

Language typology as relational measurement

Michael Cysouw
WiP, April 2008

Measurement theory

Measurement theory

- Stevens (1946)
 - ▶ from a psychological background

Measurement theory

- Stevens (1946)
 - ▶ from a psychological background
- proposed hierarchy of variables
 - ▶ nominal
 - ▶ ordinal
 - ▶ interval
 - ▶ ratio

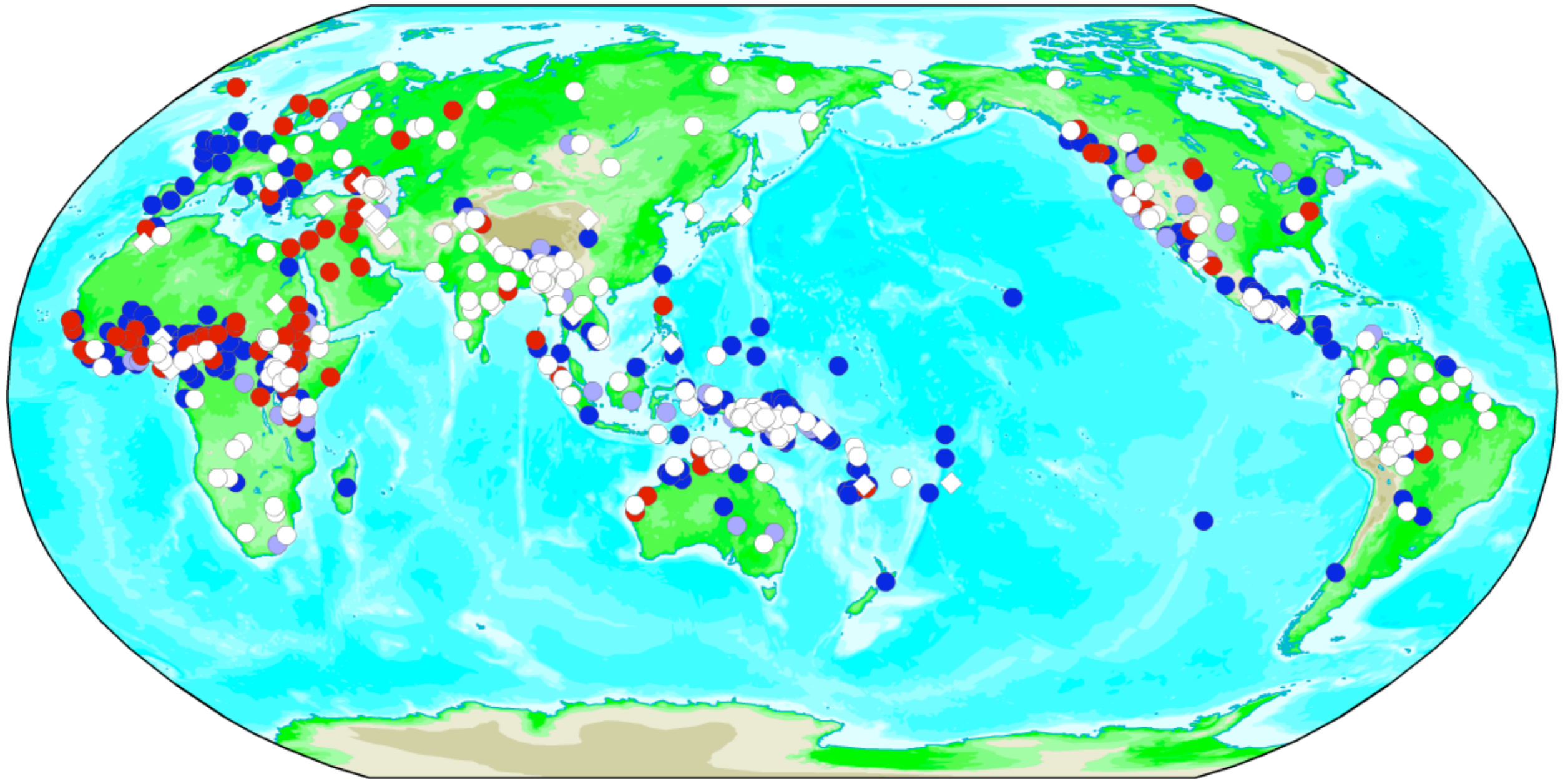
Measurement theory

- Stevens (1946)
 - ▶ from a psychological background
- proposed hierarchy of variables
 - ▶ nominal
 - ▶ ordinal
 - ▶ interval
 - ▶ ratio
- “yardstick” metaphor of measurement

Categorization


(nominal variable)

Categorization (nominal variable)



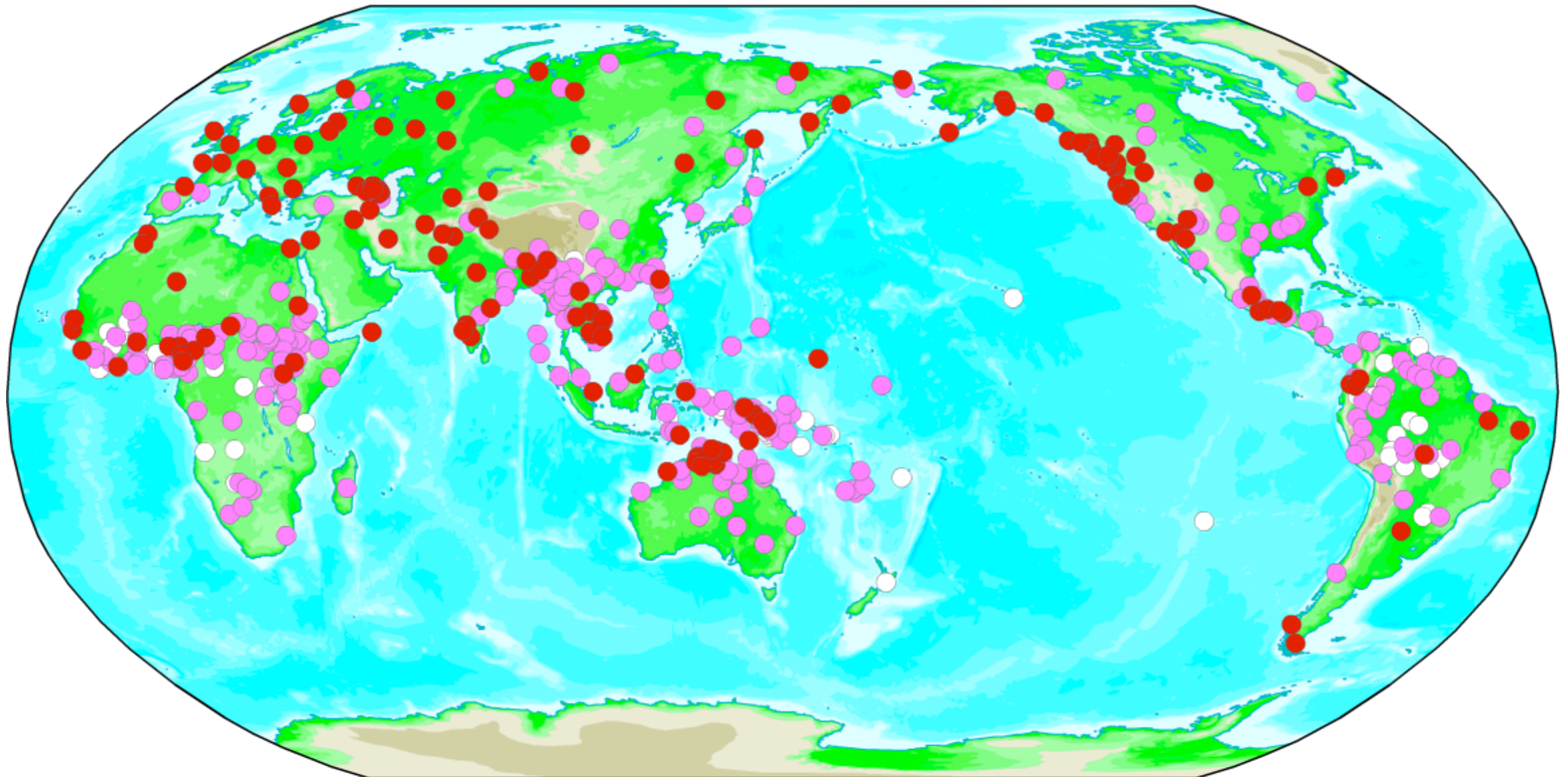
Dryer, Matthew S. (2005) 'Definite article' in: Martin Haspelmath, Matthew S. Dryer, David Gil, & Bernard Comrie (eds.) *World Atlas of Language Structures*. Oxford: Oxford University Press, 154-157.

