

The background of the slide is a stylized map of Southeast Asia, including the Malay Peninsula and the Indonesian archipelago. The landmasses are colored in shades of green and light blue, while the surrounding waters are a pale cyan. Numerous small, semi-transparent circles in various colors (pink, yellow, orange, purple, brown, and grey) are scattered across the map, representing the geographical distribution of different languages.

# Measuring Language Similarities

Michael Cysouw  
Geneling, 17 Juli 2006



# Preamble

- Measuring language similarities is independent of their interpretation
  - ▶ Genealogical Relationship ?
  - ▶ Borrowing ?
  - ▶ Human Language Universal ?
- It is a separate research questions which similarities should have what interpretation

# Method of comparison depends on datastructure

## Primary data

Elements of languages

## Secondary data

Statements about languages

Utterances

Lexemes

Descriptions

**Monolingual  
data**

Text

Dictionary

Grammar

# Measurements on Monolingual Datasources

- Textcounts, e.g.:
  - ▶ Greenberg's morphological indices
  - ▶ Givón's anaphorical distance
  - ▶ Bickel's referential density
- Gradient Grammatical Indices, e.g.:
  - ▶ Maddieson's number of consonants
  - ▶ Bickel's index of inflectional synthesis
- Dictionary counts ?

# Method of comparison depends on datastructure

## Primary data

Elements of languages

## Secondary data

Statements about languages

Utterances

Lexemes

Descriptions

**Monolingual  
data**

Text

Dictionary

Grammar

# Method of comparison depends on datastructure

## Primary data

Elements of languages

## Secondary data

Statements about languages

Utterances

Lexemes

Descriptions

**Monolingual data**

Text

Dictionary

Grammar

**Comparative data**

Contextually  
situated exemplars

Wordlists

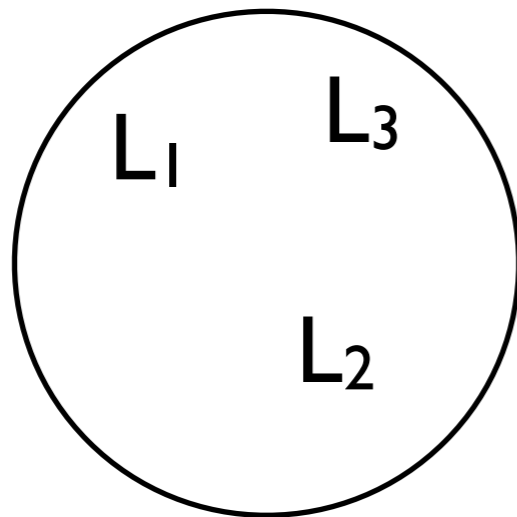
Typological  
databases

A world map with a light blue background and green landmasses. Numerous small, semi-transparent colored dots (pink, purple, red, white, and blue) are scattered across the map, representing data points for typological databases. The dots are most densely clustered in the Americas and Europe, with a notable concentration of red dots along the eastern coast of North America and the western coast of South America. Other colors are more sparsely distributed across the continents.

