

Towards phylogenetic analyses of linguistic data

Michael Cysouw

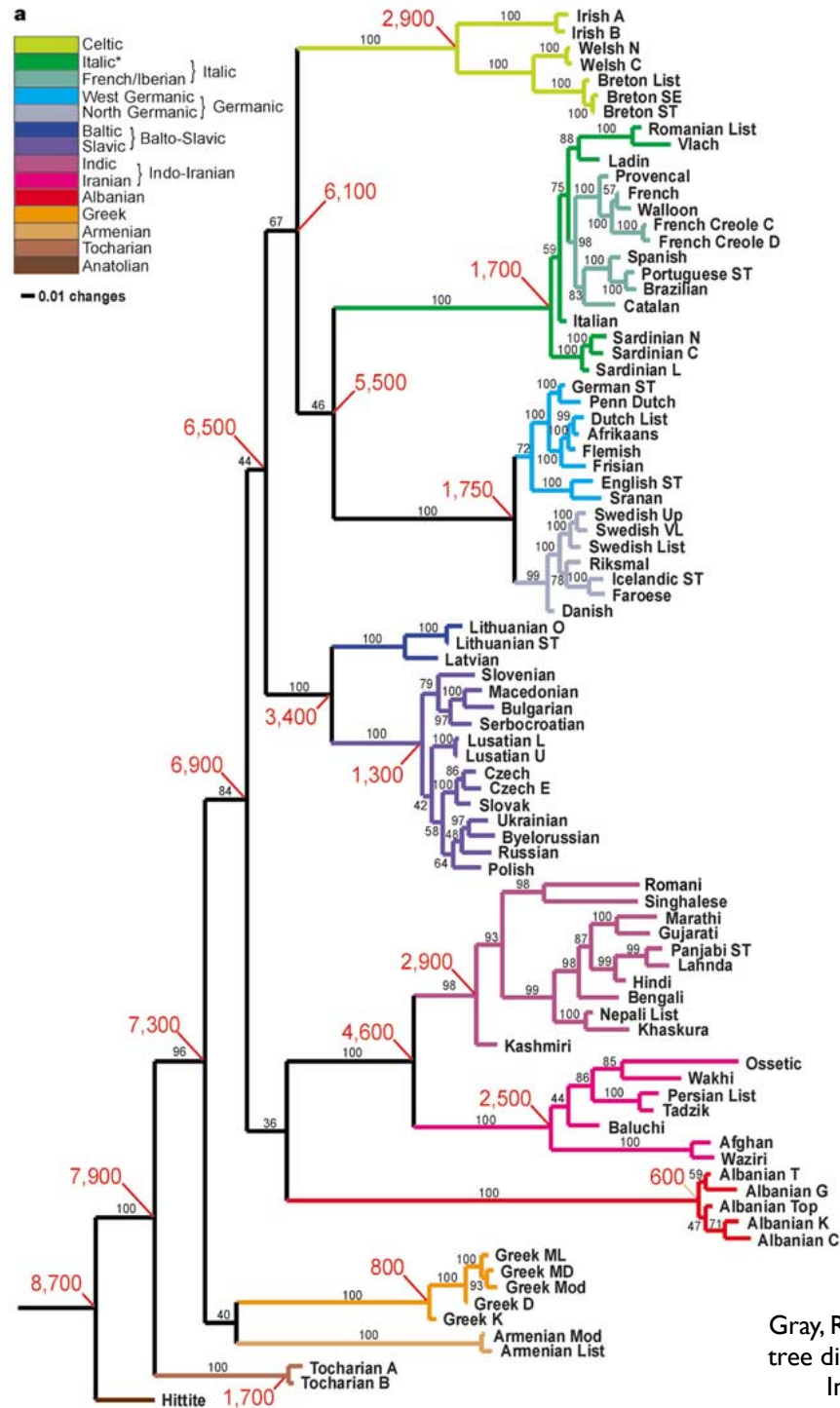


MAX-PLANCK-GESELLSCHAFT

Max Planck Institute
for Evolutionary Anthropology

I. Phylogenetic Analyses in Linguistics

Gothic	<i>fotus</i>	'foot'
Old Icelandic	<i>fo:tr</i>	'foot'
Old High German	<i>fuoz</i>	'foot'
German	<i>Fuss</i>	'foot'
Old English	<i>fo:t</i>	'foot'
English	<i>foot</i>	'foot'
Hittite	<i>pata-</i>	'foot'
Luvian	<i>pati-</i>	'foot'
Latin	<i>pe:s, pedis</i>	'foot'
Tocharian A	<i>päts</i>	'foot'
Tocharian B	<i>ptsa</i>	'foot'
Greek	<i>poús, podós</i>	'foot'
Armenian	<i>ot-n, ot-k'</i>	'foot, feet'
Sanskrit	<i>pá:t, pá:dam</i>	'foot'
Avestan	<i>pad-</i>	'foot'
Lithuanian	<i>pa~das</i>	'shoe'
Latvian	<i>acu-pedius</i>	'swift-footed'
Old Church Slavic	<i>podš</i>	'ground'
Proto-Indo-European	<i>*p^het'-</i>	



Gray, Russel & Quentin D. Atkinson (2003). Language-tree divergence times support the Anatolian theory of Indo-European origin. *Nature* 426: 435-439.

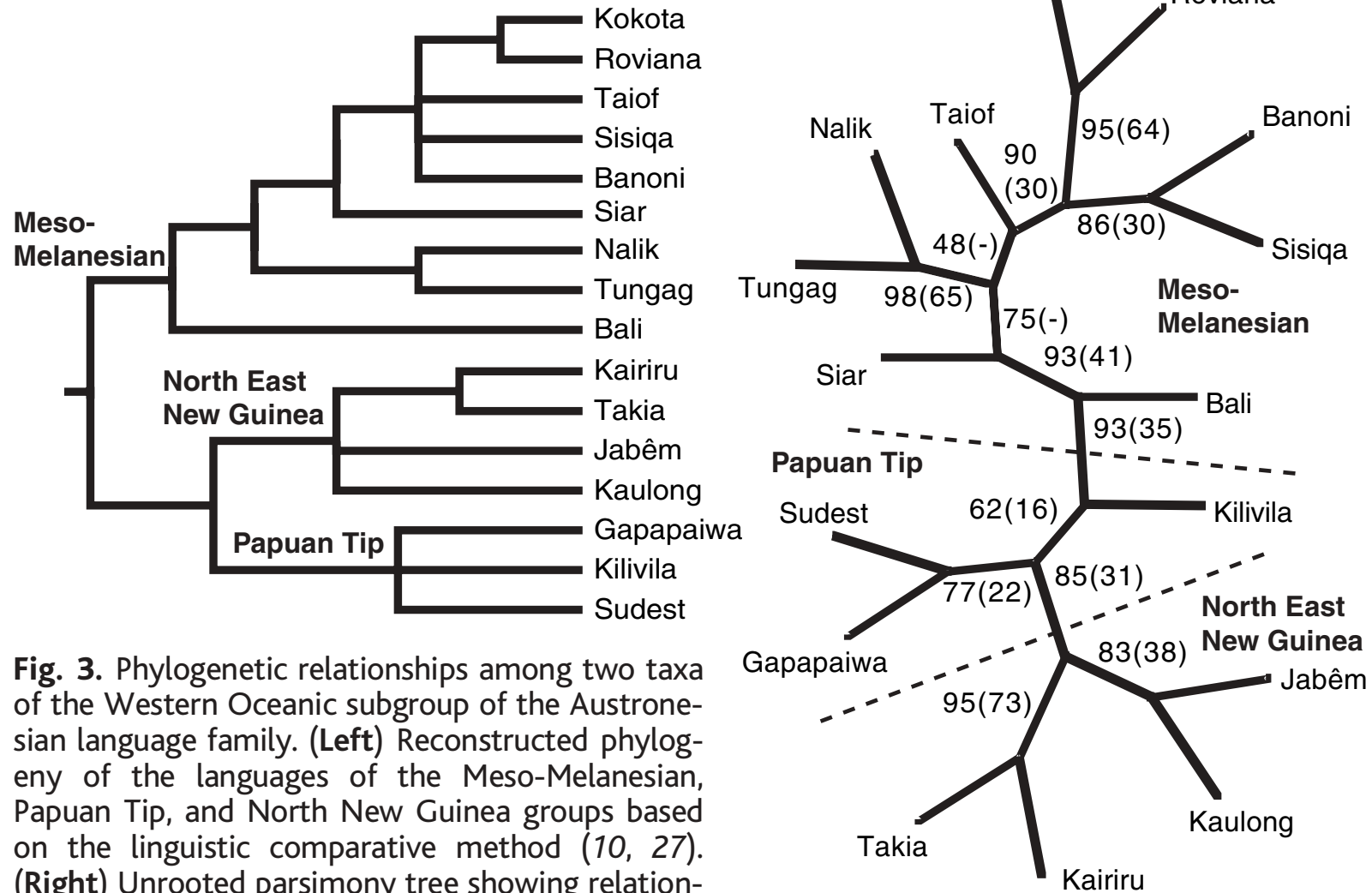


Fig. 3. Phylogenetic relationships among two taxa of the Western Oceanic subgroup of the Austronesian language family. (Left) Reconstructed phylogeny of the languages of the Meso-Melanesian, Papuan Tip, and North New Guinea groups based on the linguistic comparative method (10, 27). (Right) Unrooted parsimony tree showing relationships among the Meso-Melanesian and Papuan Tip groups based on grammatical traits only (that is, discarding abundant lexical evidence) (the figure shows reweighted and raw bootstrap values). The two trees show a high degree of concordance, with monophyly in both major taxa and the similar geographical structuring of within-taxon diversity.

2. Introducing Linguistic Typology



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

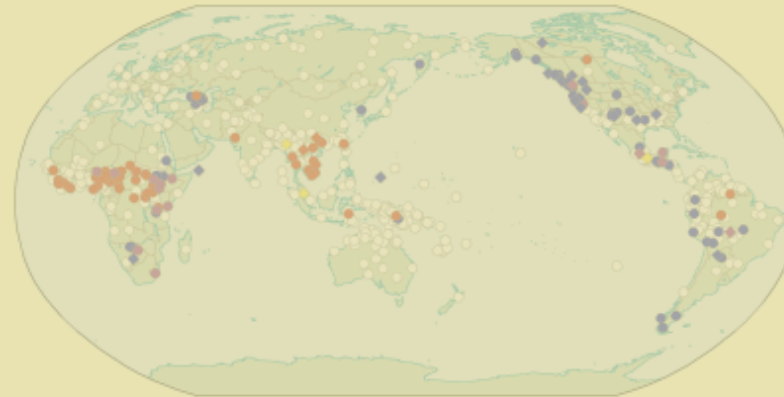
FEATURE VIEWER

LANGUAGE VIEWER

COMPOSER

ABOUT THIS PROGRAM

GUIDED TOUR



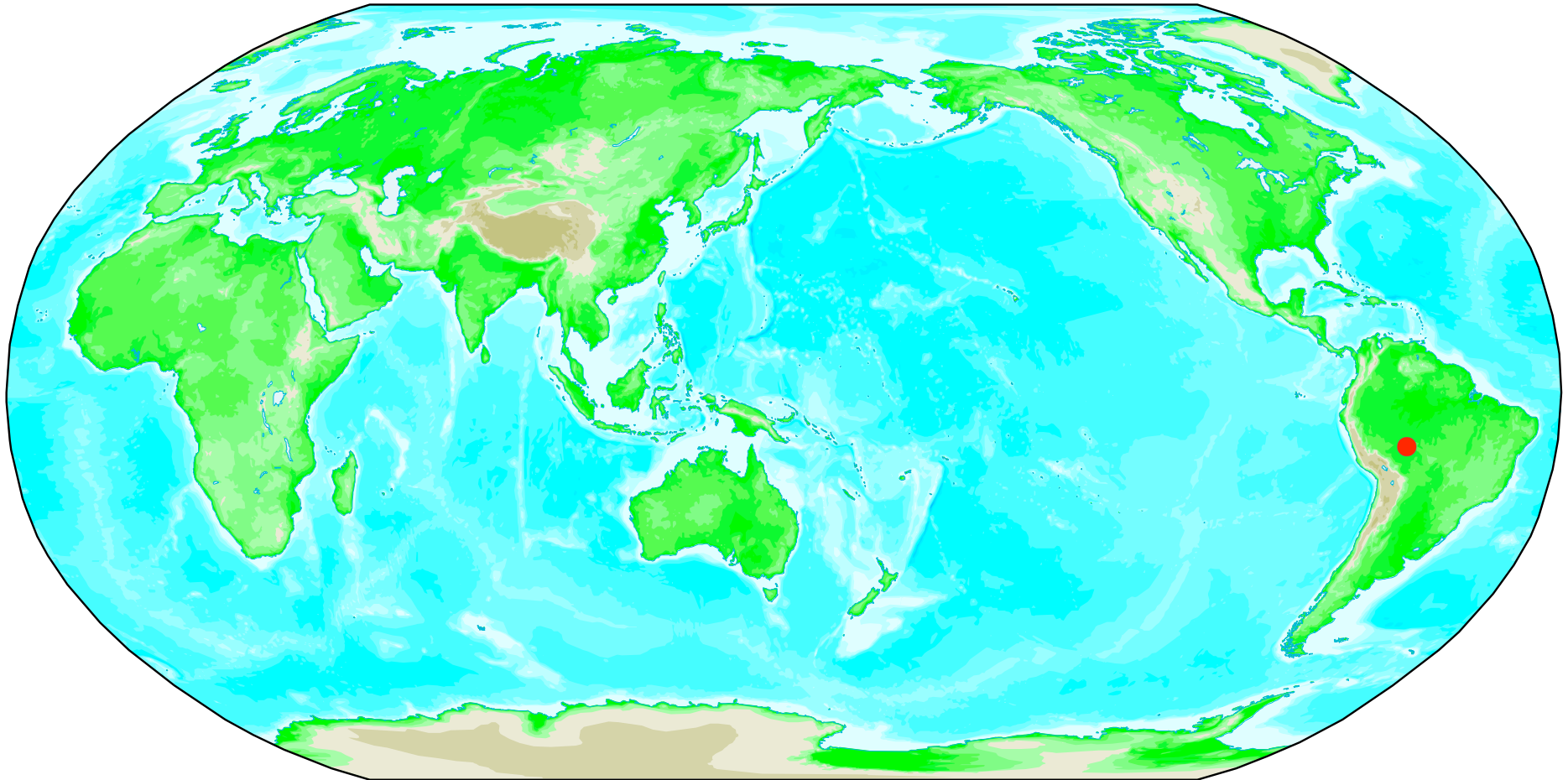
3. Investigating unusual characteristics