Categorization (nominal variable)

- 
- A world map showing the distribution of five categories of definite and indefinite articles across various languages. The map is overlaid with colored dots and symbols representing different linguistic features. The categories are:
1. Definite word distinct from demonstrative
 2. Demonstrative word used as definite article
 3. Definite affix
 4. No definite, but indefinite article
 5. No definite or indefinite article

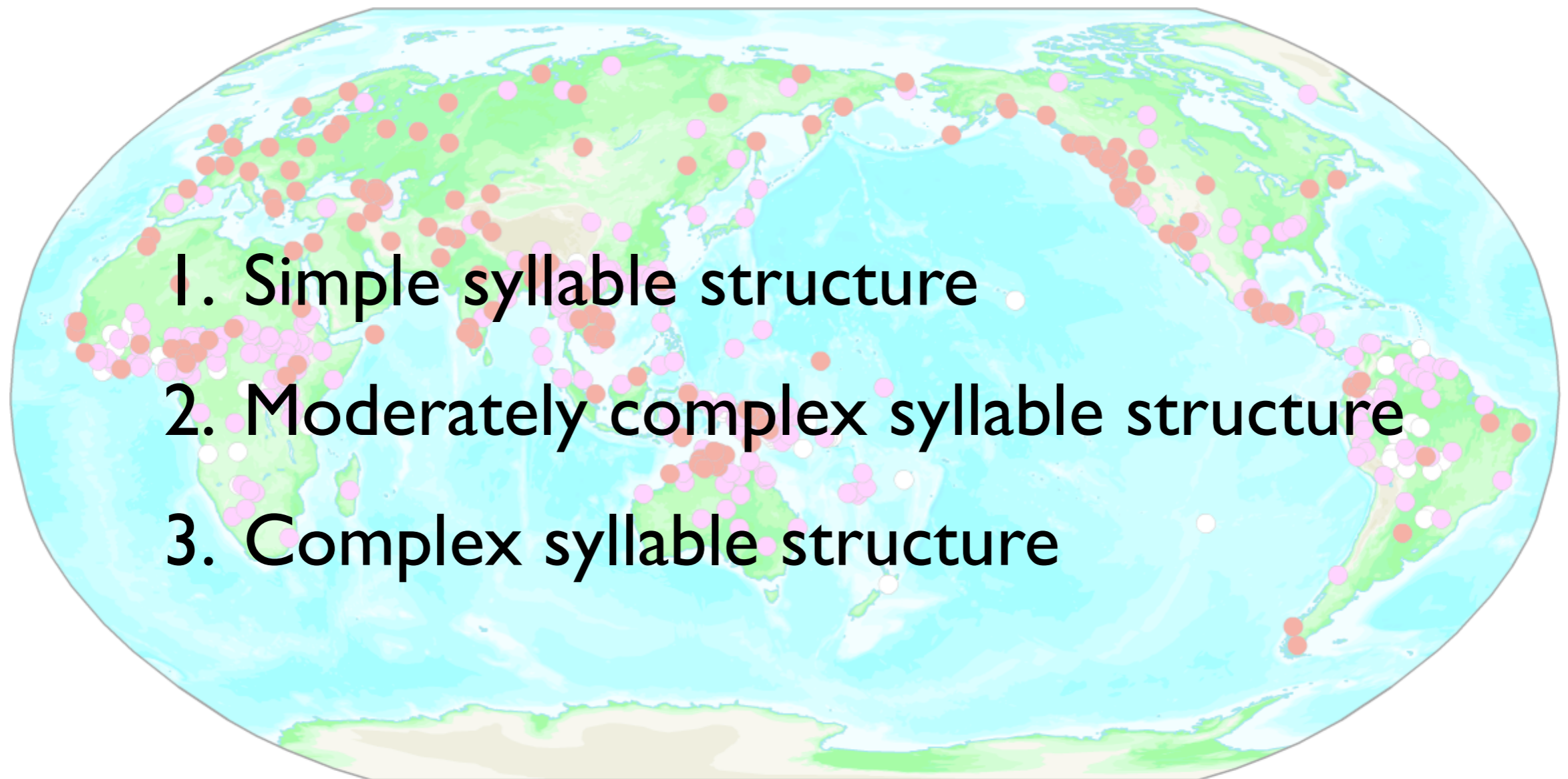
Linearly ordered categorization (interval variable)

Linearly ordered categorization (interval variable)



Maddieson, Ian (2005) 'Syllable structure' in: Martin Haspelmath, Matthew S. Dryer, David Gil, & Bernard Comrie (eds.) *World Atlas of Language Structures*. Oxford: Oxford University Press, 54-57.

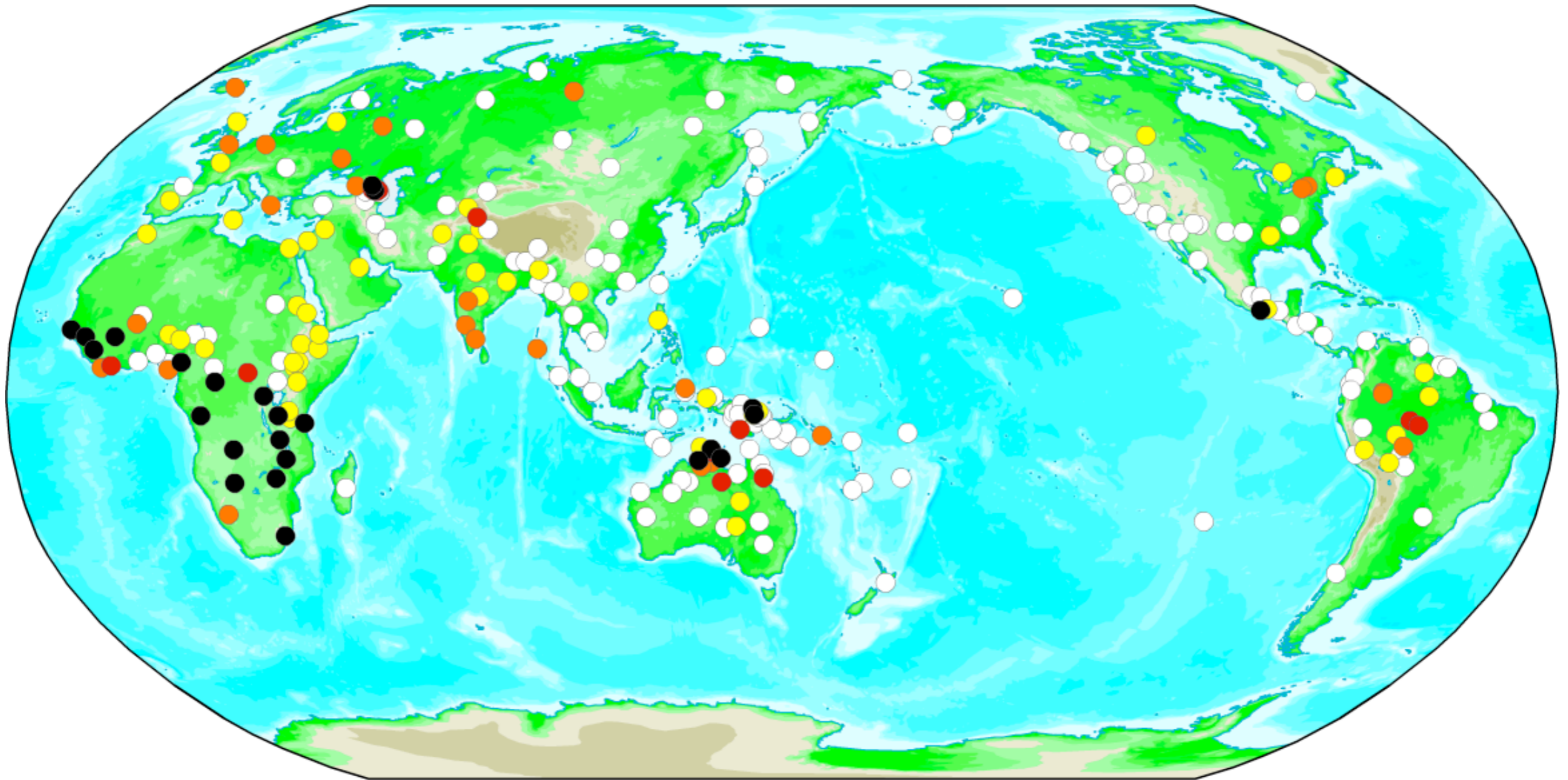
Linearly ordered categorization (interval variable)