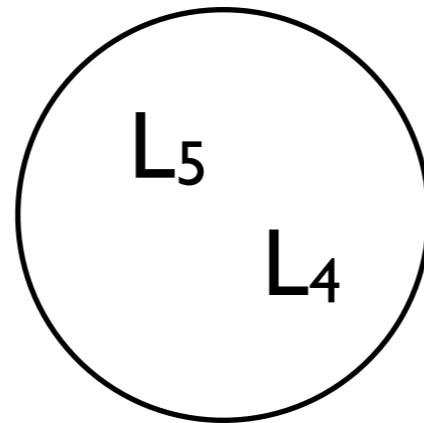
# **Typological Databases**

# Typology: Grouping of Languages into Types

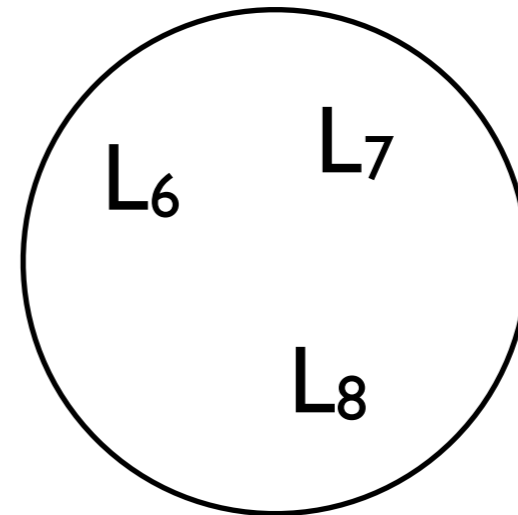
Type A



Type B

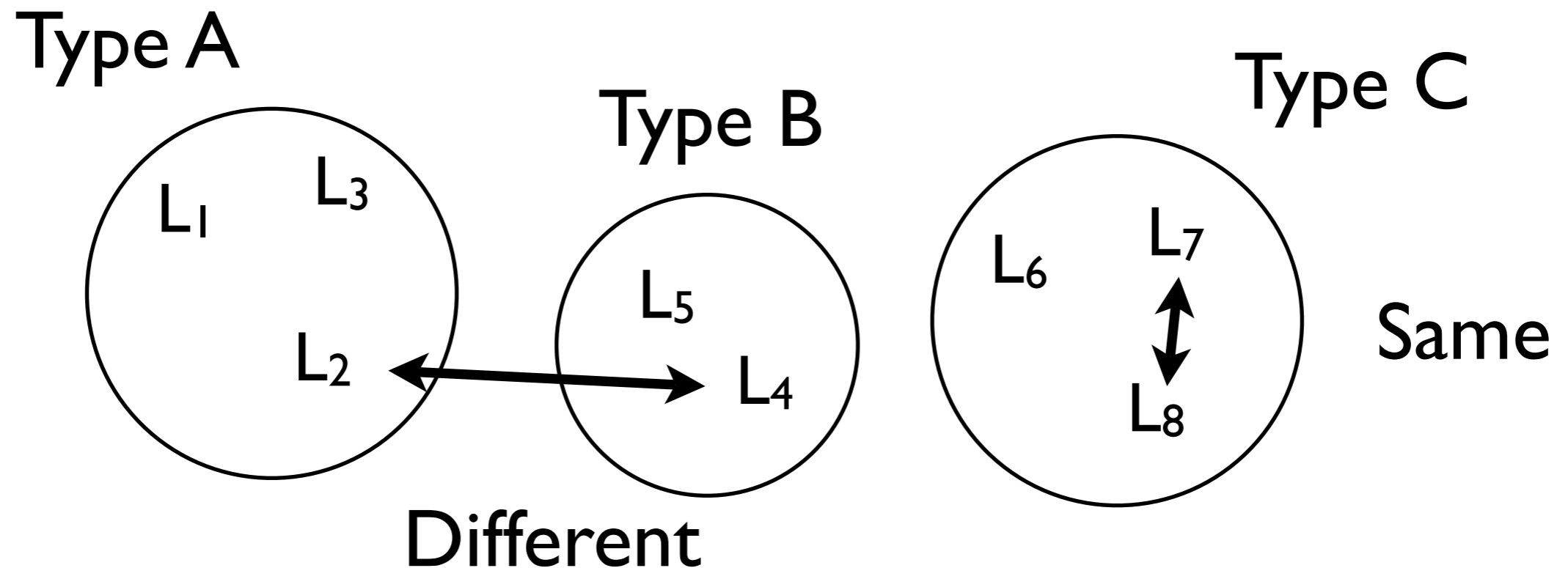