And the winners are:

In the category:

**‘Most Unusual
Individual Language’**

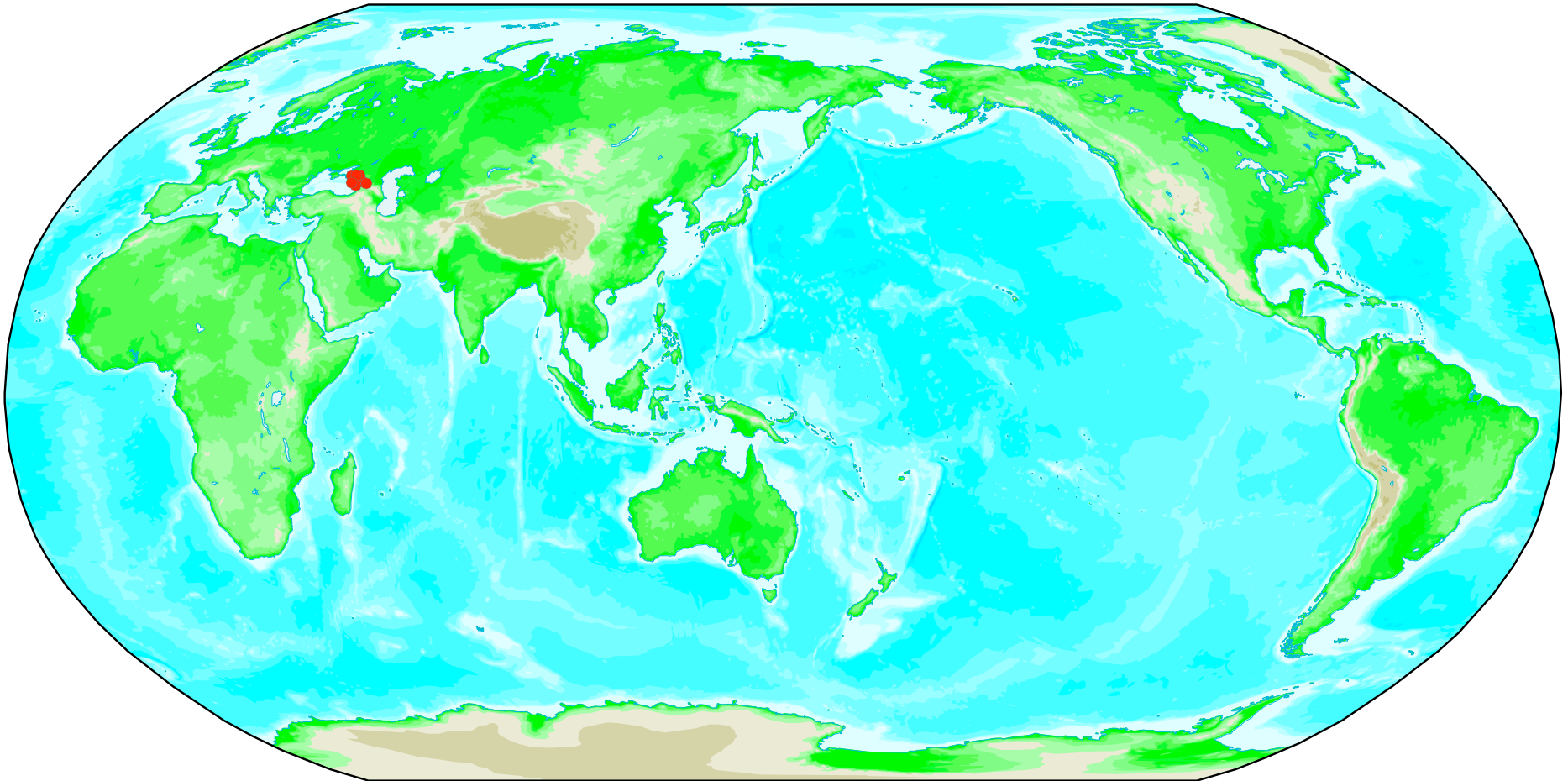
Wari'



In the category:

**‘Most Unusual
Genealogical Group’**

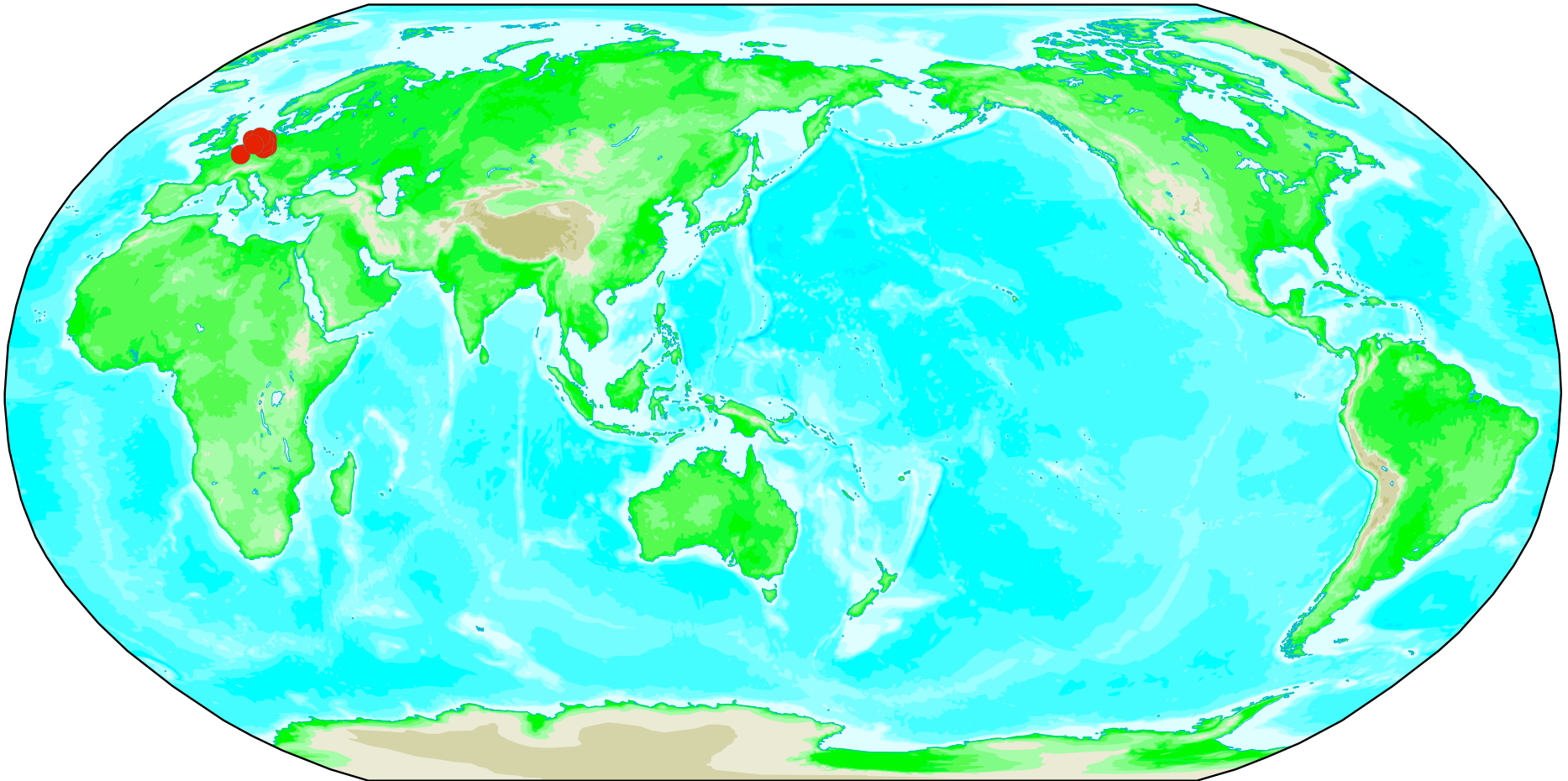
Northwest Caucasian



In the category:

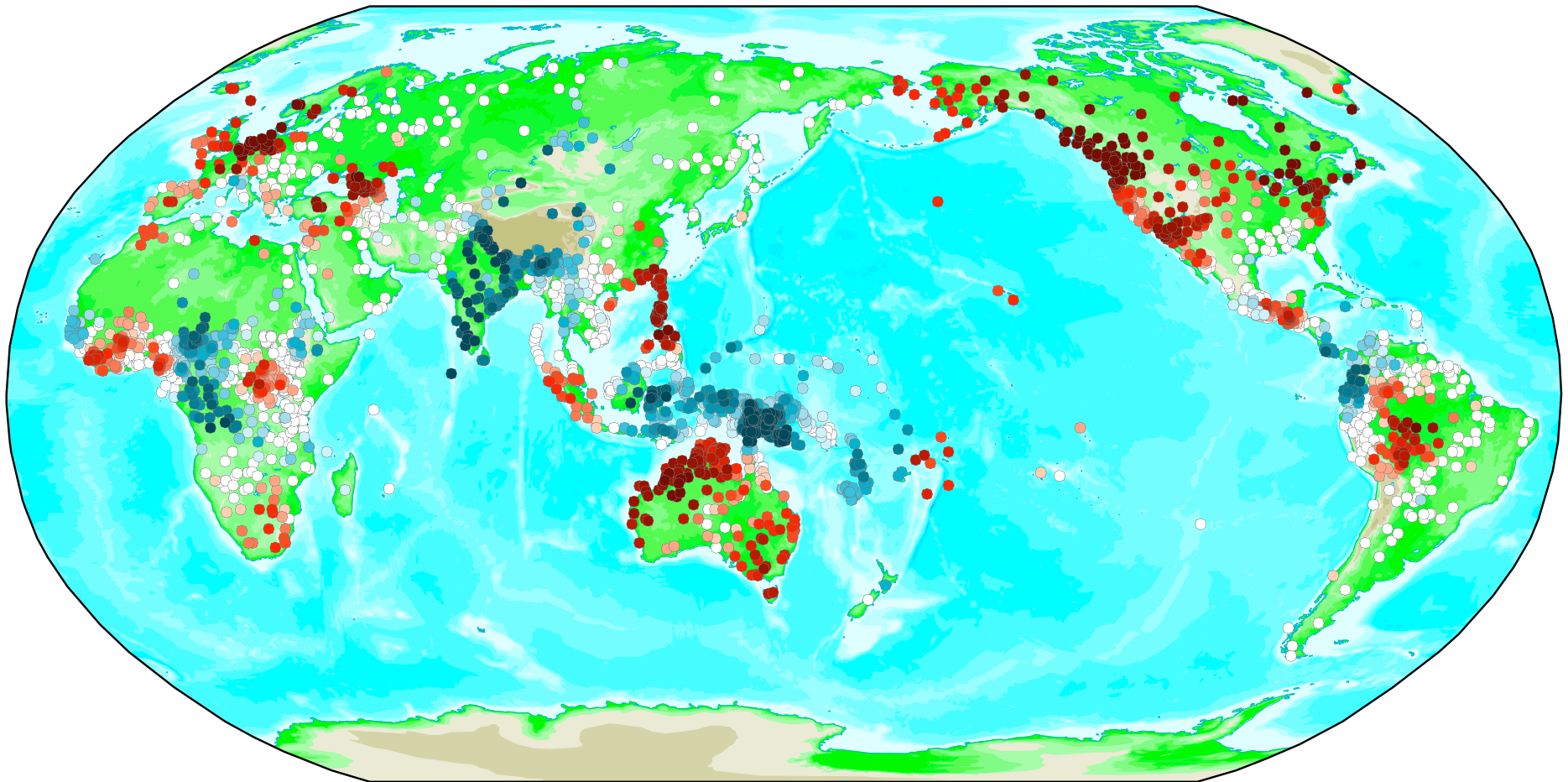
**‘Most Unusual
Geographical Area’**

Northwest Continental Europe



All languages

(red = rare, blue = common)