Count
(ratio variable)

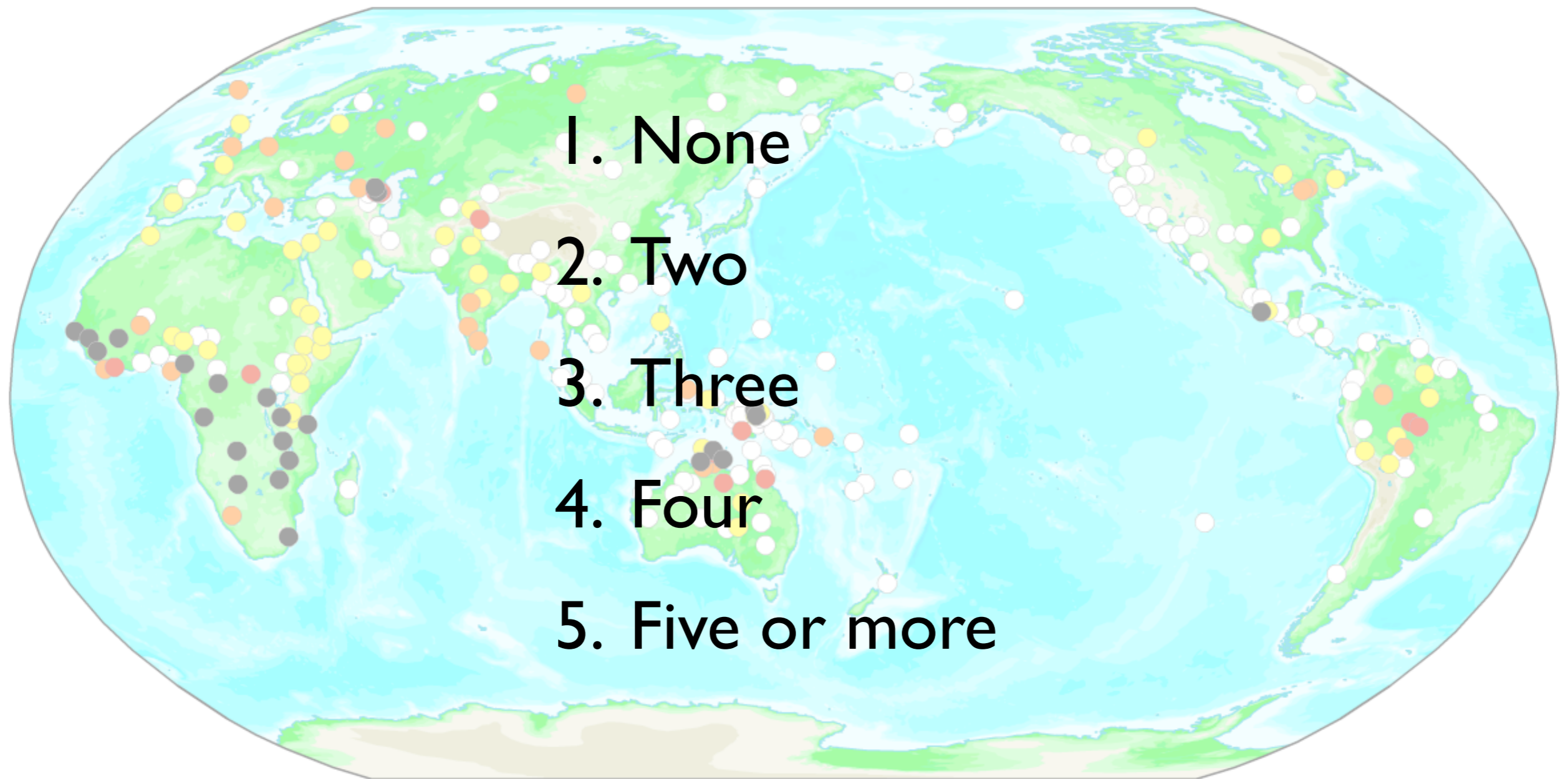
Count

(ratio variable)



Corbett, Greville G. (2005) 'Number of genders' in: Martin Haspelmath, Matthew S. Dryer, David Gil, & Bernard Comrie (eds.) *World Atlas of Language Structures*. Oxford: Oxford University Press, 126-129.

Count (ratio variable)



Continuum (ratio variable)

Continuum

(ratio variable)

- Not common in language comparison

Continuum

(ratio variable)

- Not common in language comparison
- Examples:
 - ▶ physical characteristics of speech
 - ▶ averages of corpus counts

Continuum (ratio variable)

Language	Average wordlength
Hmong Nua	3.72
English	5.05
German	6.23
Cashinahua	6.42
Bugis	6.45
Inuktitut	14.99

Continuum

(ratio variable)

- Not common in language comparison
- Examples:
 - ▶ physical characteristics of speech
 - ▶ averages of corpus counts
- Watch out with the interpretation of combinations of various (*a priori*) independent counts (e.g. sum or fraction)

Problems

Problems

- More measurements wanted

Problems

- More measurements wanted
 - ▶ more specification in categorization

Problems

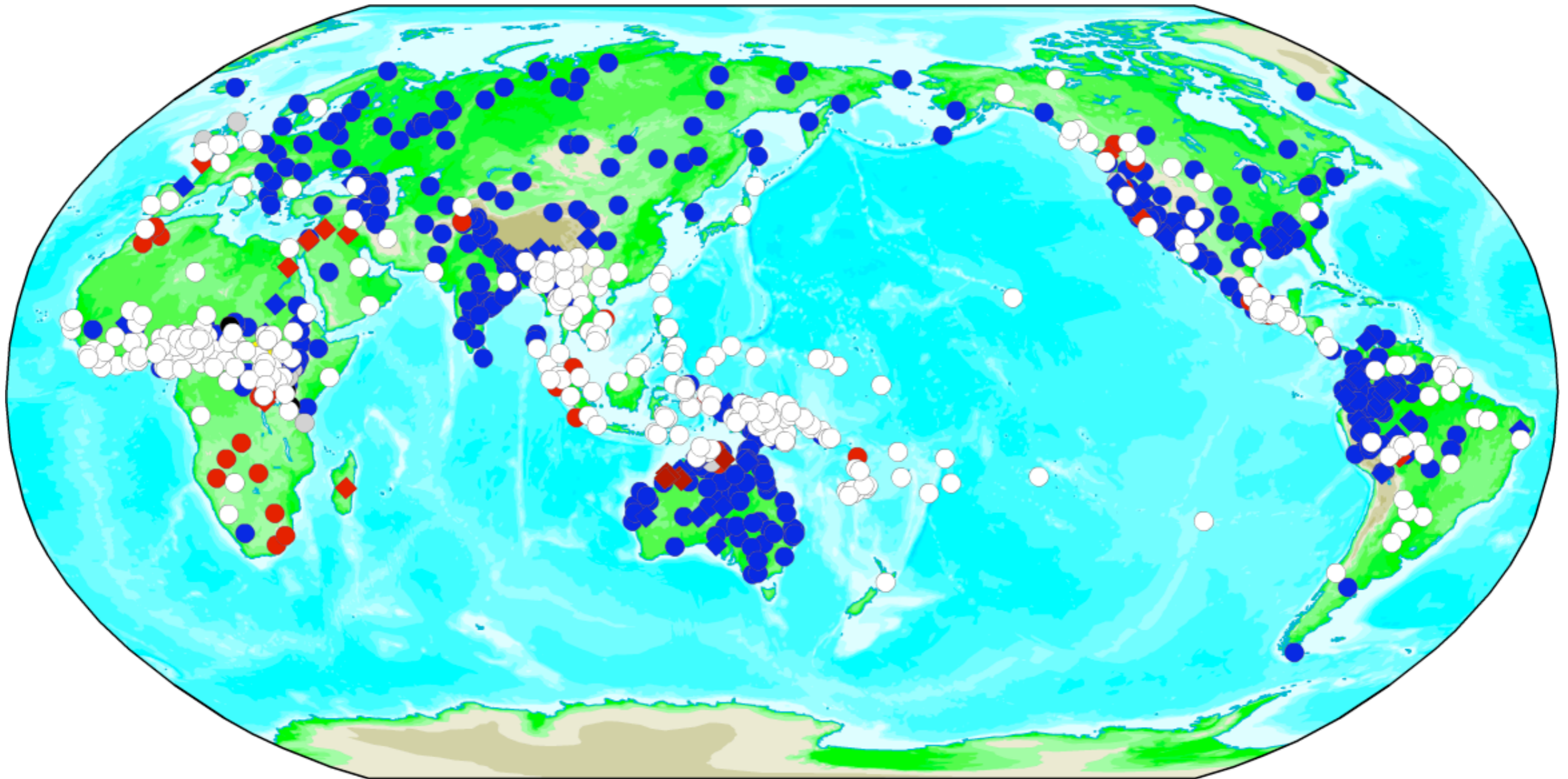
- More measurements wanted
 - ▶ more specification in categorization
 - ▶ full pairwise comparisons