Type C





# Little Information ...



*We want more fine-grained measurements of similarity !*

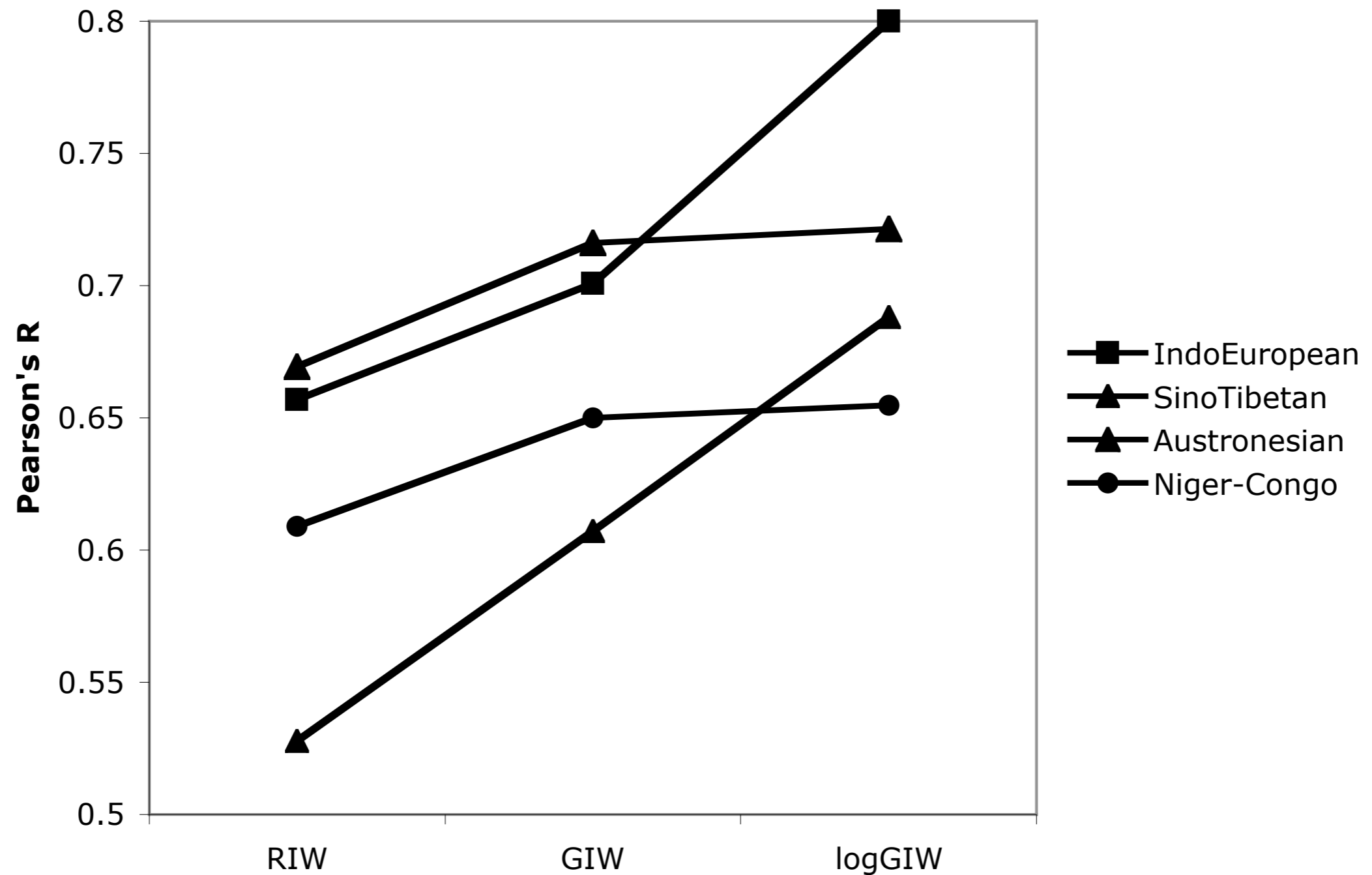
# Solution: Use many Typologies (think: WALS)