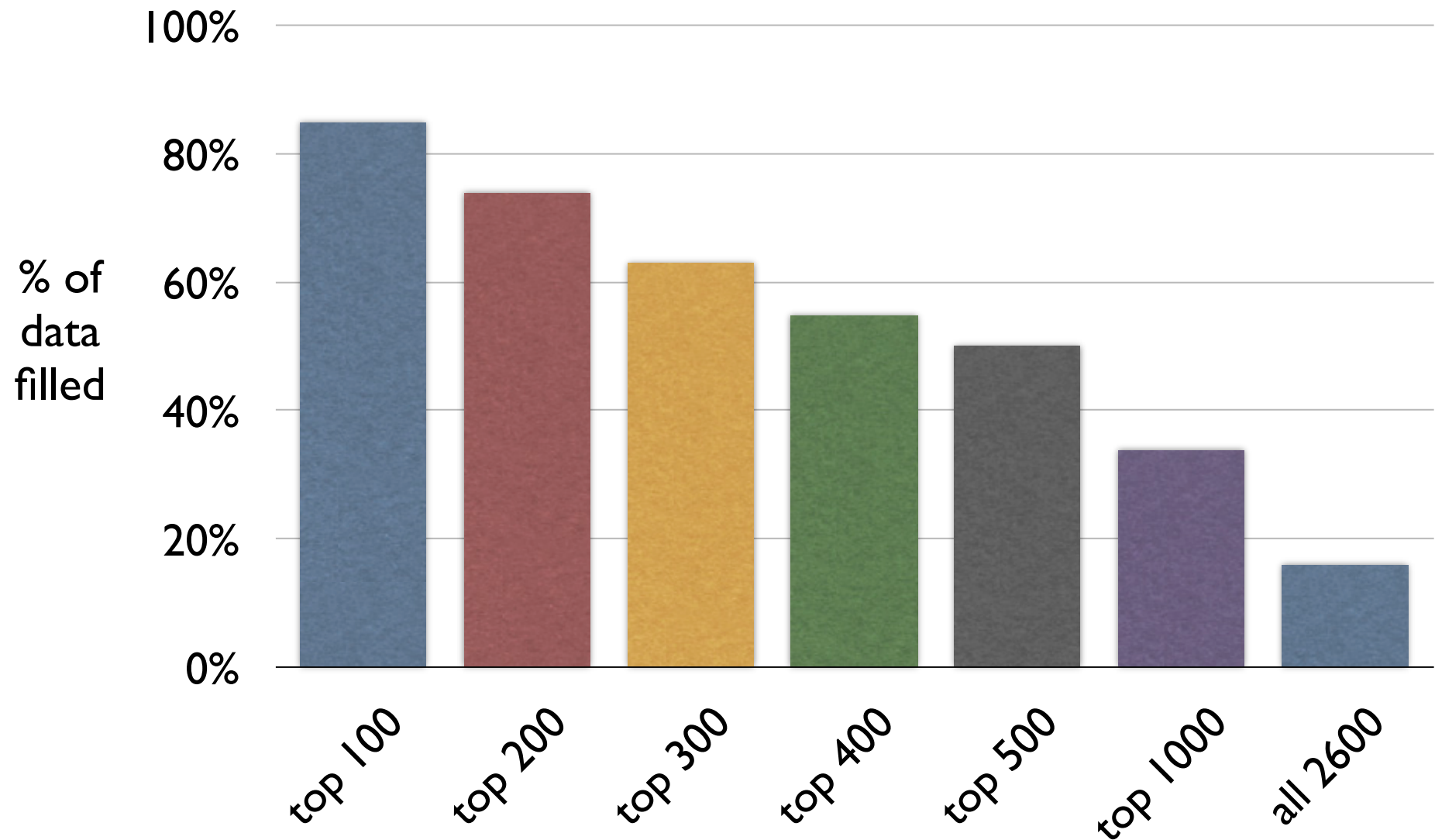


4. Towards Phylogenetic Analyses with WALS

A wealth of data (for linguistics)

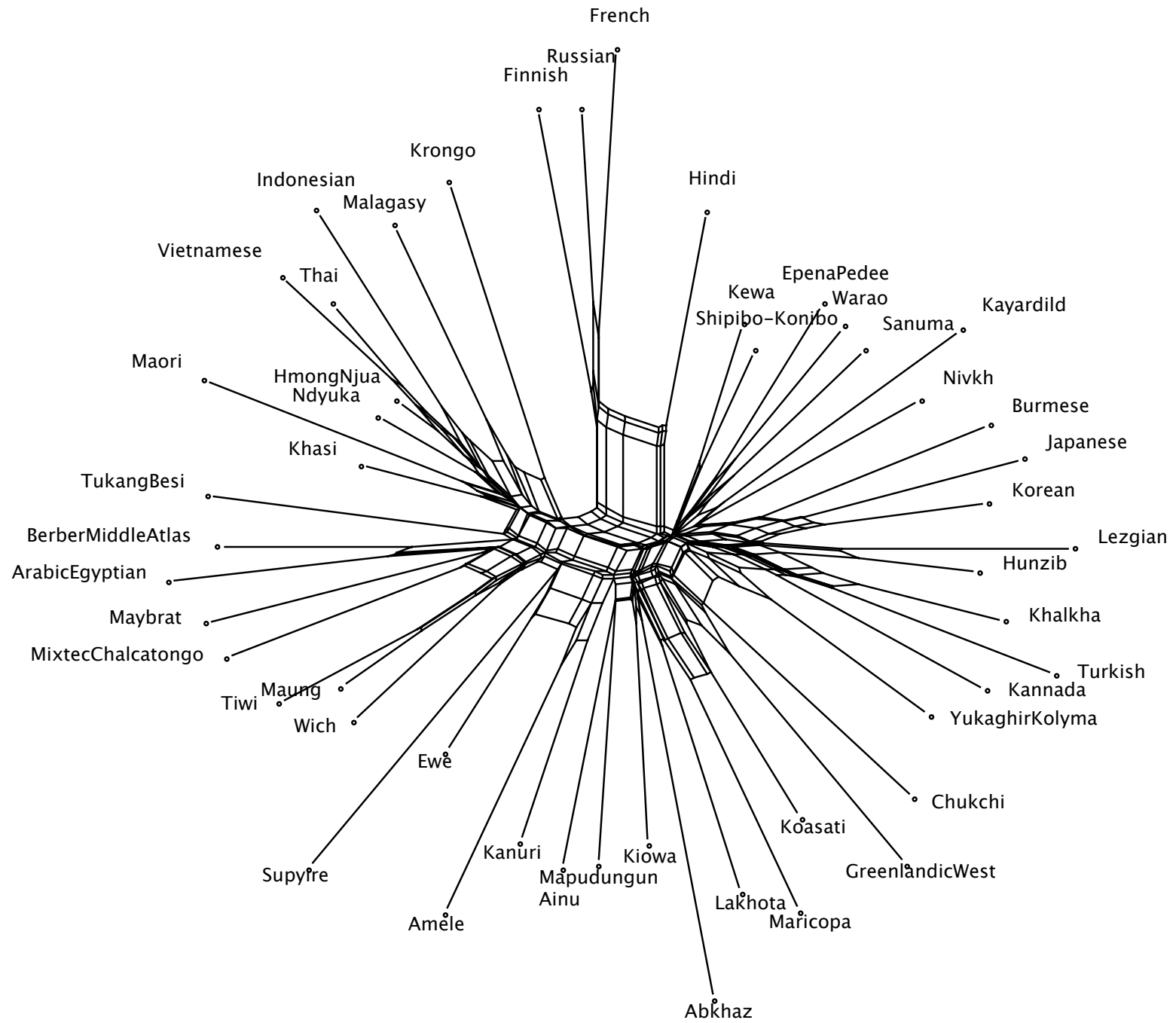
- 2,600 languages
- 140 characteristics
- almost 60,000 datapoints
- but: $60,000/2,600*140 = 0.165$
- the datatable is only 16.5 % filled !

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 %)



5. Improving Typology

5. Improving Typology

- Finer-grained Coding
- 'Deconstructing' Typology
- Select Suitable Characteristics

WALS the Feature Viewer

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

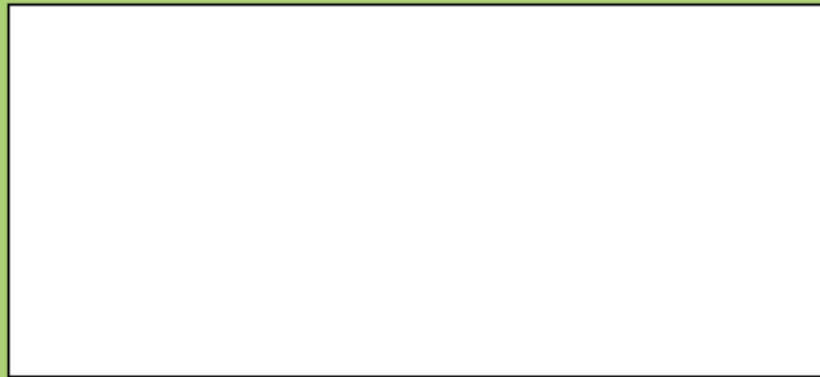
symbol: include: [click to list languages below](#) [no. of lgs : of genera : of families]