Problems

- More measurements wanted
 - ▶ more specification in categorization
 - ▶ full pairwise comparisons
- Difficult to combine measurements of different kinds

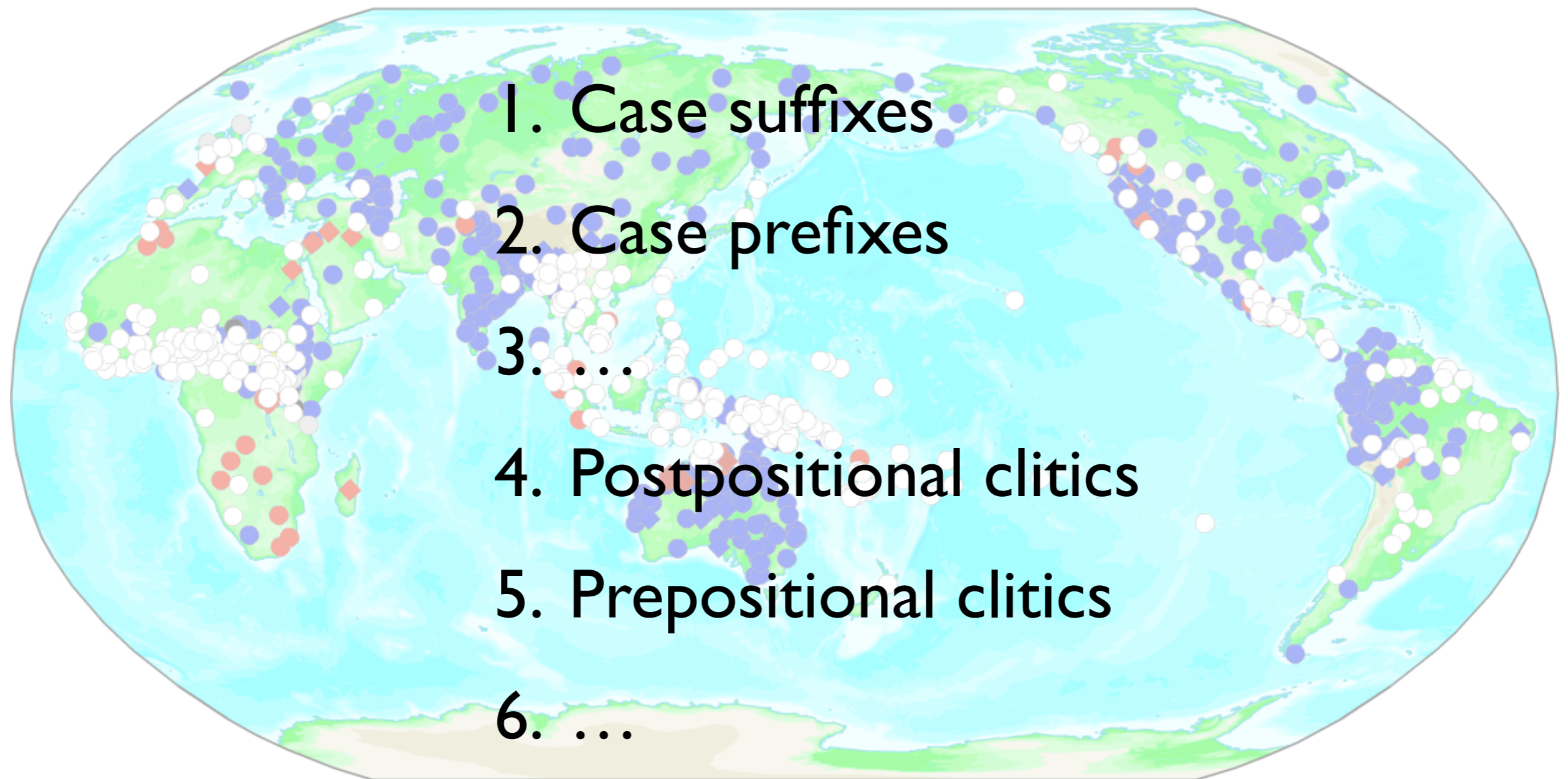
More specification for categorizations

More specification for categorizations



Dryer, Matthew S. (2005) 'Position of case affixes' in: Martin Haspelmath, Matthew S. Dryer, David Gil, & Bernard Comrie (eds.) *World Atlas of Language Structures*. Oxford: Oxford University Press, 210-213.