- Measures of similarity between  $L_1$  and  $L_2$  for a set of  $W$  typologies:
  - ▶  $S(L_1, L_2) = \# \text{ Similarities}$
  - ▶  $S(L_1, L_2) = W - \# \text{ Differences}$
  - ▶  $S(L_1, L_2) = \# \text{ Similarities} / \# \text{ Data available for both}$
  - ▶  $S(L_1, L_2) = \# \text{ Similarities} / \# \text{ Similarities} + \# \text{ Differences}$
  - ▶ ‘Relativer Identitätswert’ (Goebel 1984)

# Improvement for Similarities

- Unusual types are more indicative of similarity:
  - ▶ Instead of counting every similarity as '1'
  - ▶ use:  $1 - (\text{fraction having this type})$
  - ▶ 'Gewichteter Identitätswert' (Goebel 1984)
  - ▶ This idea is related to statistical information
  - ▶ use:  $-\log(\text{fraction having this type})$

# Specifying Similarities



# Improving Differences

- $S(L_1, L_2) = \frac{\# \text{ Similarities}}{\# \text{ Similarities} + \# \text{ Differences}}$
- Instead of counting every differences as '1', take into account that some types are more similar to each other than others

LANGUAGE VIEWER

COMPOSER

# WALS the Feature Viewer

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**

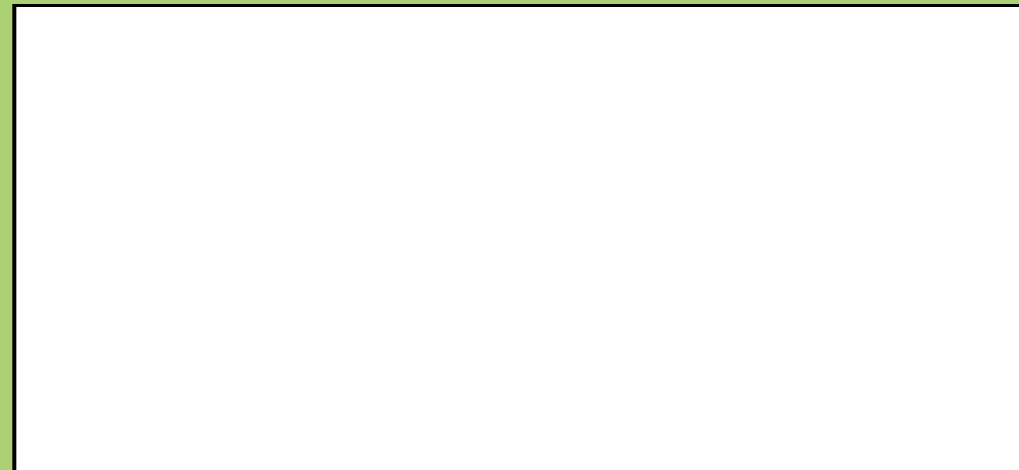
symbol: include: [click to list languages below](#) [no. of lgs : of genera : of families]

<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 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"/>