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

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

<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

DESCRIPTION

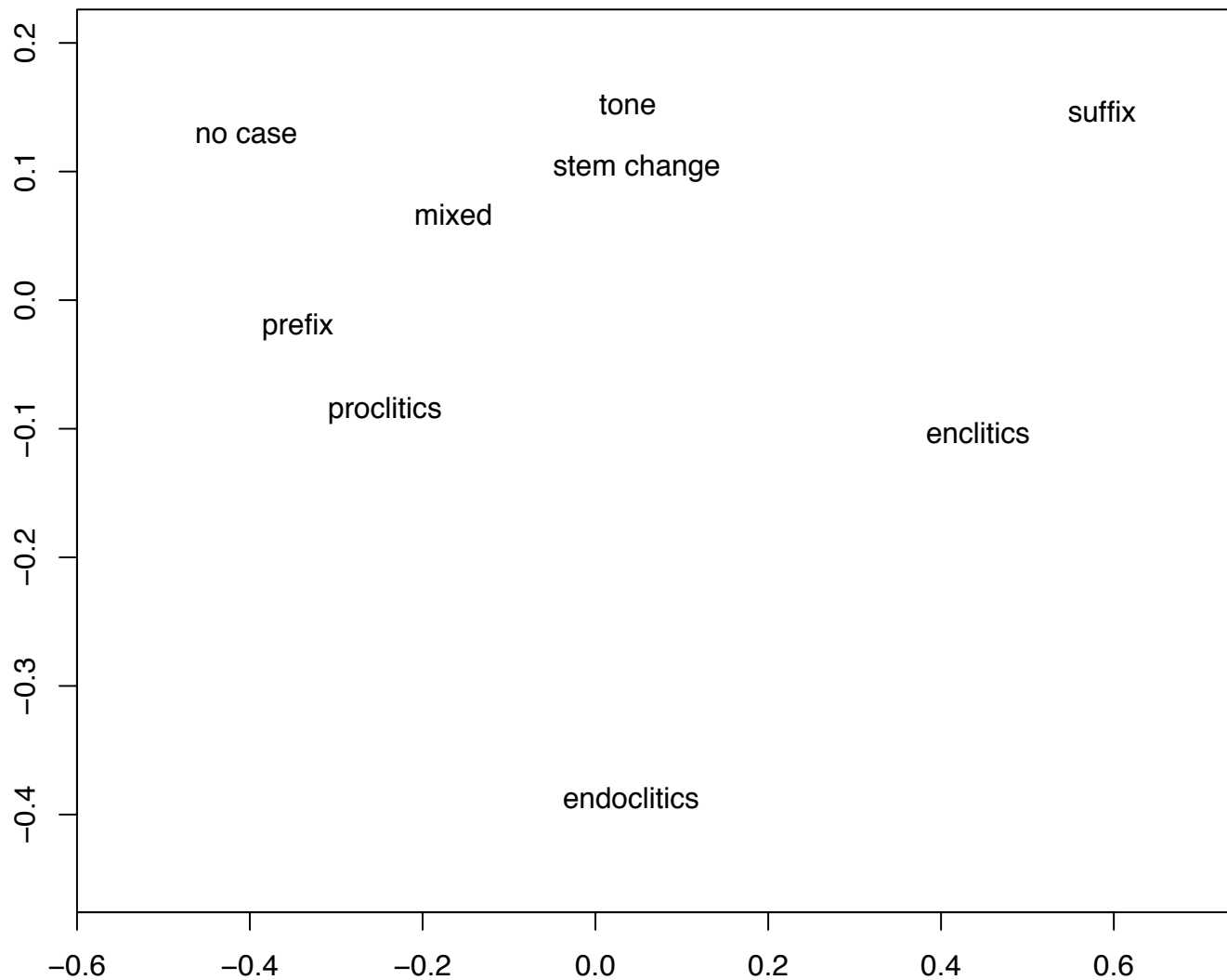


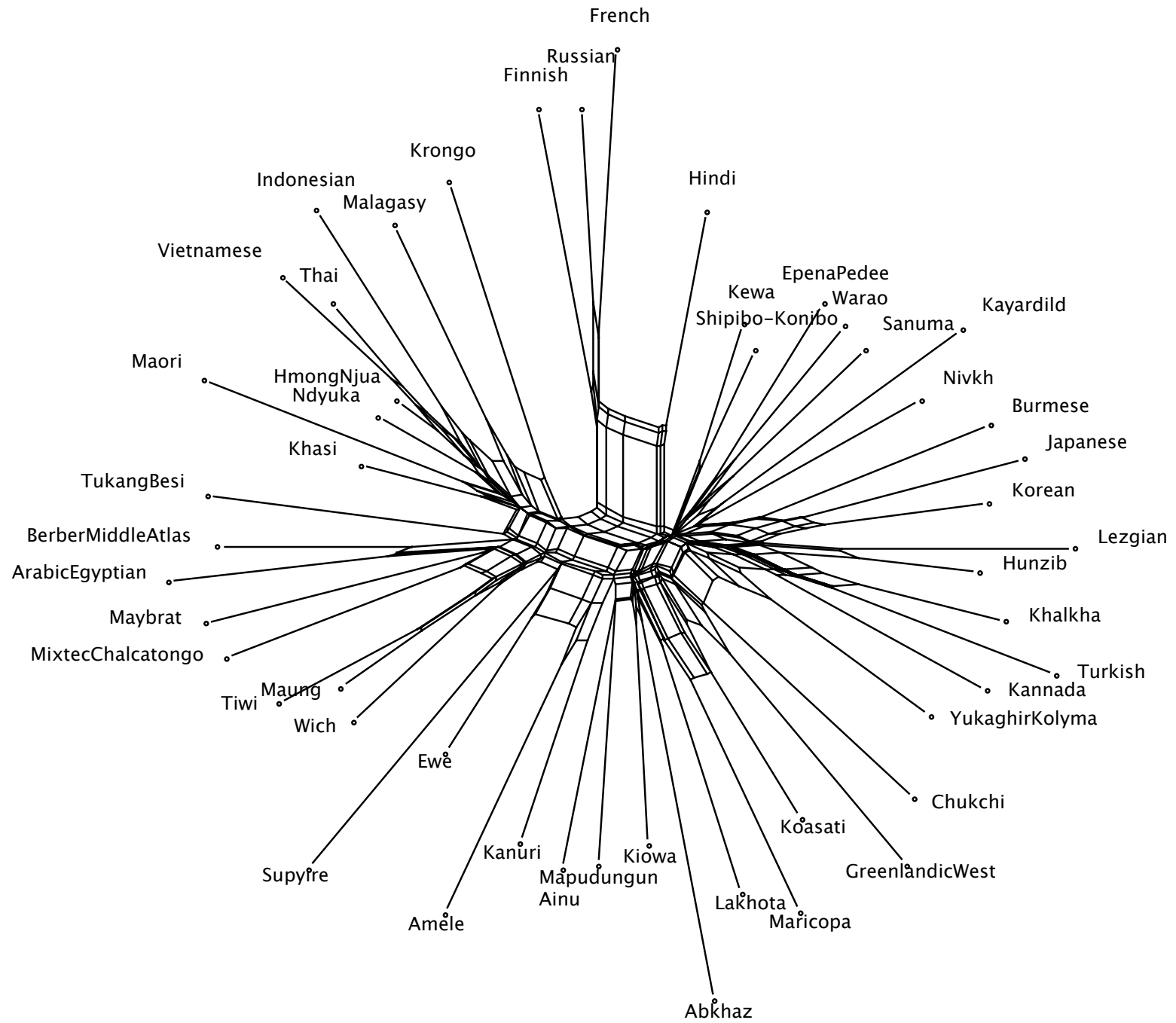
arrange the languages by

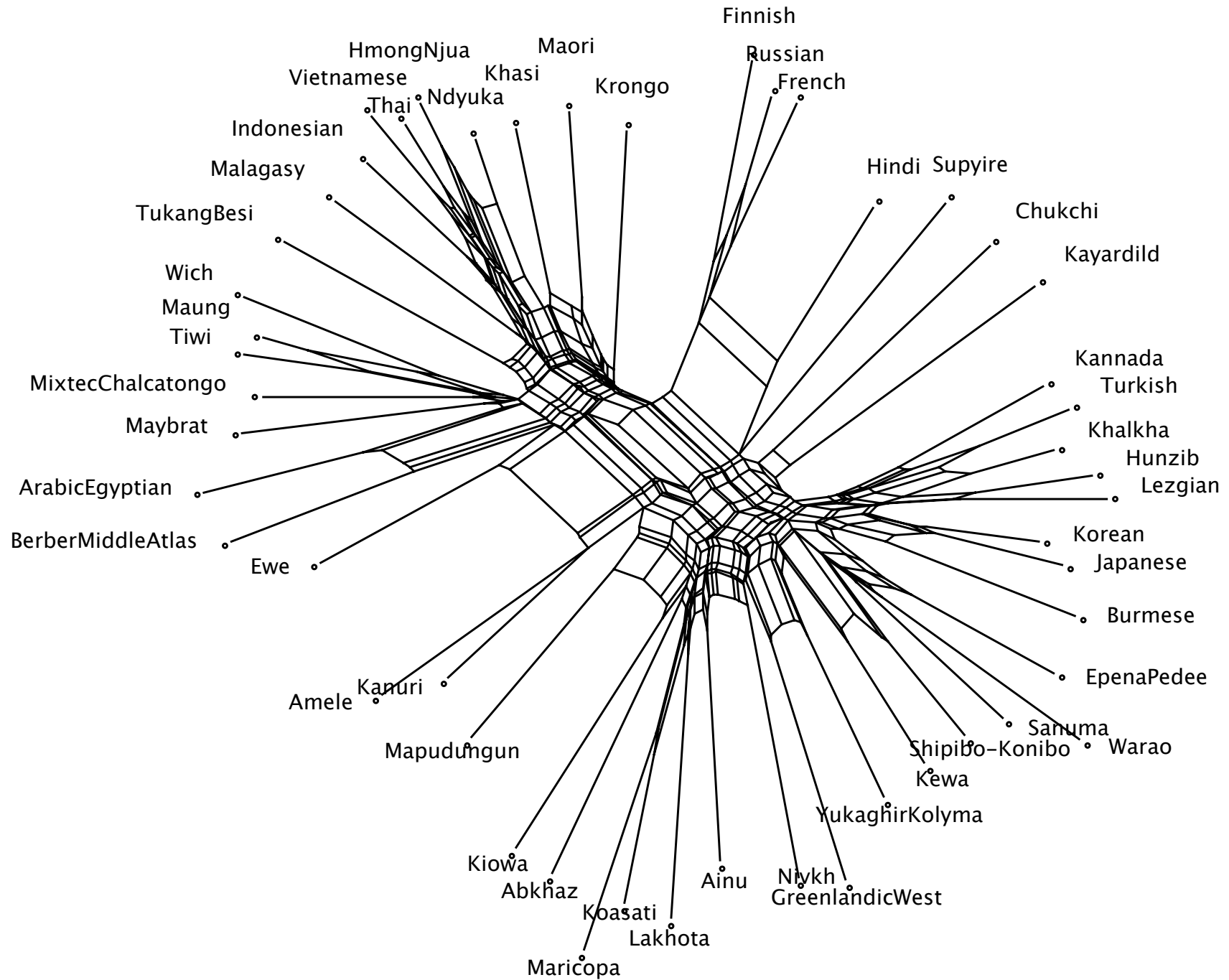
languages

COPYLIST

Estimate **character-similarities** from cooccurrence in low-level genetic groups







5. Improving Typology

- Finer-grained Coding
- **'Deconstructing' Typology**
- Select Suitable Characteristics

	L ₁	L ₂	L ₃	L ₄	L ₅	L ₆	L ₇	L ₈	...
L ₁	1	1	1	0	0	0	0	0	
L ₂	1	1	1	0	0	0	0	0	
L ₃	1	1	1	0	0	0	0	0	
L ₄	0	0	0	1	1	0	0	0	
L ₅	0	0	0	1	1	0	0	0	
L ₆	0	0	0	0	0	1	1	1	
L ₇	0	0	0	0	0	1	1	1	
L ₈	0	0	0	0	0	1	1	1	
...									

Undifferentiated Typology

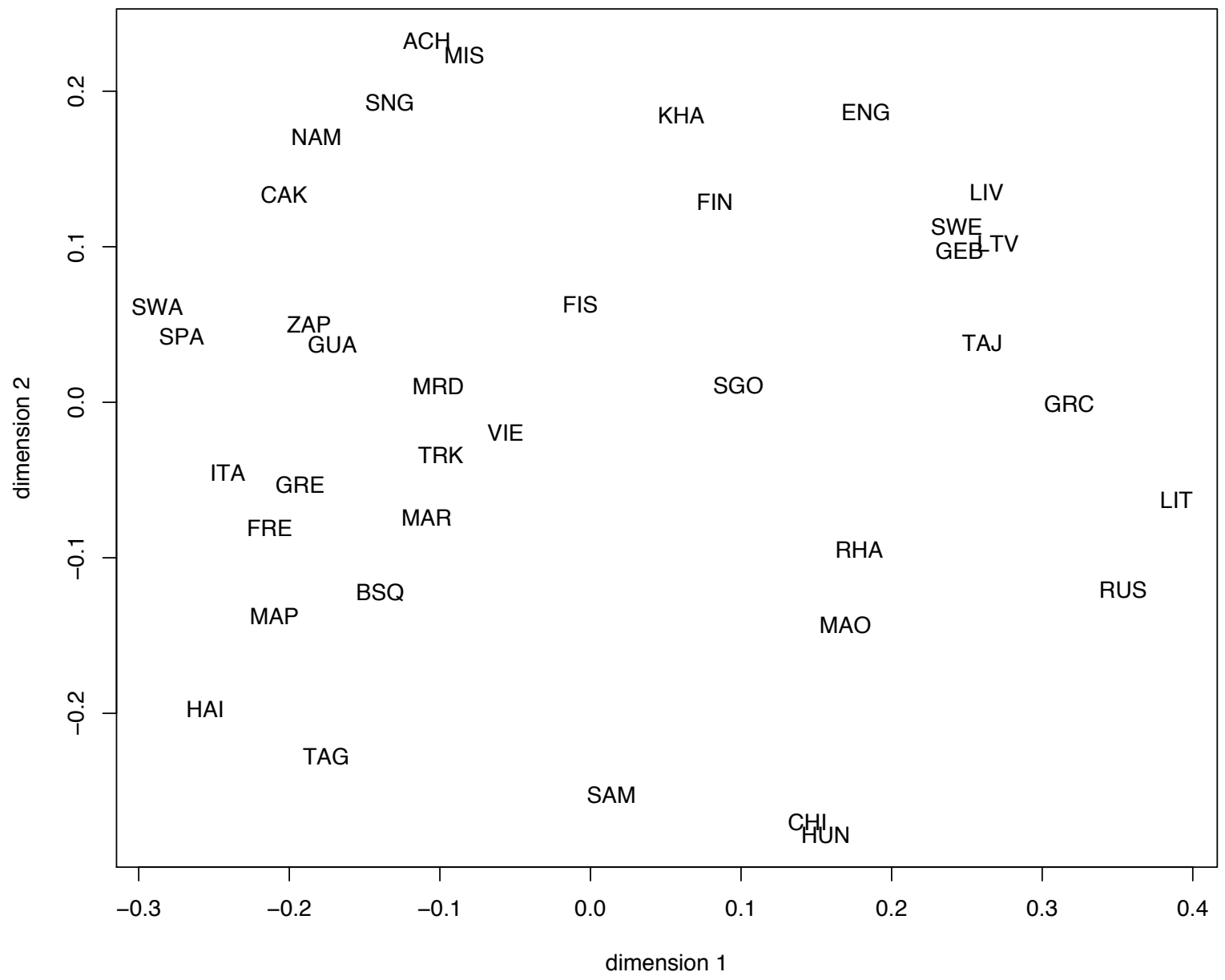
	L ₁	L ₂	L ₃	L ₄	L ₅	L ₆	L ₇	L ₈	...
L ₁	1	1	1	0.37	0.37	0.28	0.28	0.28	
L ₂	1	1	1	0.37	0.37	0.28	0.28	0.28	
L ₃	1	1	1	0.37	0.37	0.28	0.28	0.28	
L ₄	0.37	0.37	0.37	1	1	0.58	0.58	0.58	
L ₅	0.37	0.37	0.37	1	1	0.58	0.58	0.58	
L ₆	0.28	0.28	0.28	0.58	0.58	1	1	1	
L ₇	0.28	0.28	0.28	0.58	0.58	1	1	1	
L ₈	0.28	0.28	0.28	0.58	0.58	1	1	1	
...									