More specification for categorizations



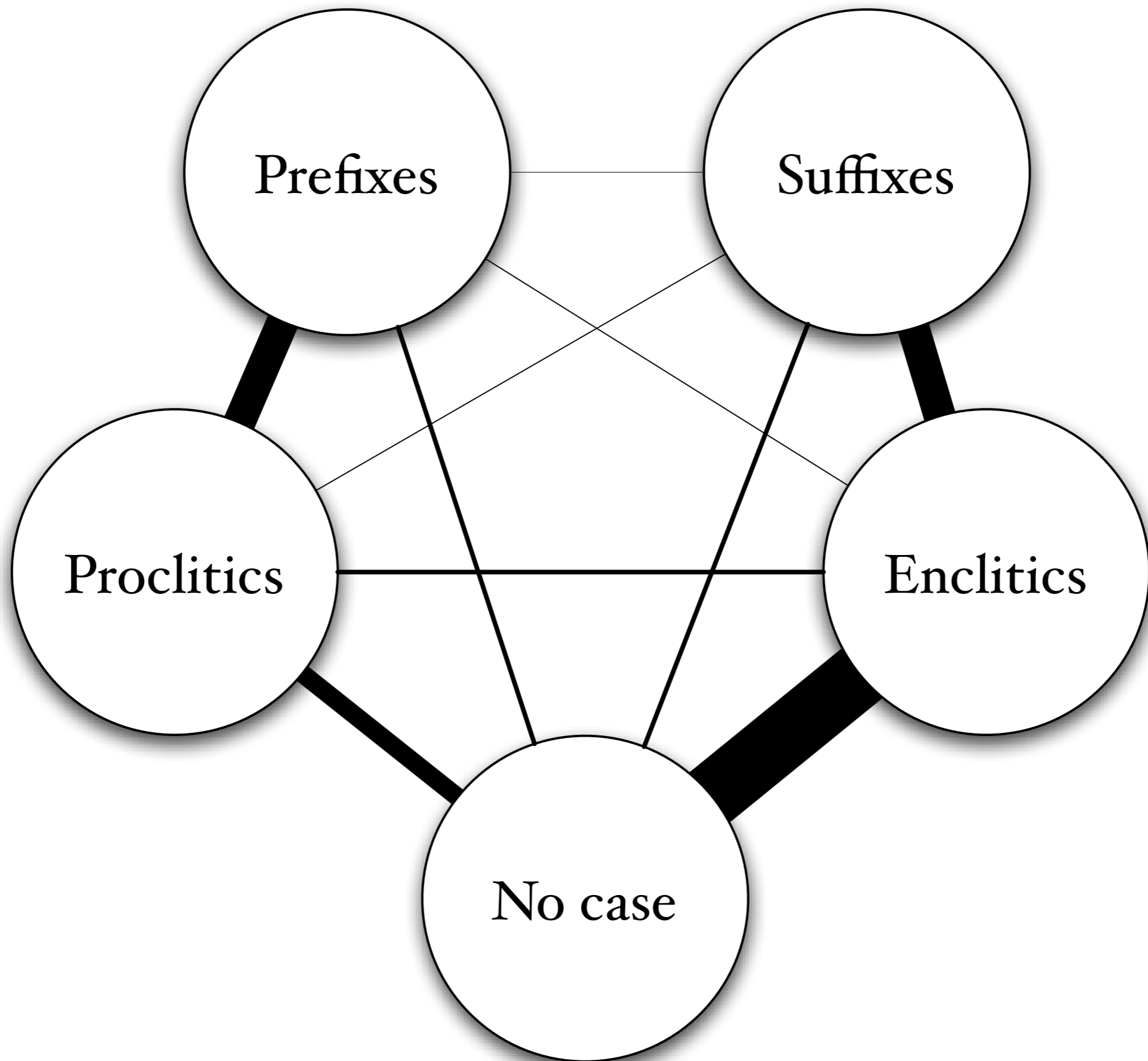
Prefixes

Suffixes

Proclitics

Enclitics

No case



Relational metaphor of measurement

Relational metaphor of measurement

- Express typology as pairwise language-to-language similarities

Relational metaphor of measurement

- Express typology as pairwise language-to-language similarities
- Such a typology consists of data with separate interpretation of the meaning of the data

L₁

L₃

L₂

L₅

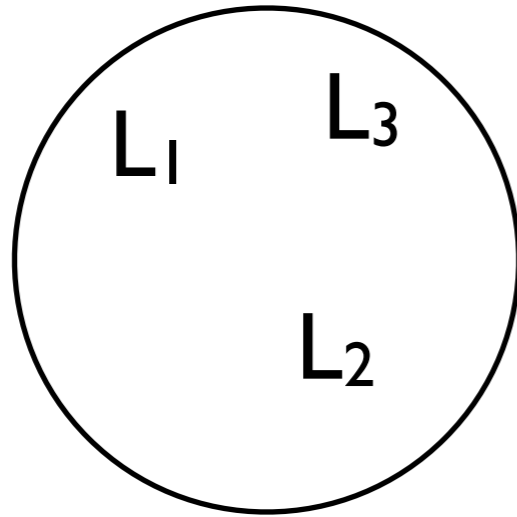
L₄

L₆

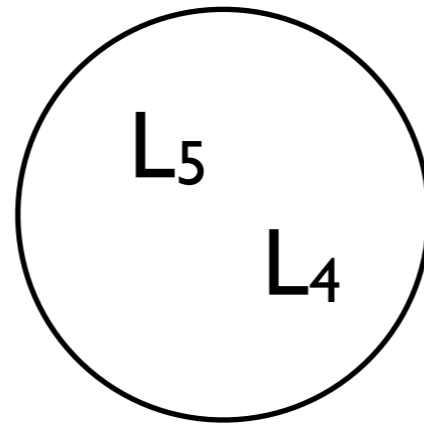
L₇

L₈

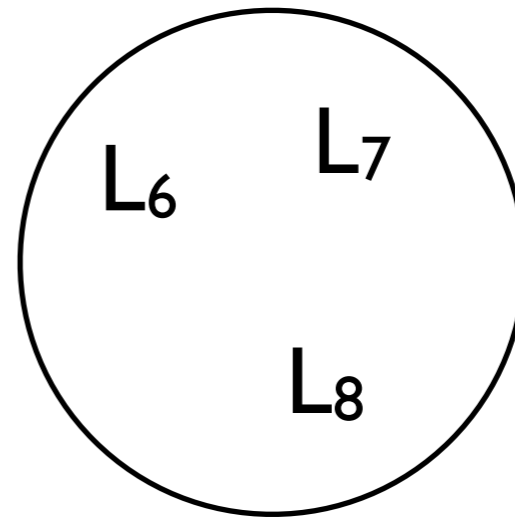
Type A



Type B



Type C



Type A

Type B

Type C

	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	
...									

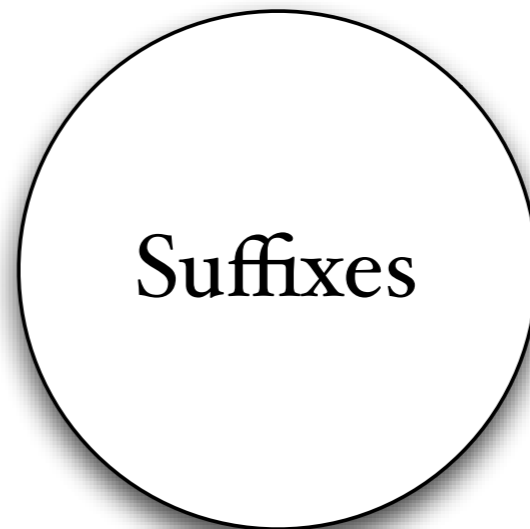
Undifferentiated Categorization

TYPE Δ

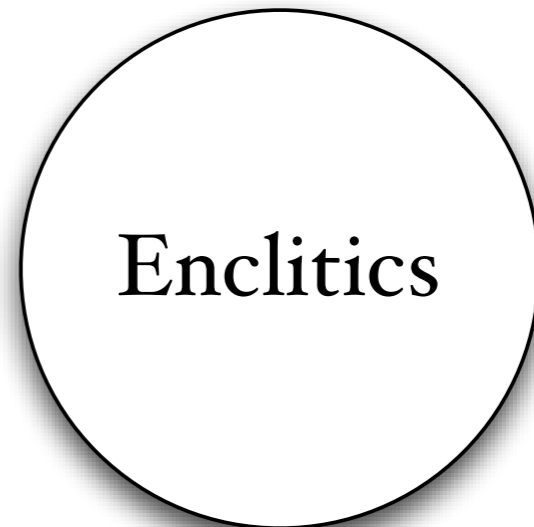
TYPE R

TYPE C

L ₁
L ₂
L ₃
L ₄
L ₅
L ₆
L ₇
L ₈
...



...



Undifferentiated Categorization

Type A

Type B

Type C

	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 Categorization

Type A

Type B

Type C

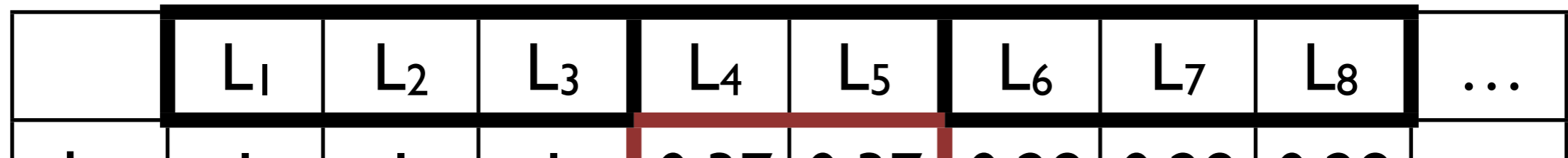
	L ₁	L ₂	L ₃	L ₄	L ₅	L ₆	L ₇	L ₈	...
L ₁				0.37	0.37	0.28	0.28	0.28	
L ₂				0.37	0.37	0.28	0.28	0.28	
L ₃				0.37	0.37	0.28	0.28	0.28	
L ₄	0.37	0.37	0.37			0.51	0.51	0.51	
L ₅	0.37	0.37	0.37			0.51	0.51	0.51	
L ₆	0.28	0.28	0.28	0.51	0.51				
L ₇	0.28	0.28	0.28	0.51	0.51				
L ₈	0.28	0.28	0.28	0.51	0.51				
...									