DESCRIPTION



arrange the languages by

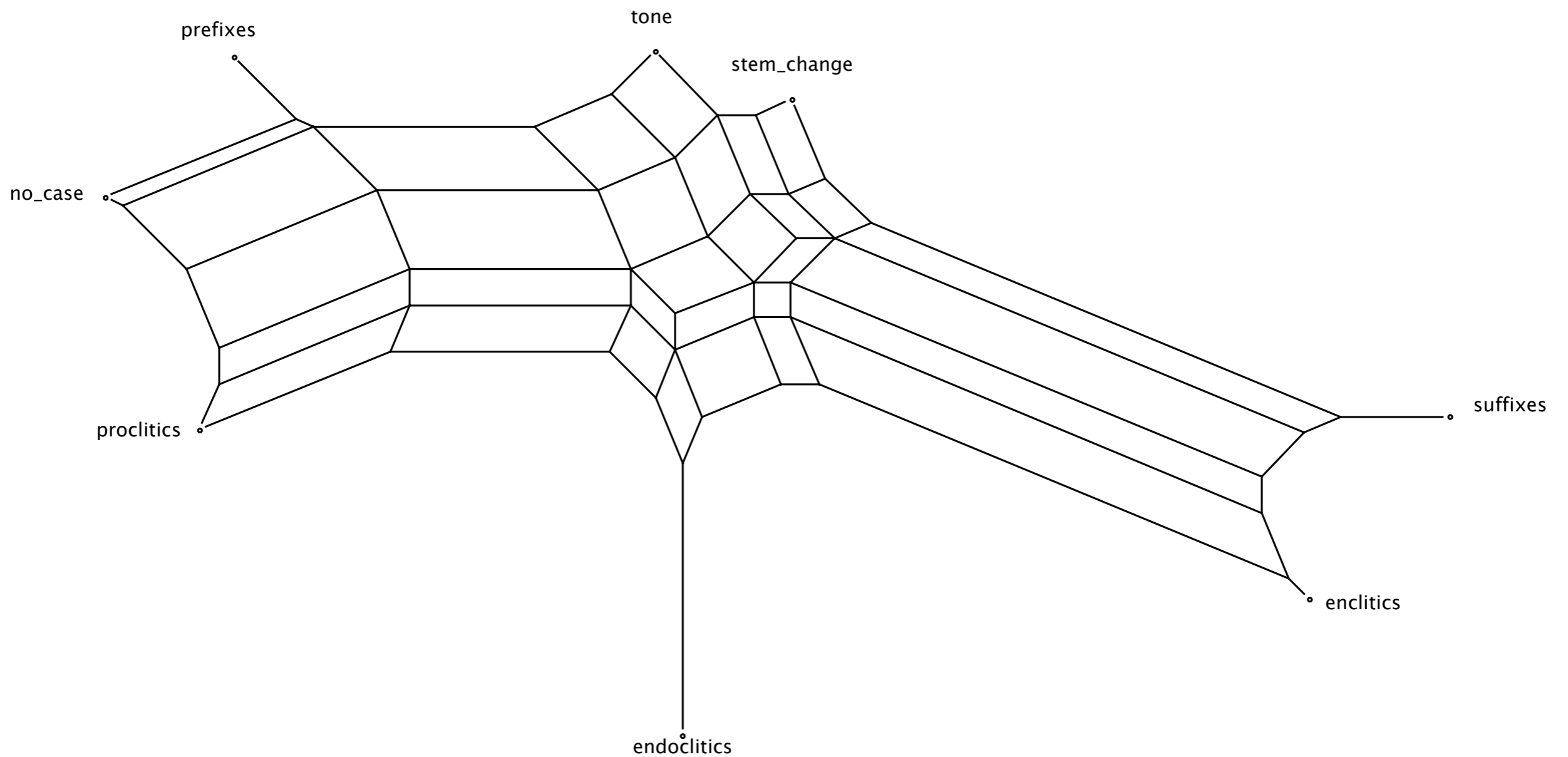
languages

COPY LIST

# Type Similarities

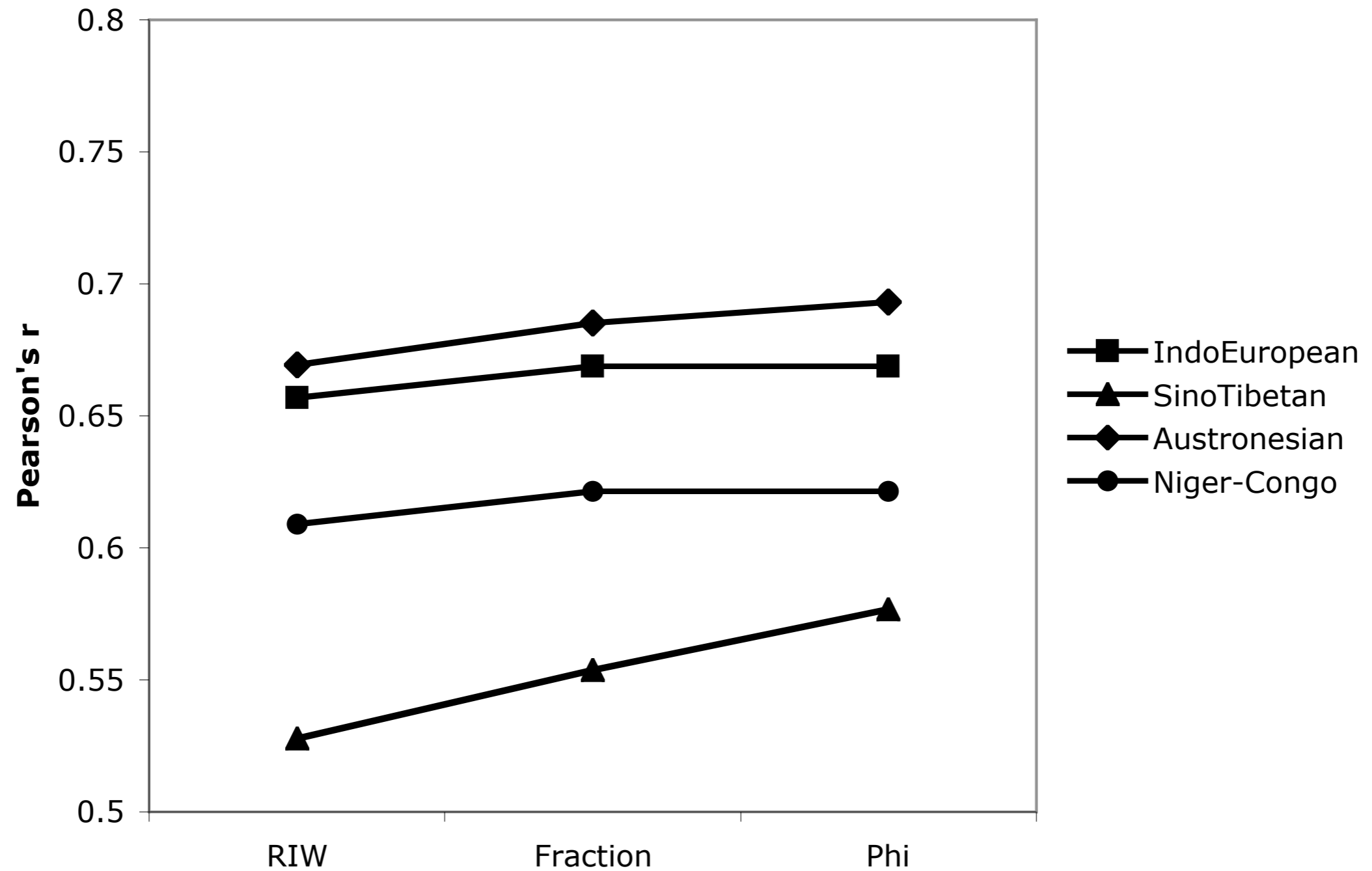
	suffix	prefix	tone	stem	mixed	enclit	proclit	endo	no
suffixes	1	0.04	0.67	0.61	0.36	0.80	0.18	0.26	0.00
prefixes	0.04	1	0.72	0.74	0.88	0.21	0.62	0.48	0.88
tone	0.67	0.72	1	0.77	0.75	0.43	0.52	0.53	0.74
stem_change	0.61	0.74	0.77	1	0.76	0.48	0.53	0.53	0.64
mixed	0.36	0.88	0.75	0.76	1	0.29	0.96	0.52	0.86
enclitics	0.80	0.21	0.43	0.48	0.29	1	0.26	0.43	0.18
proclitics	0.18	0.62	0.52	0.53	0.96	0.26	1	0.52	0.89
endoclitics	0.26	0.48	0.53	0.53	0.52	0.43	0.52	1	0.29
no_case	0.00	0.88	0.74	0.64	0.86	0.18	0.89	0.29	1