Inter-type similarities

	L ₁	L ₂	L ₃	L ₄	L ₅	L ₆	L ₇	L ₈	...
L ₁	1	0.55	0.72	0.31	0.70	0.61	0.50	0.58	
L ₂	0.55	1	0.55	0.31	0.40	0.44	0.31	0.48	
L ₃	0.72	0.55	1	0.29	0.53	0.51	0.48	0.60	
L ₄	0.31	0.31	0.29	1	0.38	0.36	0.26	0.27	
L ₅	0.70	0.40	0.53	0.38	1	0.64	0.51	0.46	
L ₆	0.61	0.44	0.51	0.36	0.64	1	0.57	0.43	
L ₇	0.50	0.31	0.48	0.26	0.51	0.57	1	0.47	
L ₈	0.58	0.48	0.60	0.27	0.46	0.43	0.47	1	
...									

'Deconstructed' Typology

	MRD	LIT	ENG	FRE
1050	sams	eiti	go	aller
1070	sams	eiti	come	venir
1090	sams	eiti	come	venir
1104	lisems	kopti	come	sortir
1105	valgoms	zengti	descend	descendre
1114	–	–	come	se faire entendre
1120	vetjams	varyti	drive	pousser
1140	sams	eiti	come	se rendre
1160	jutams	eiti	walk	marcher



5. Improving Typology

- Finer-grained Coding
- 'Deconstructing' Typology
- **Select Suitable Characteristics**

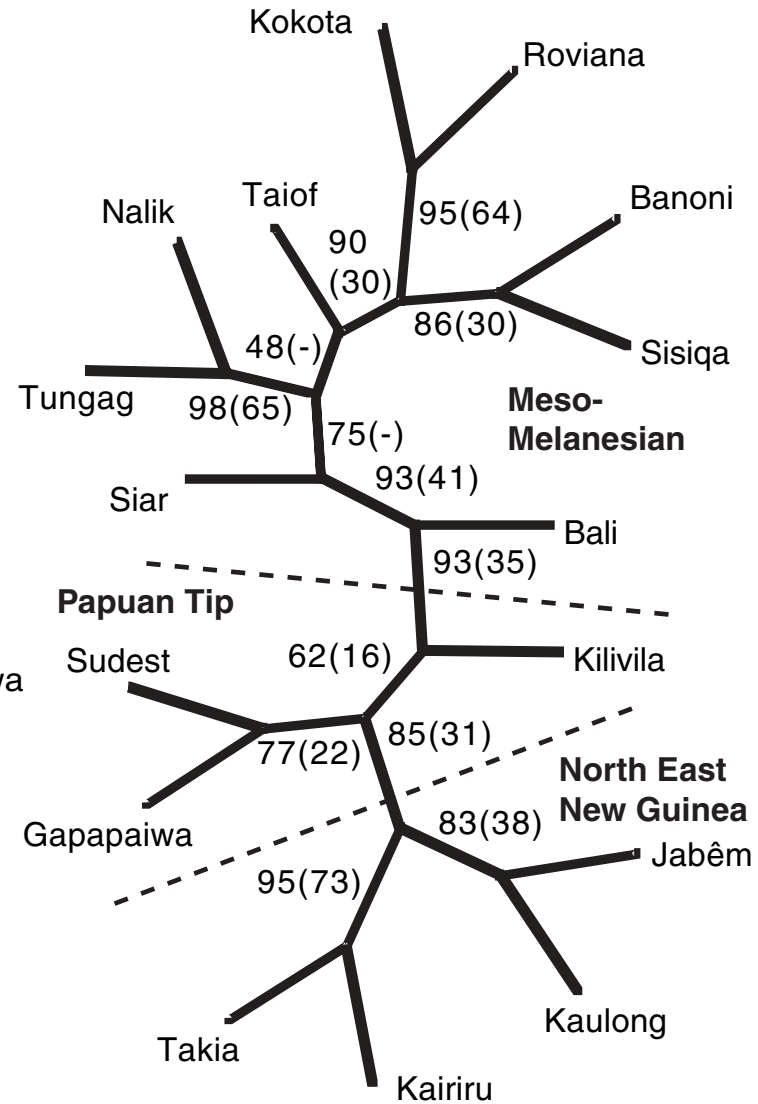
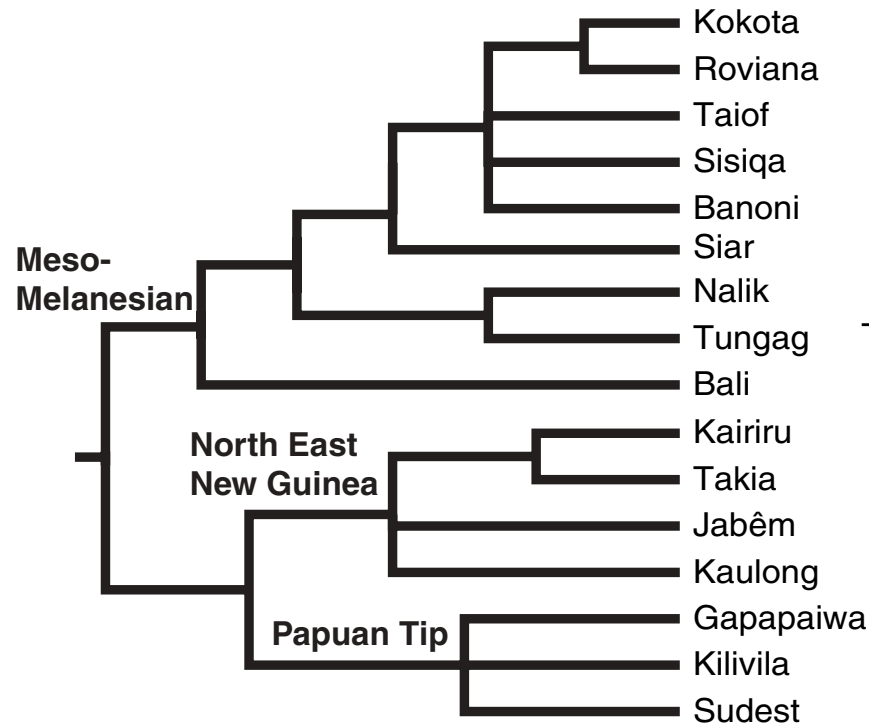
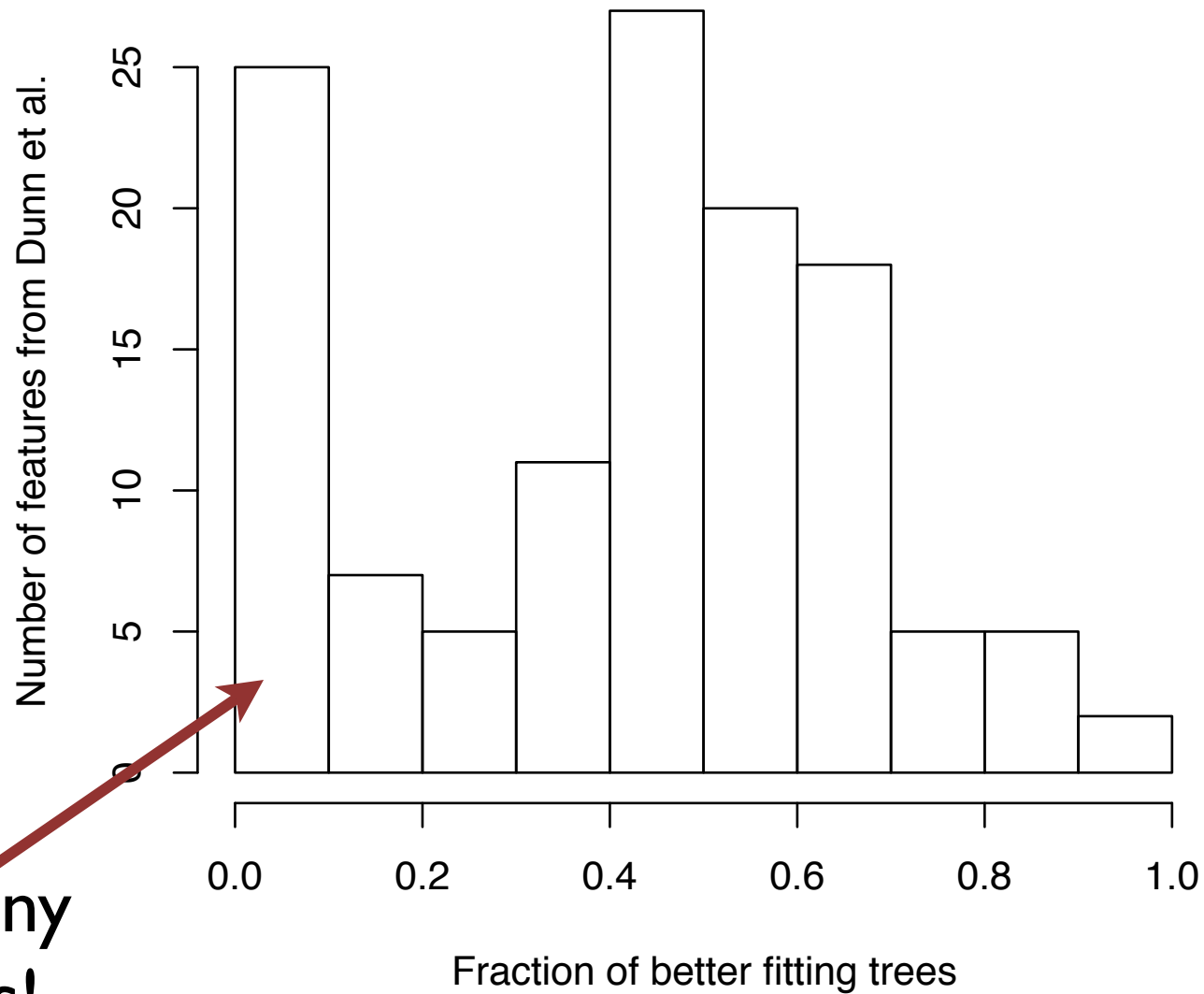


Fig. 3. Phylogenetic relationships among two taxa of the Western Oceanic subgroup of the Austronesian language family. **(Left)** Reconstructed phylogeny of the languages of the Meso-Melanesian, Papuan Tip, and North New Guinea groups based on the linguistic comparative method (10, 27). **(Right)** Unrooted parsimony tree showing relationships among the Meso-Melanesian and Papuan Tip groups based on grammatical traits only (that is, discarding abundant lexical evidence) (the figure shows reweighted and raw bootstrap values). The two trees show a high degree of concordance, with monophyly in both major taxa and the similar geographical structuring of within-taxon diversity.