Differentiated Categorization

Type A

Type B

Type C



A. Simple syllable structure

B. Moderately complex syllable structure

C. Complex syllable structure

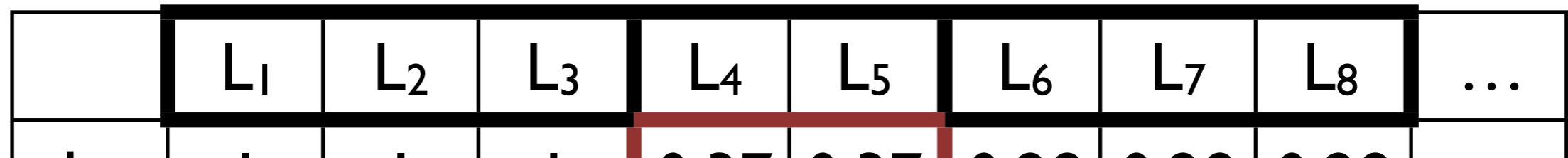


Differentiated Categorization

Type A

Type B

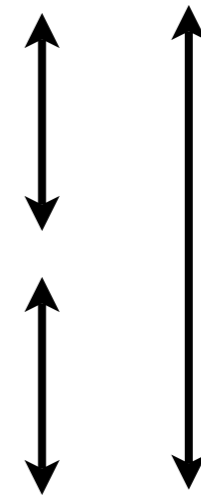
Type C



A. Simple syllable structure

B. Moderately complex syllable structure

C. Complex syllable structure



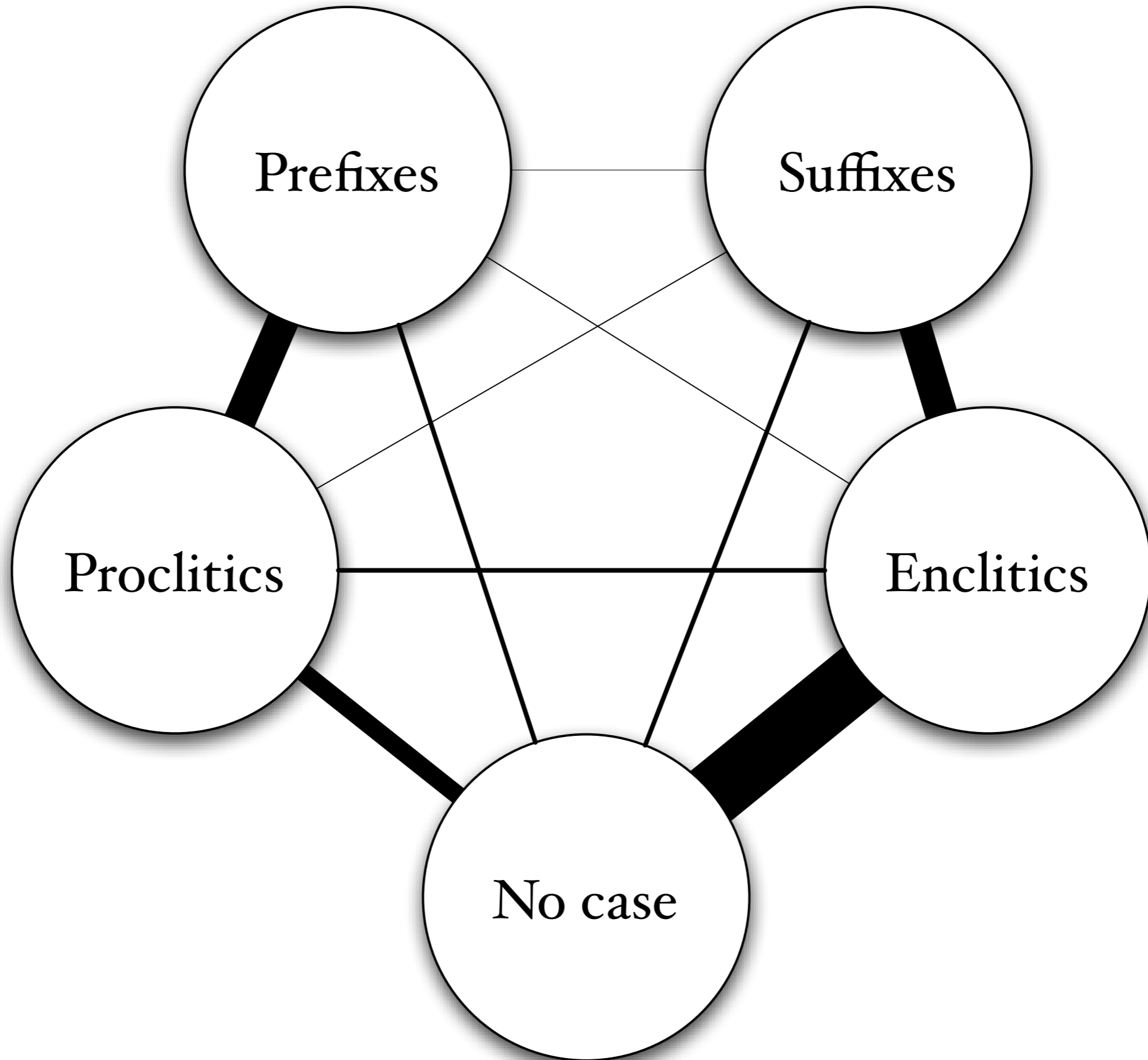
Differentiated Categorization

TYPE Δ

TYPE R

TYPE C

L ₁
L ₂
L ₃
L ₄
L ₅
L ₆
L ₇
L ₈
...



...

Differentiated Categorization

Type A

Type B

Type C

	L ₁	L ₂	L ₃	L ₄	L ₅	L ₆	L ₇	L ₈	...
L ₁				0.37	0.37	0.28	0.28	0.28	
L ₂				0.37	0.37	0.28	0.28	0.28	
L ₃				0.37	0.37	0.28	0.28	0.28	
L ₄	0.37	0.37	0.37			0.51	0.51	0.51	
L ₅	0.37	0.37	0.37			0.51	0.51	0.51	
L ₆	0.28	0.28	0.28	0.51	0.51				
L ₇	0.28	0.28	0.28	0.51	0.51				
L ₈	0.28	0.28	0.28	0.51	0.51				
...									

Differentiated Categorization

	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	
...									

