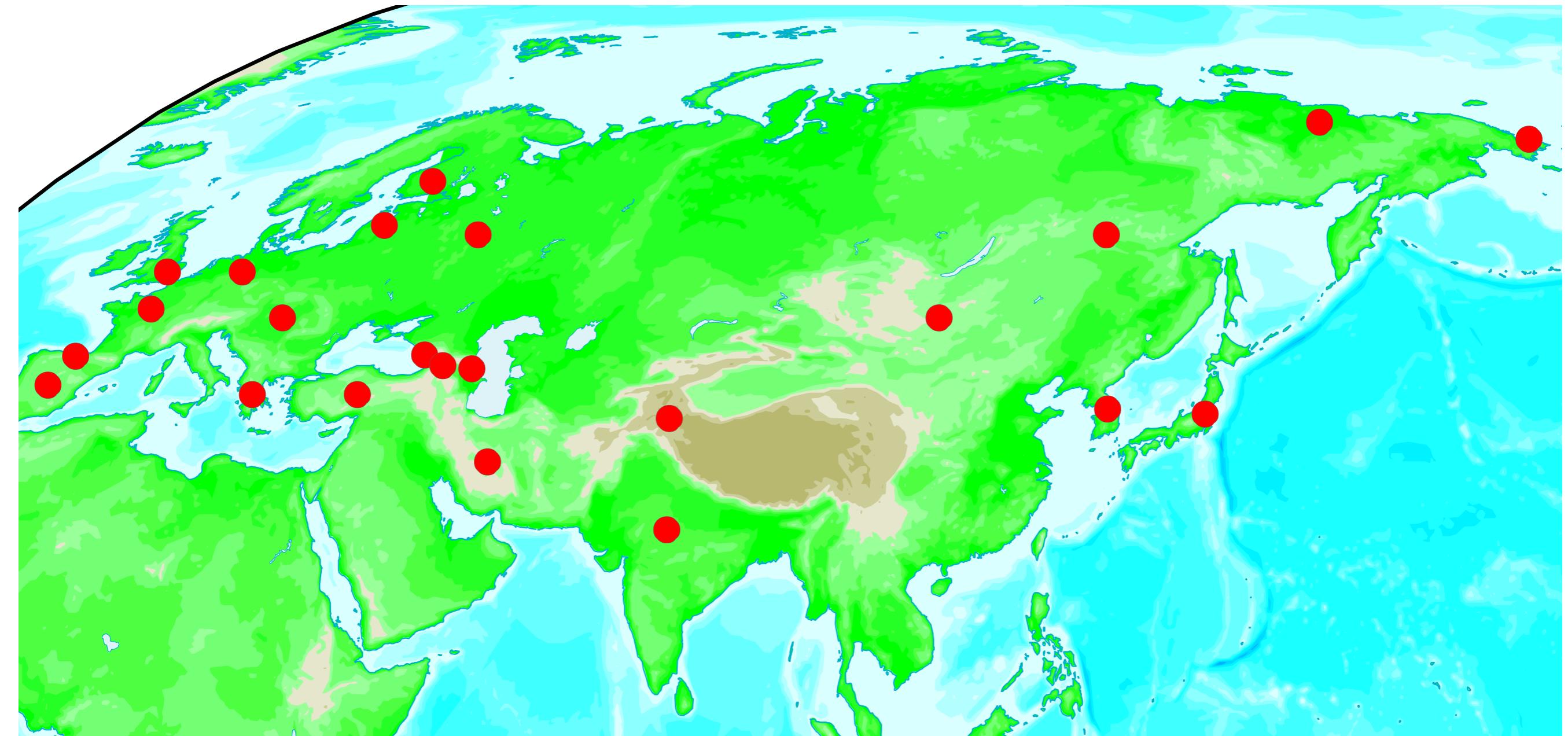


MDS of typological distances

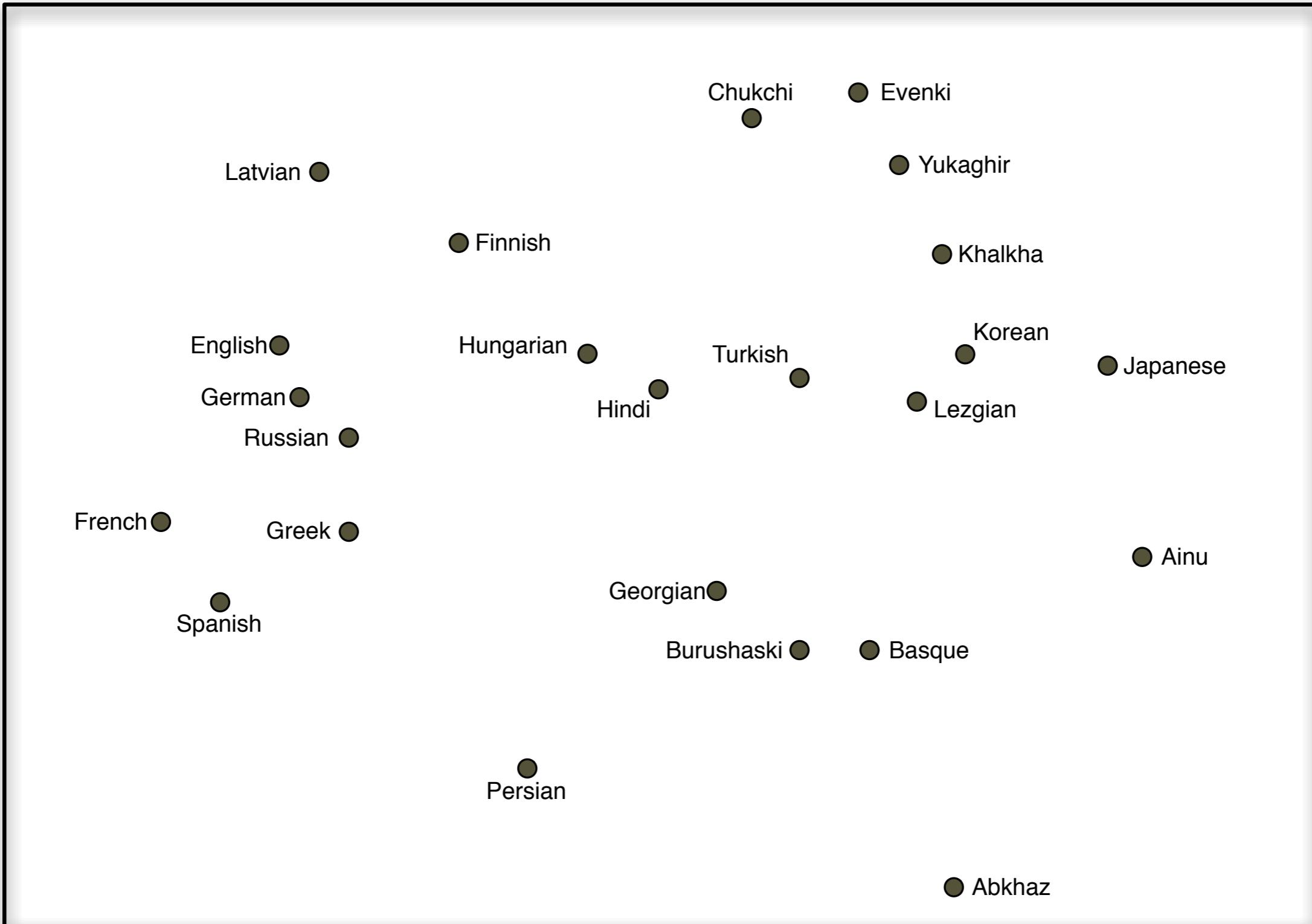




A closer look at geography: the case of Eurasia



MDS of typological distances



WALS forever !