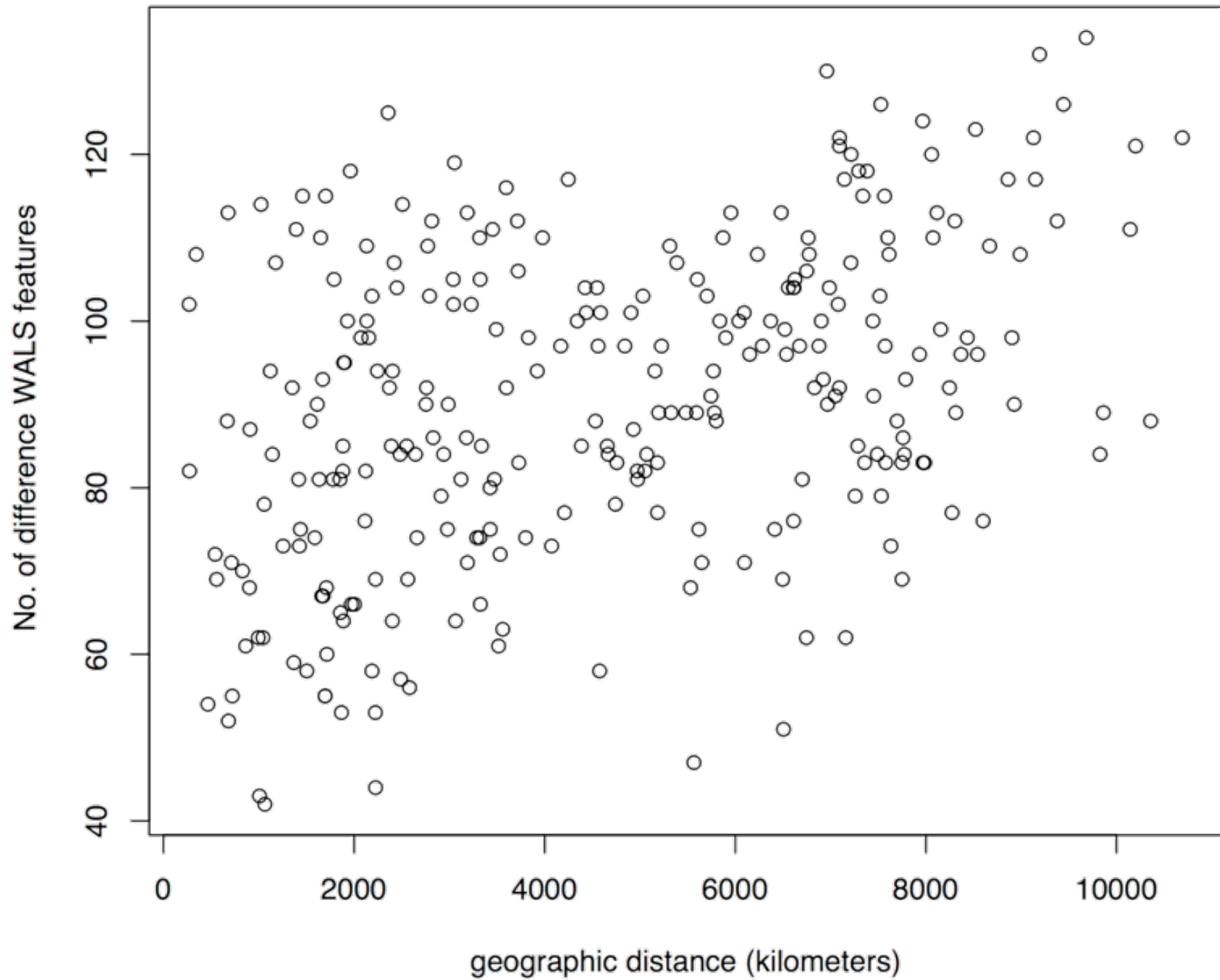
Fit for all 125 features

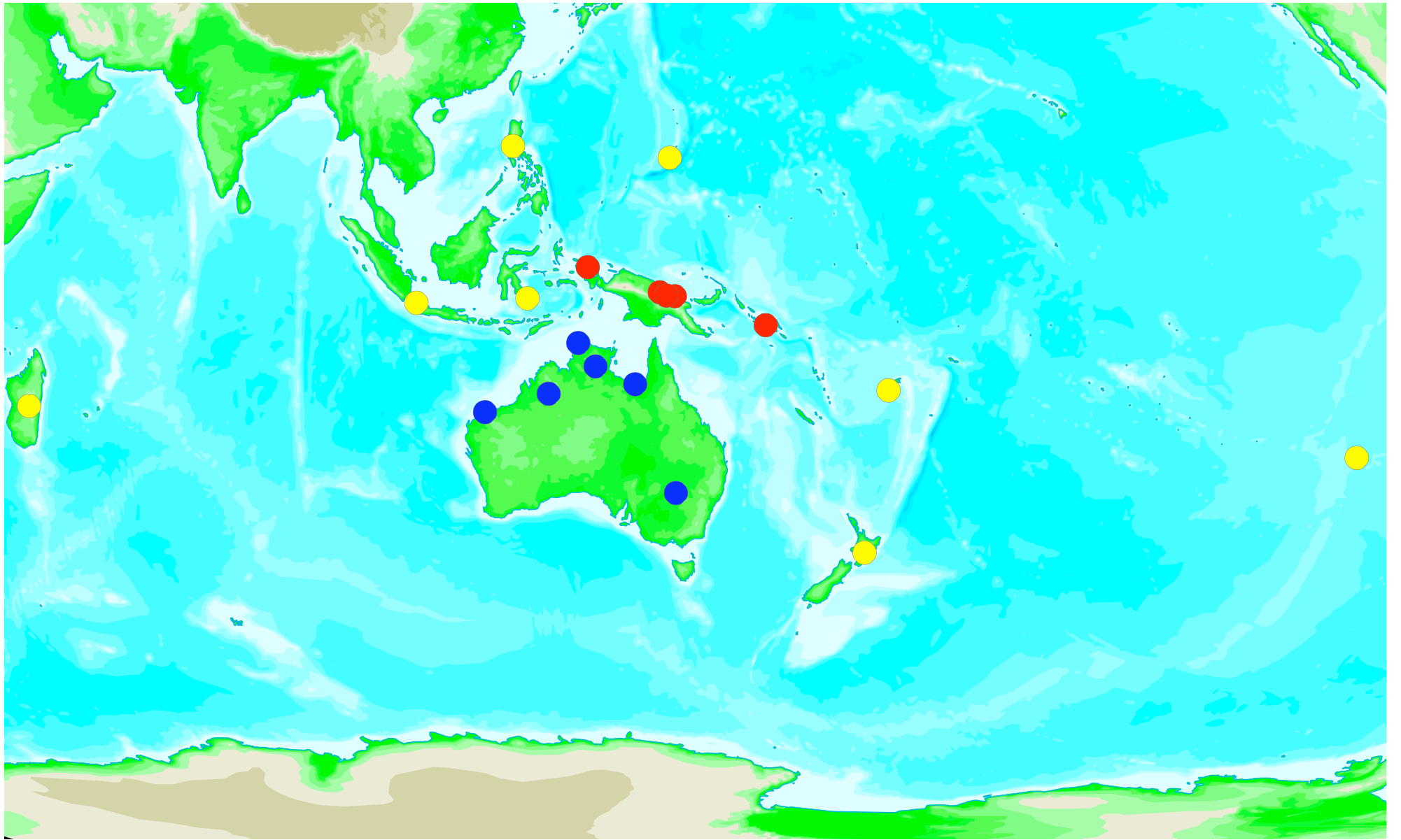


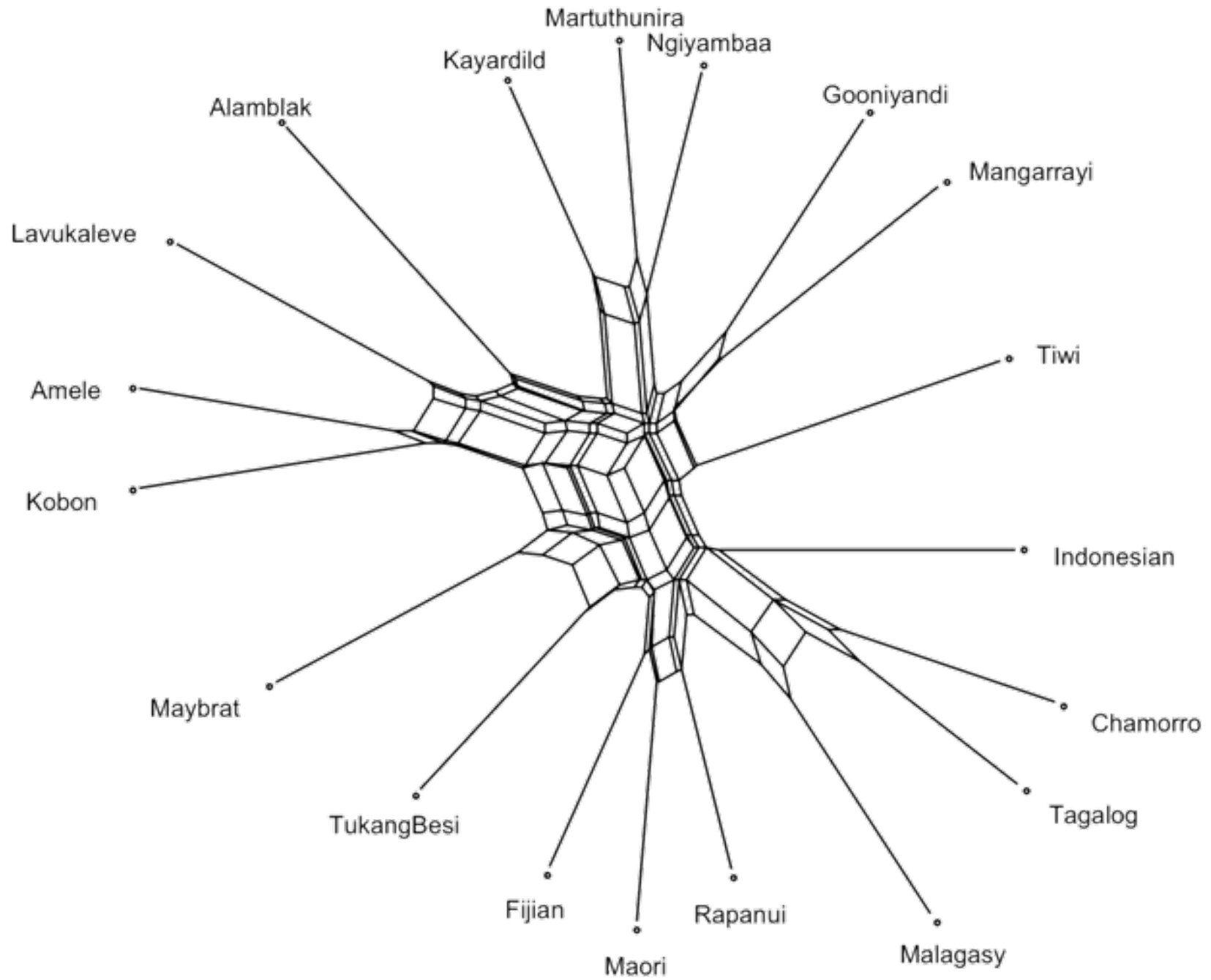
(Too) many
good fits!

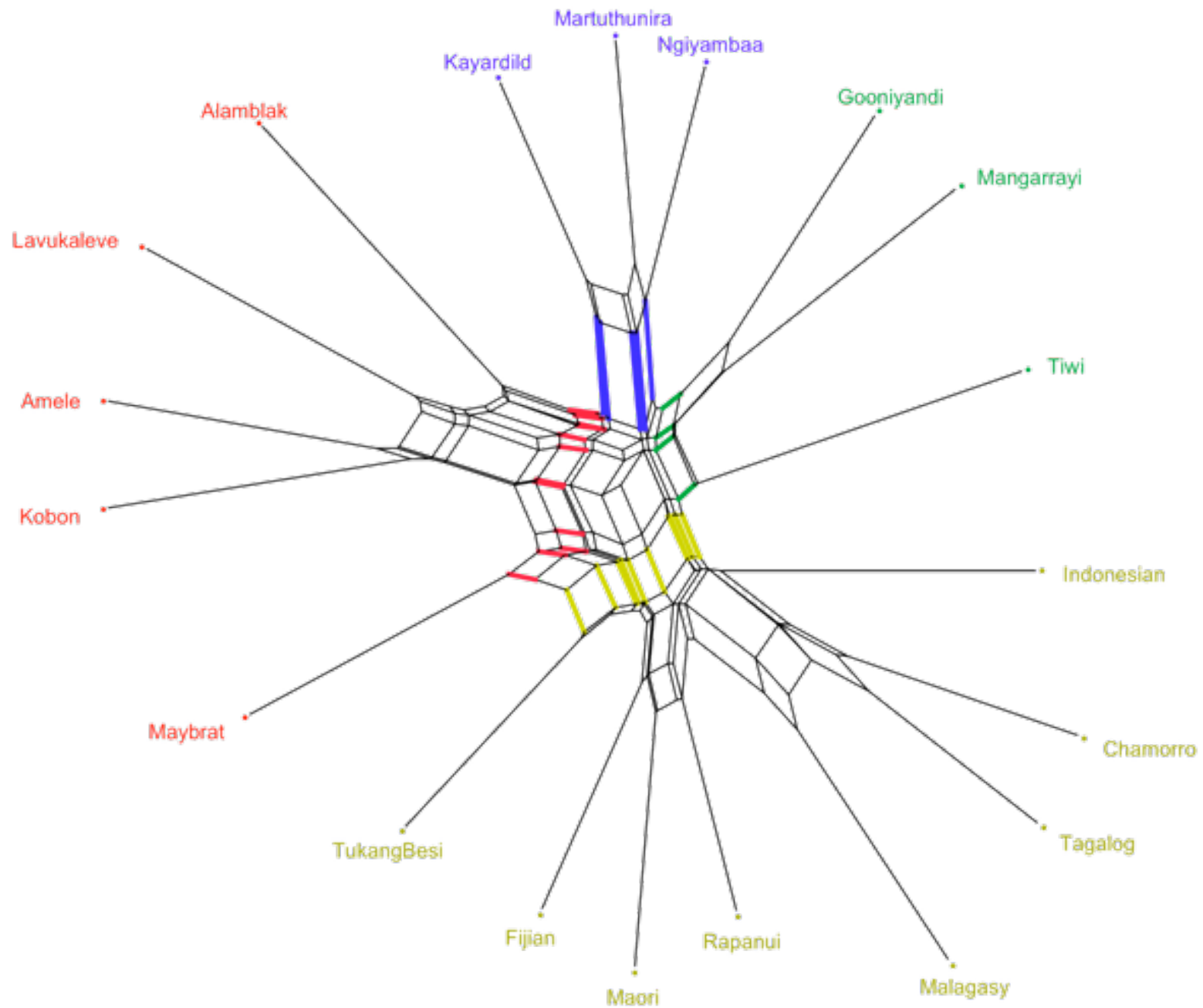
6. The Problem of Geographical Similarities