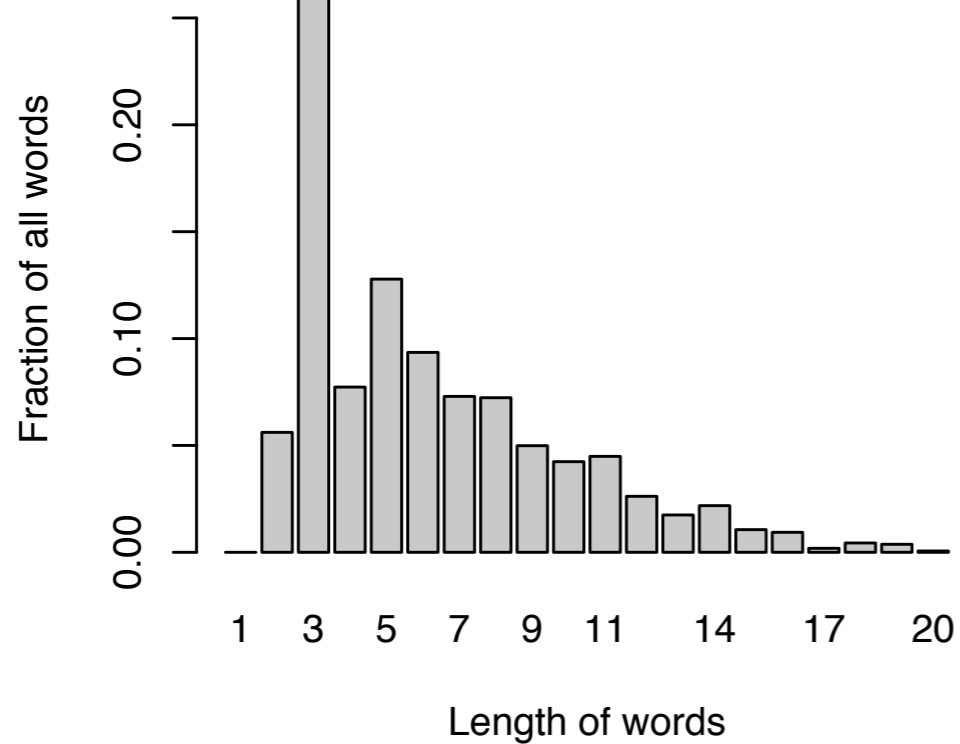
‘Deconstructed’ Typology

Pairwise Comparison

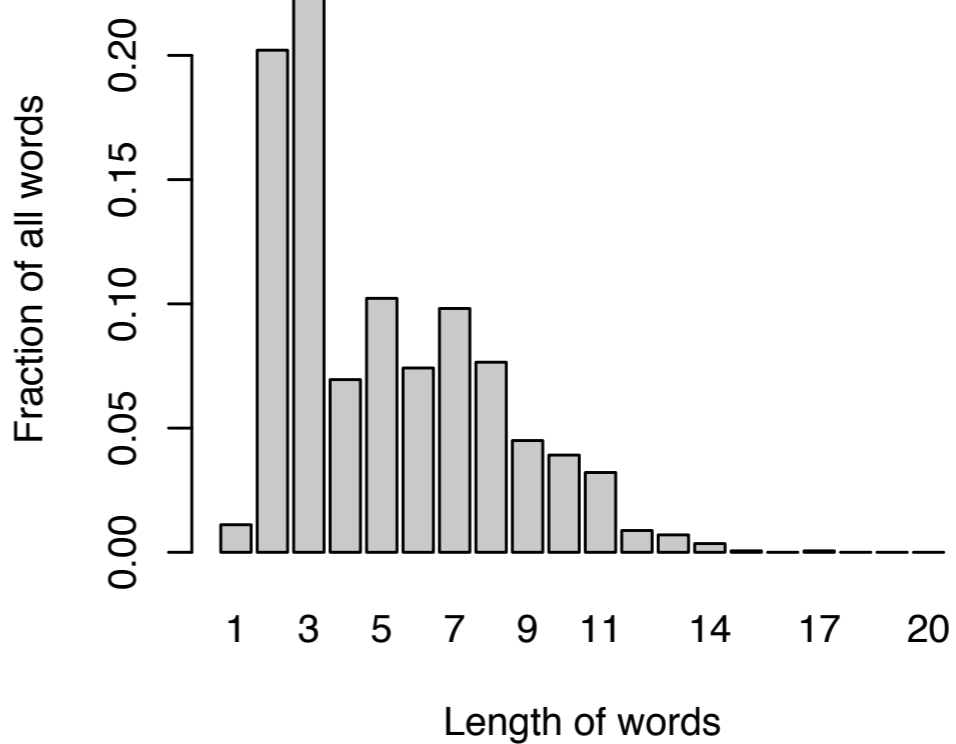
Pairwise Comparison

Language	Average wordlength
Hmong Nua	3.72
English	5.05
German	6.23
Cashinahua	6.42
Bugis	6.45
Inuktitut	14.99