# Network of Types





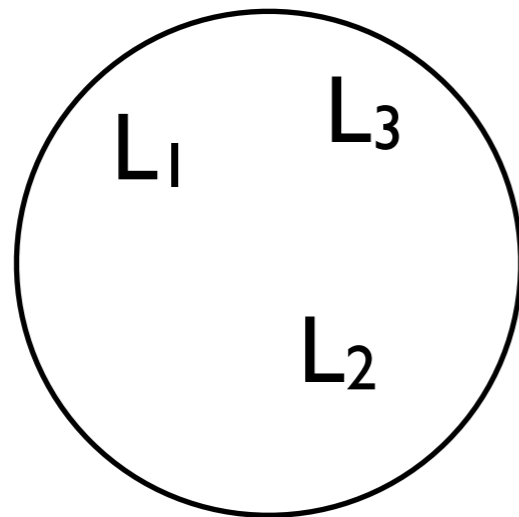
# Specifying Differences



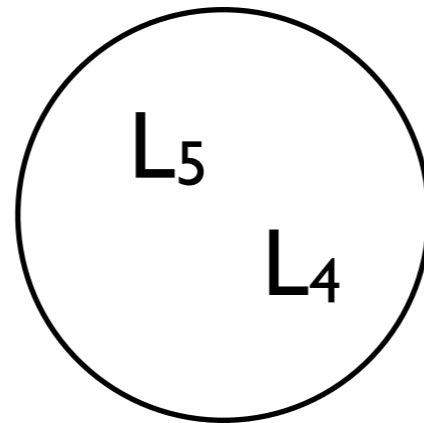
# Further Improvements

- Combine both of these improvement (?)
- Investigate which typologies are most strongly linked to goal at hand:
  - ▶ genealogically stable features
  - ▶ less likely borrowable features
- Make finer grained typologies!

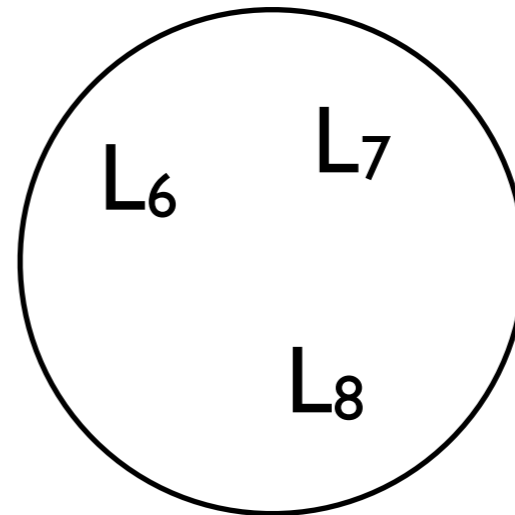
Type A



Type B



Type C











Type A

Type B

Type C

	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>	L <sub>5</sub>	L <sub>6</sub>	L <sub>7</sub>	L <sub>8</sub>	...
L <sub>1</sub>	1	1	1	0	0	0	0	0	
L <sub>2</sub>	1	1	1	0	0	0	0	0	
L <sub>3</sub>	1	1	1	0	0	0	0	0	
L <sub>4</sub>	0	0	0	1	1	0	0	0	
L <sub>5</sub>	0	0	0	1	1	0	0	0	
L <sub>6</sub>	0	0	0	0	0	1	1	1	
L <sub>7</sub>	0	0	0	0	0	1	1	1	
L <sub>8</sub>	0	0	0	0	0	1	1	1	
...									