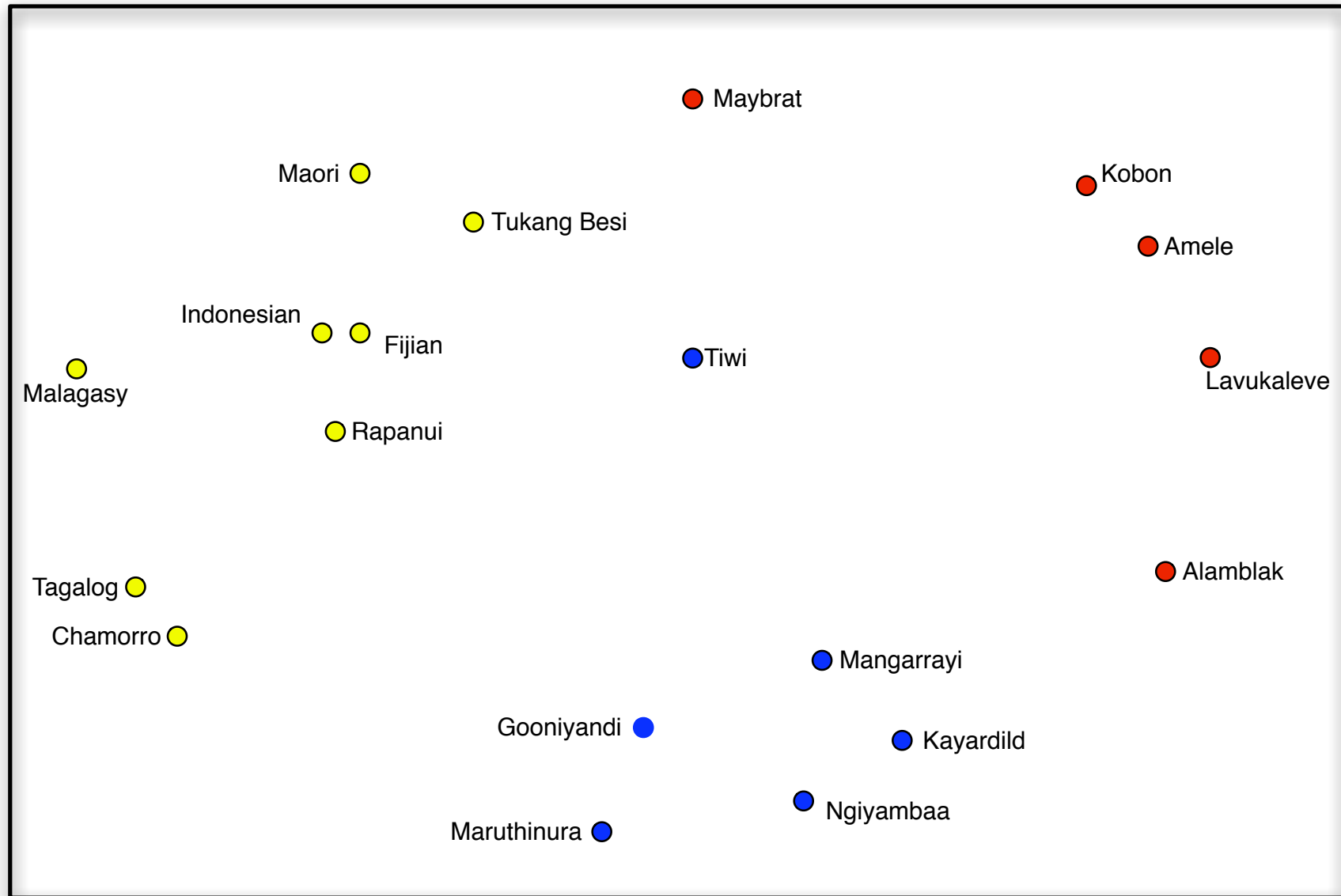


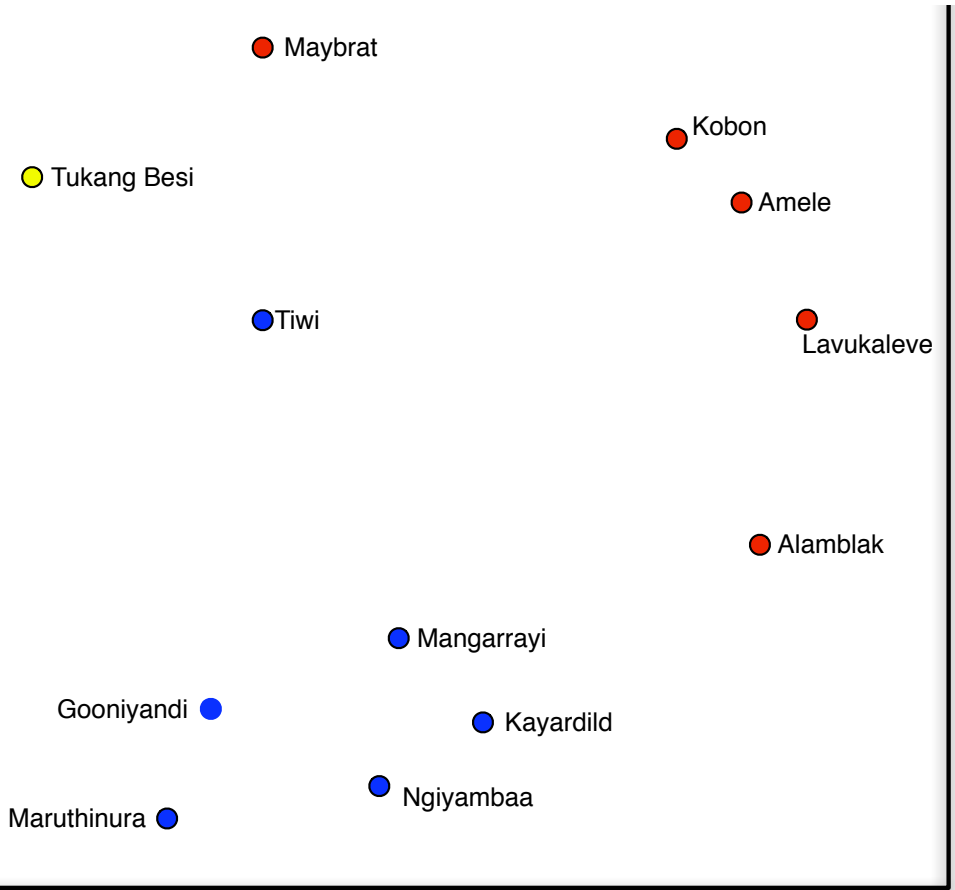




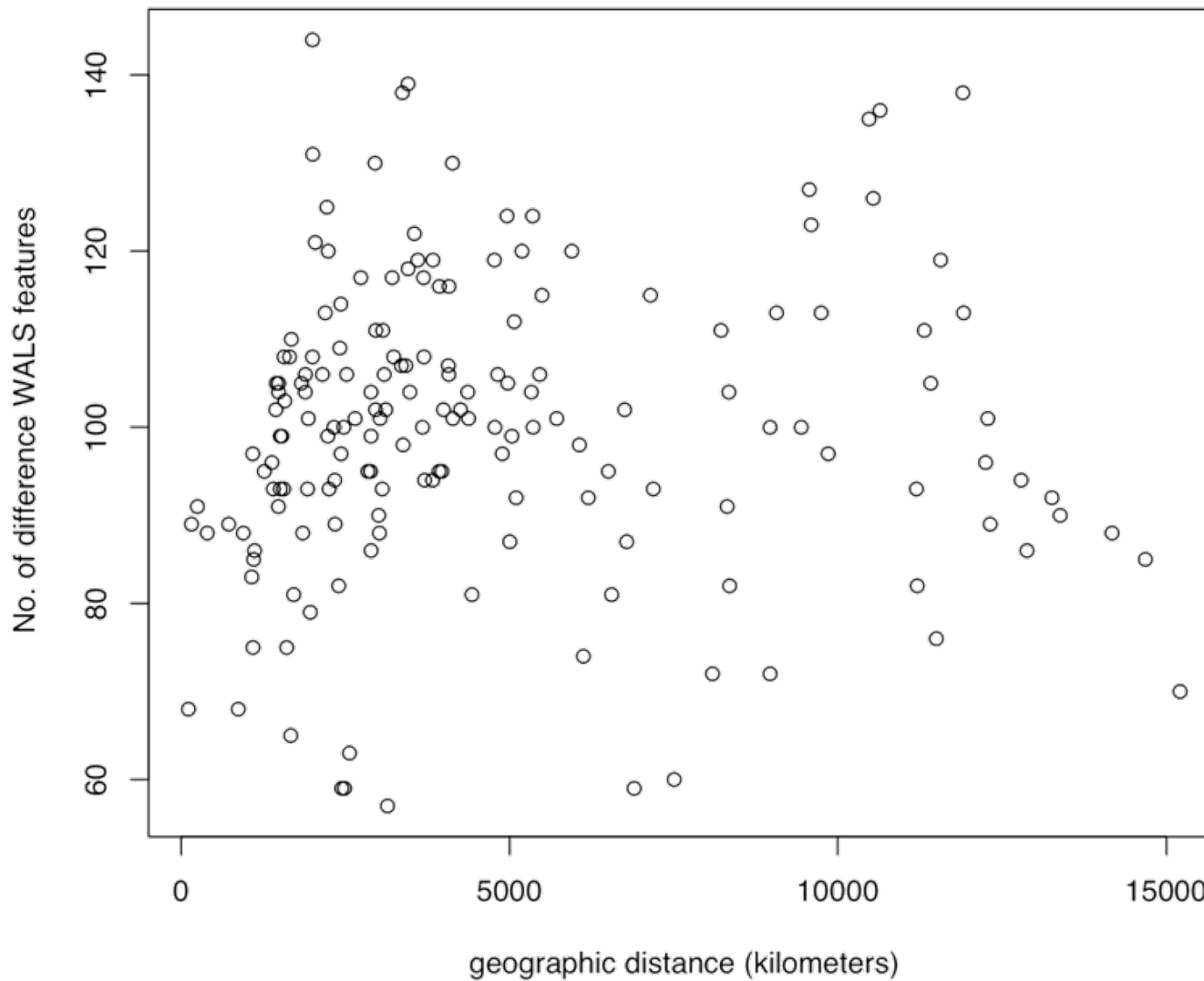


MDS of typological distances





Typology/geography correlation

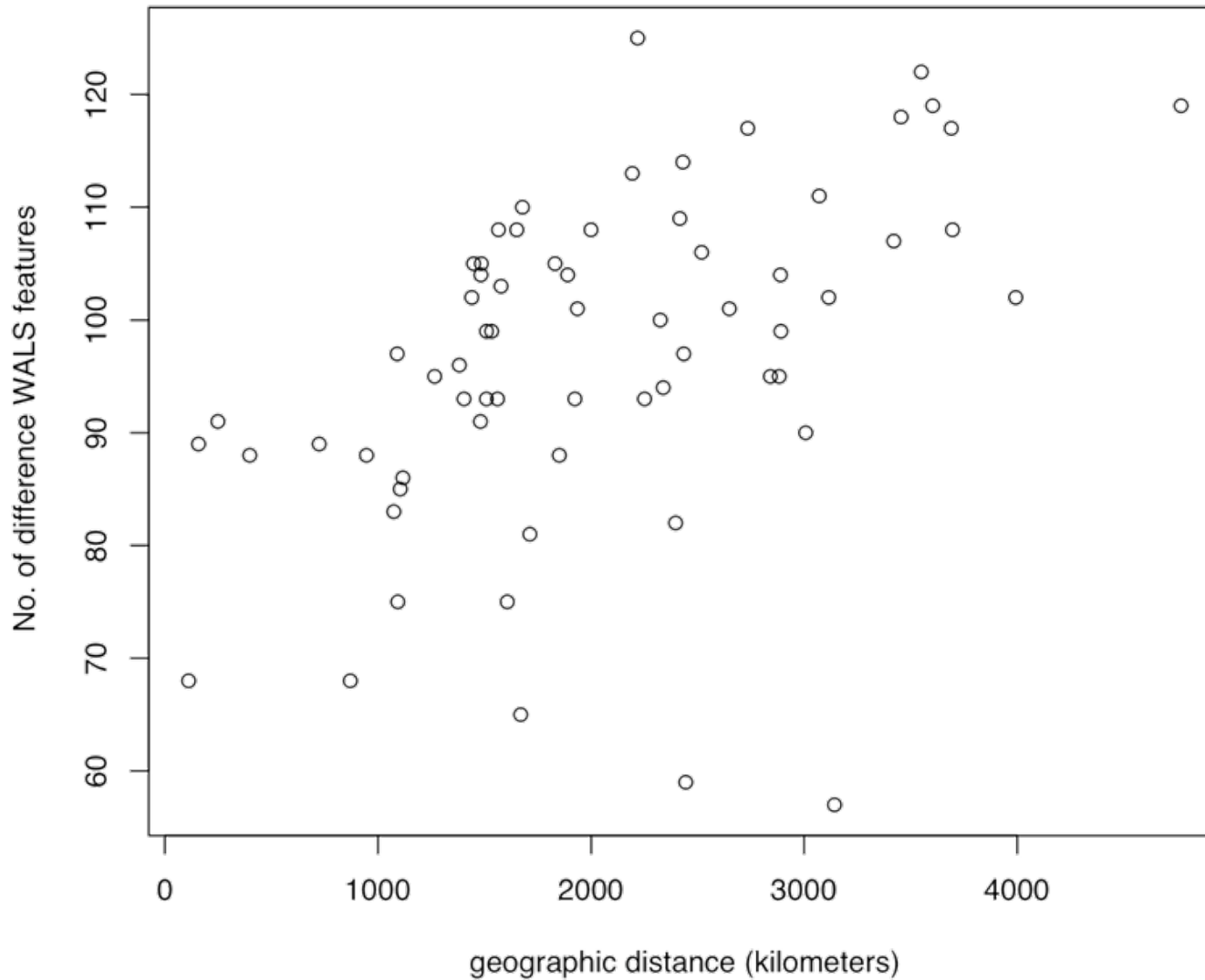


Mantel test
 $p = .349$

When does correlation improve?