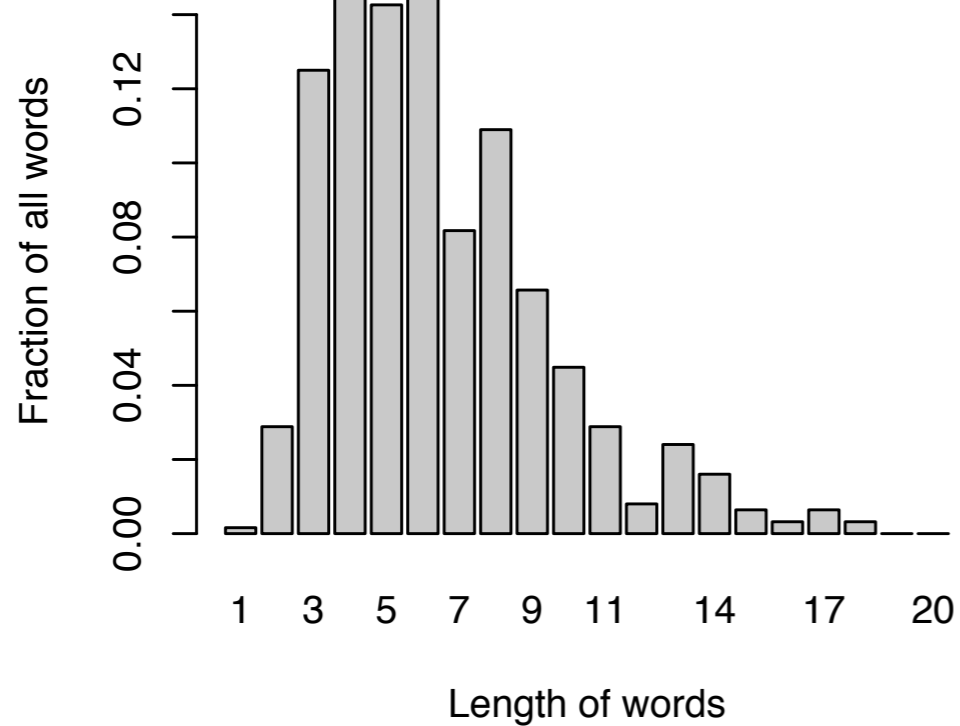
German



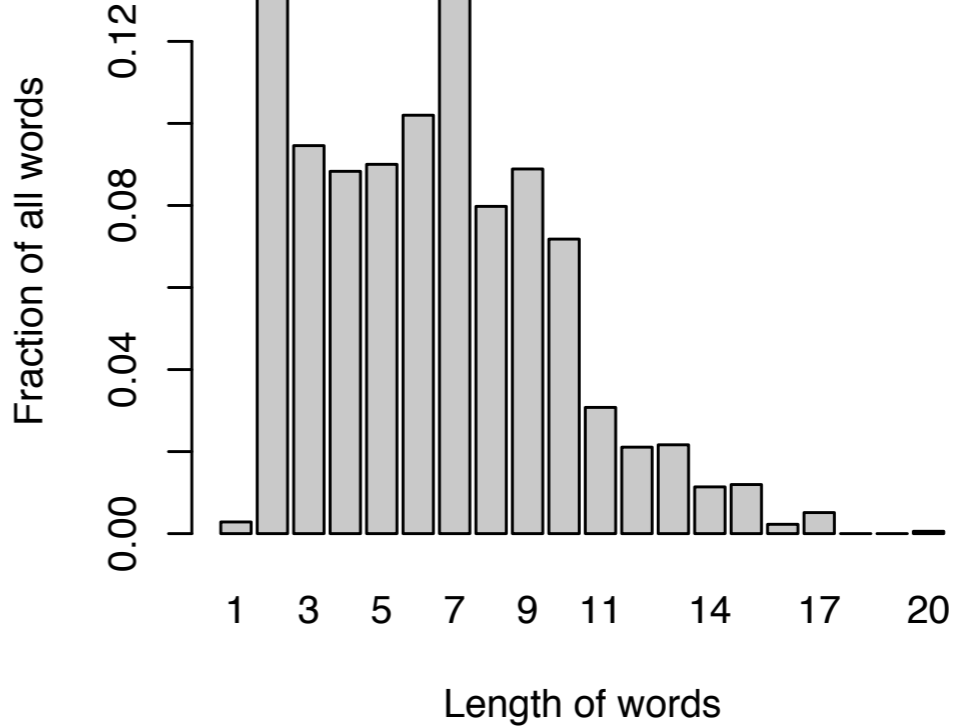
English



Cashinahua

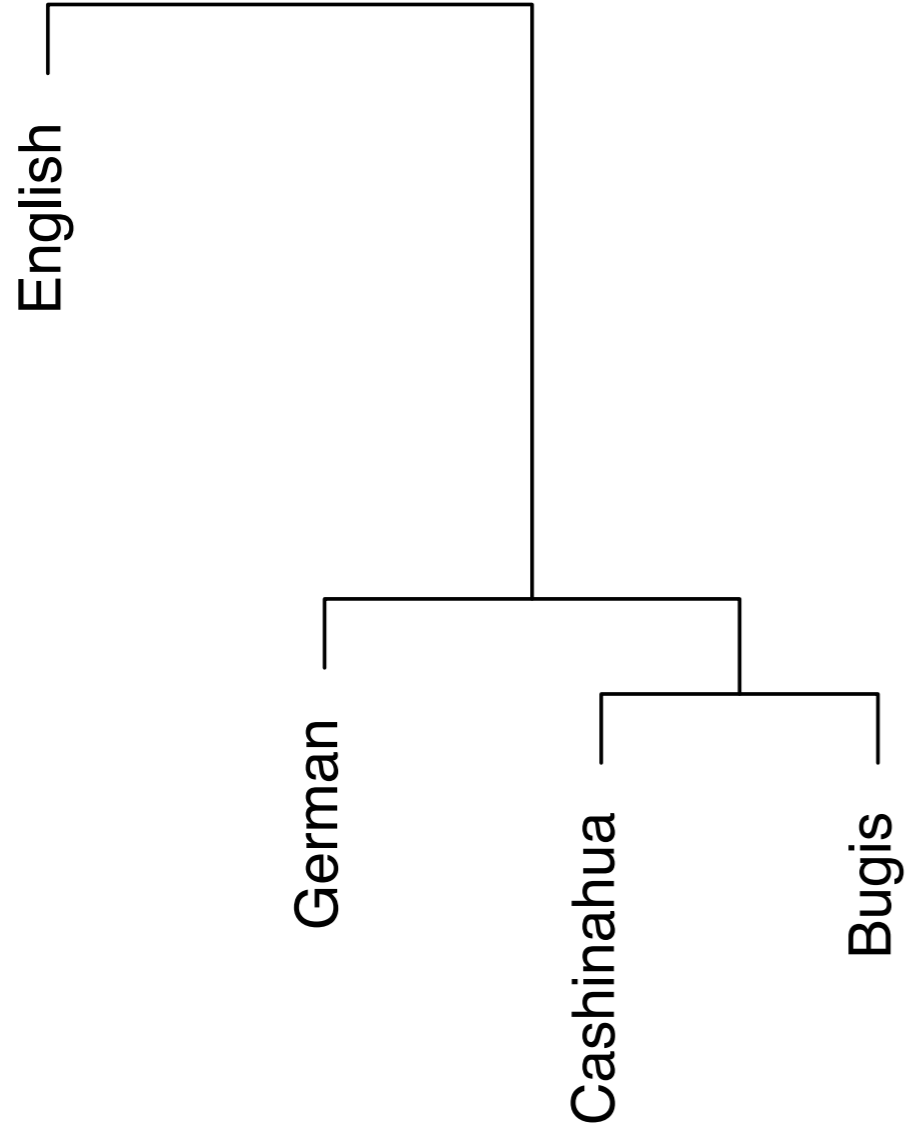


Bugis

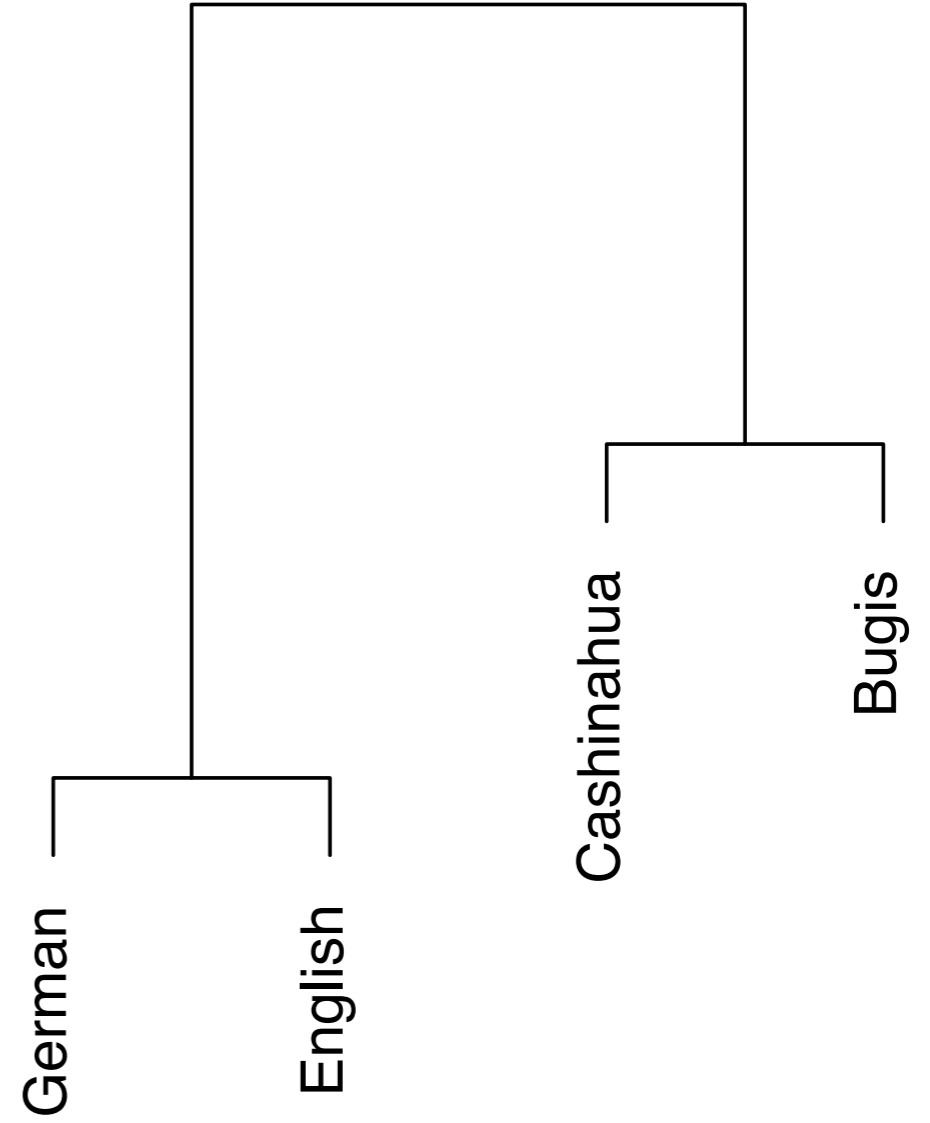


Hmong Nua	0	0.60	0.53	0.58	0.74	1
English	0.60	0	0.19	0.32	0.23	0.74
German	0.53	0.19	0	0.23	0.27	0.66
Cashinahua	0.58	0.32	0.23	0	0.25	0.70
Bugis	0.74	0.23	0.27	0.25	0	0.68
Inuktitut	1	0.74	0.66	0.70	0.68	0

Average wordlength



Wordlength distribution



Data

Similarity

	Data	Similarity
undifferentiated categorization	types of languages	implicit (all types are equally dissimilar to each other)

	Data	Similarity
undifferentiated categorization	types of languages	implicit (all types are equally dissimilar to each other)
count, continuum	values measured	implicit similarity function (typically absolute difference)

	Data	Similarity
undifferentiated categorization	types of languages	implicit (all types are equally dissimilar to each other)
count, continuum	values measured	implicit similarity function (typically absolute difference)
differentiated categorization	types of languages	specify similarity between types (empirically or not)

	Data	Similarity
undifferentiated categorization	types of languages	implicit (all types are equally dissimilar to each other)
count, continuum	values measured	implicit similarity function (typically absolute difference)
differentiated categorization	types of languages	specify similarity between types (empirically or not)
count, continuum	values measured	specific similarity function (e.g. stressing low differences)

	Data	Similarity
undifferentiated categorization	types of languages	implicit (all types are equally dissimilar to each other)
count, continuum	values measured	implicit similarity function (typically absolute difference)
differentiated categorization	types of languages	specify similarity between types (empirically or not)
count, continuum	values measured	specific similarity function (e.g. stressing low differences)
'deconstructed' typology	whatever (e.g. collection of word lengths)	whatever (e.g. histogram similarity)

Language similarities ?!

- Similarities between languages do not follow automatically from the data !
- It has to be explicitly stated how the similarities are arrived at
- Different kinds of similarities are possible with the same data

Typology as language similarities

Typology as language similarities

- Separating data and similarity is both a curse and a blessing

Typology as language similarities

- Separating data and similarity is both a curse and a blessing
- **Curse:** it is necessary to be much more precise in what it takes for two languages to be similar

Typology as language similarities

- Separating data and similarity is both a curse and a blessing
- **Curse:** it is necessary to be much more precise in what it takes for two languages to be similar
- **Blessing:** Such precision results in much more consistent and fine-grained language typologies