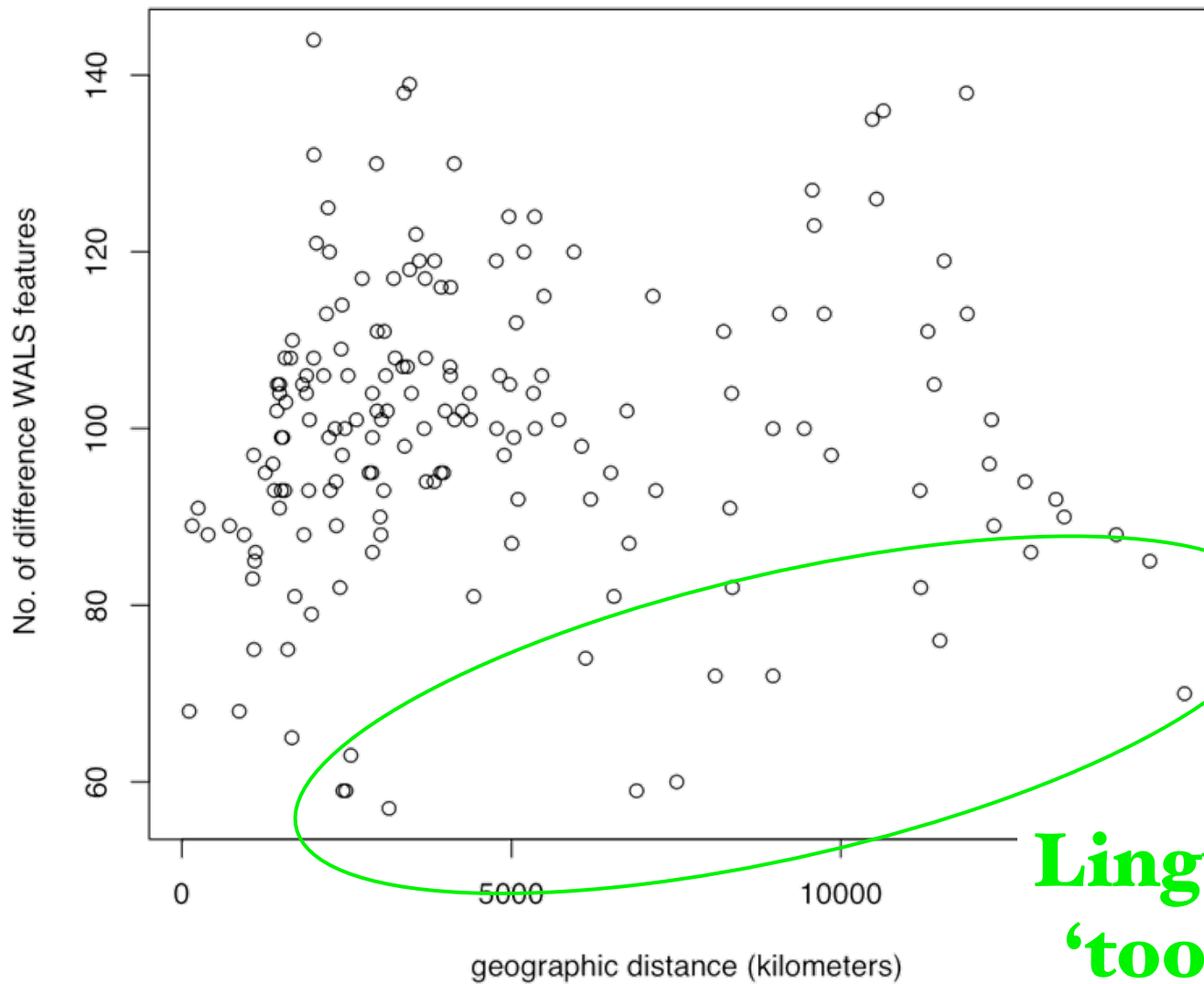
	Pearson's <i>r</i>
Nothing removed	.035
Rapanui	.186
Chamorro	.086
Indonesian	.076
Fijian	.073
Tagalog	.071
Maori	.062
Tukang Besi	.048

Correlation for selection



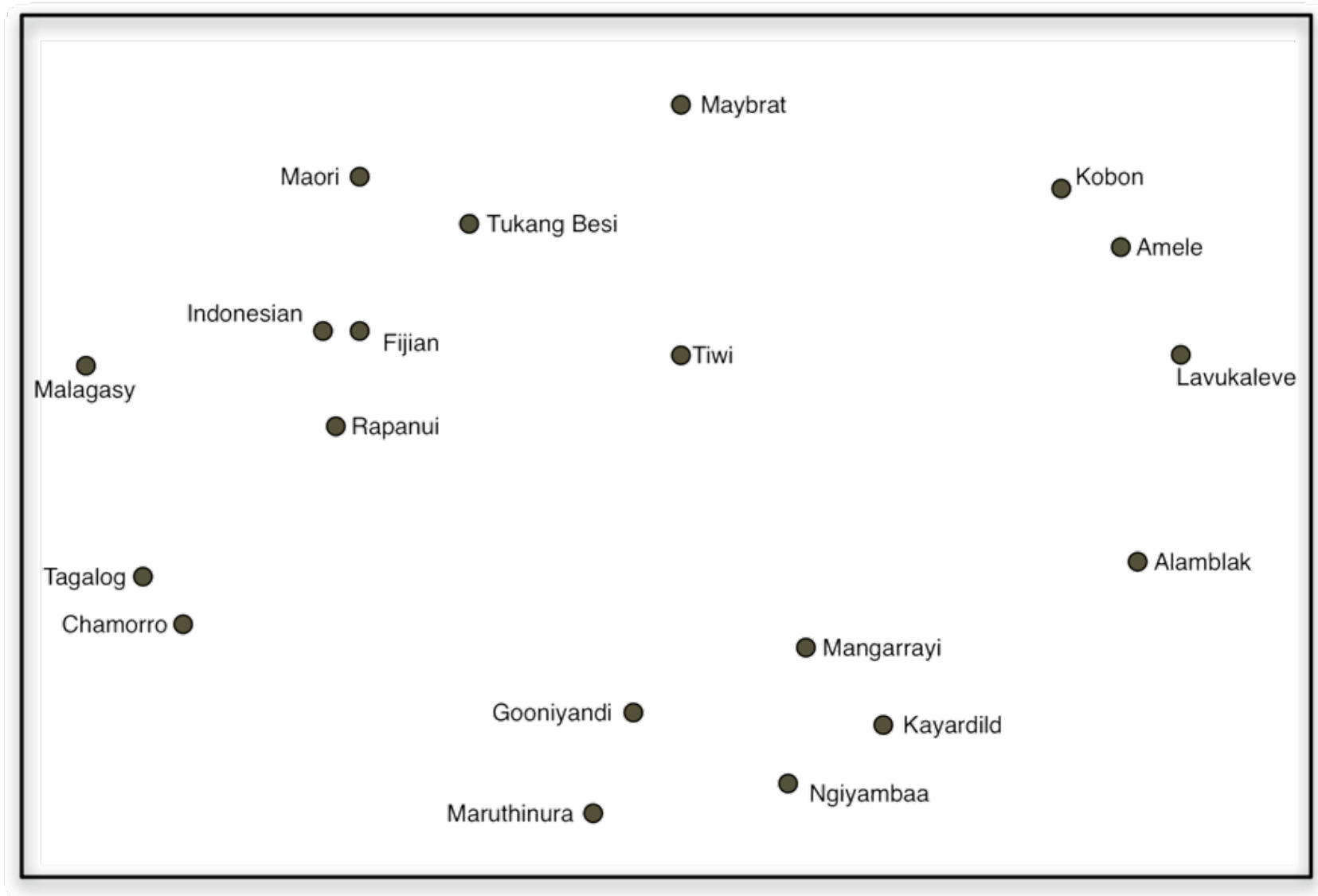
Mantel test p
= .001

Investigation typology/geography relation

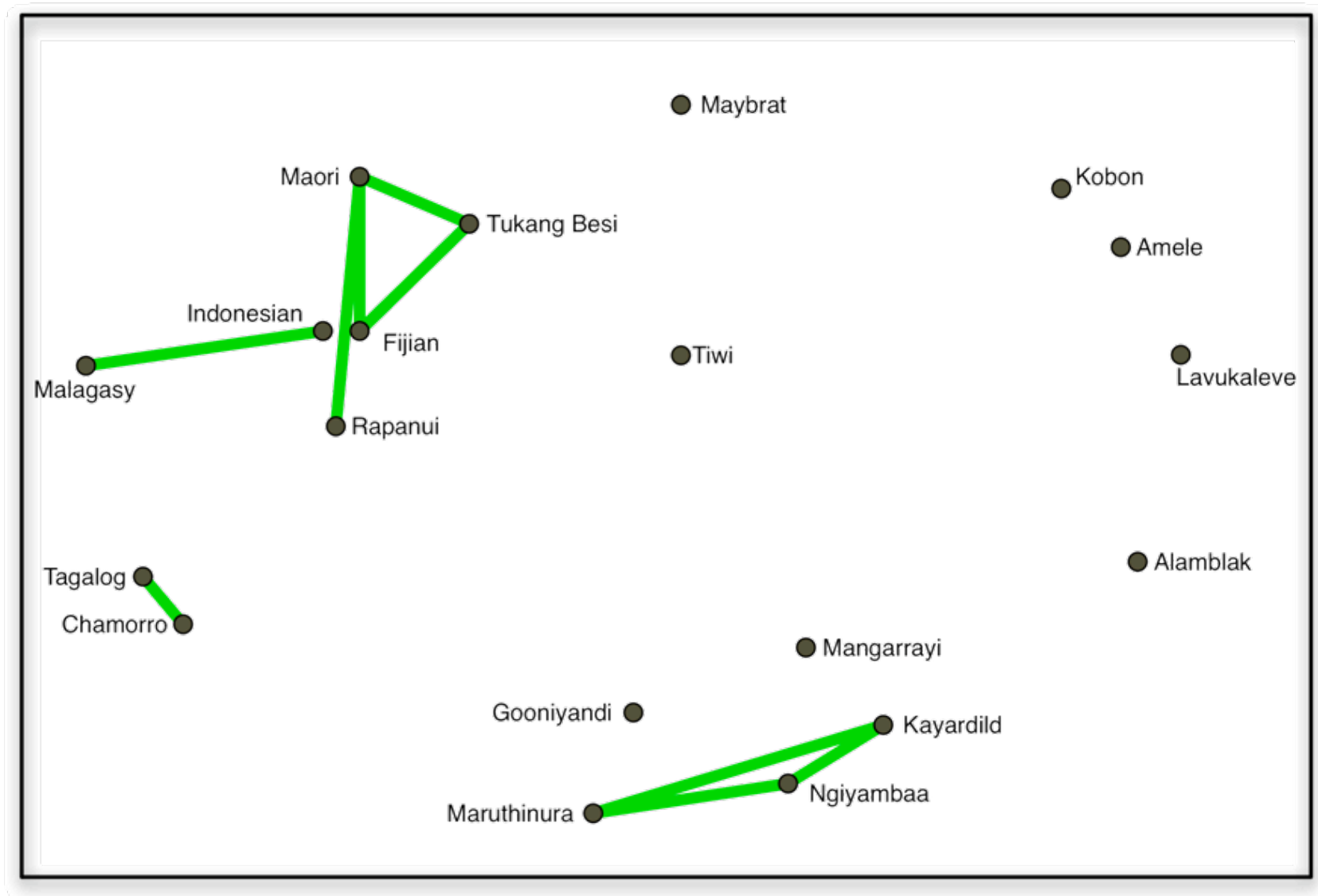


**Linguistically
'too similar'**

Linguistically 'too similar'



Linguistically 'too similar'





MAX-PLANCK-GESellschaft

**Max Planck Institute
for Evolutionary Anthropology**