Undifferentiated Typology



Type A

Type B

Type C

	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>	L <sub>5</sub>	L <sub>6</sub>	L <sub>7</sub>	L <sub>8</sub>	...
L <sub>1</sub>				0.37	0.37	0.28	0.28	0.28	
L <sub>2</sub>				0.37	0.37	0.28	0.28	0.28	
L <sub>3</sub>				0.37	0.37	0.28	0.28	0.28	
L <sub>4</sub>	0.37	0.37	0.37			0.58	0.58	0.58	
L <sub>5</sub>	0.37	0.37	0.37			0.58	0.58	0.58	
L <sub>6</sub>	0.28	0.28	0.28	0.58	0.58				
L <sub>7</sub>	0.28	0.28	0.28	0.58	0.58				
L <sub>8</sub>	0.28	0.28	0.28	0.58	0.58				
...									

Differentiated Typology

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

**‘Deconstructed’ Typology, or a ‘Typology without Types’**



# **Contextually Situating Exemplars**

# Contextually Situated Exemplars ?

- Choose a domain to make a typology
- Choose a (large) set of **situations in context** that fit into the domain
- Investigate these situations in all languages to compare
- Easy way: use existing translations ('parallel texts')

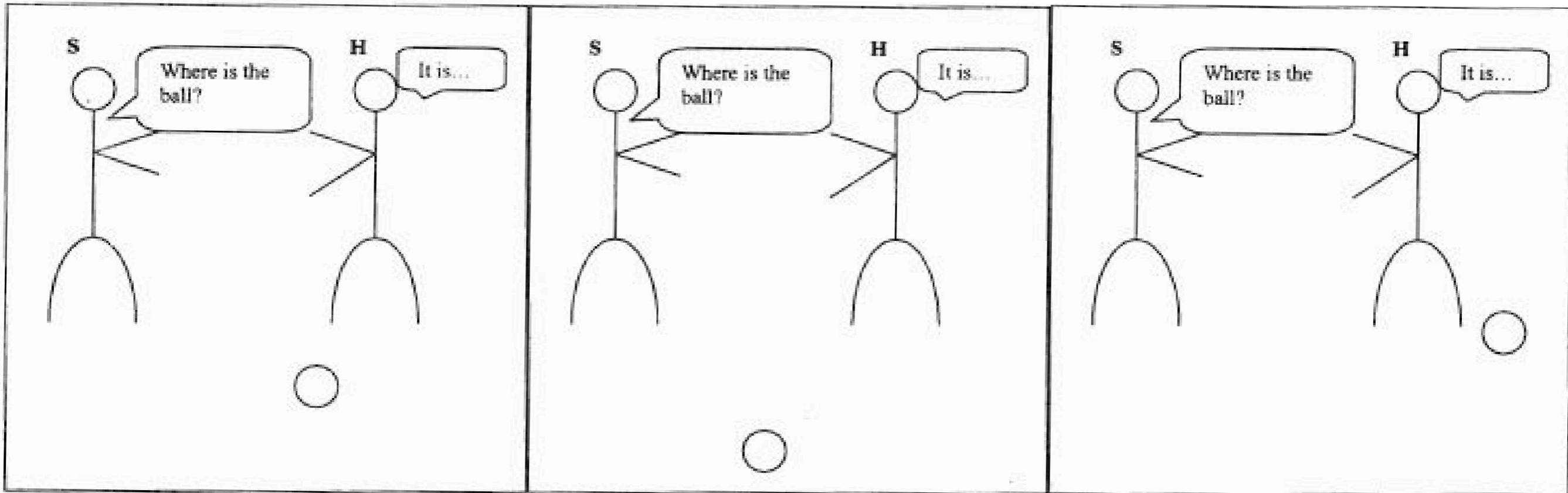
## Appendix

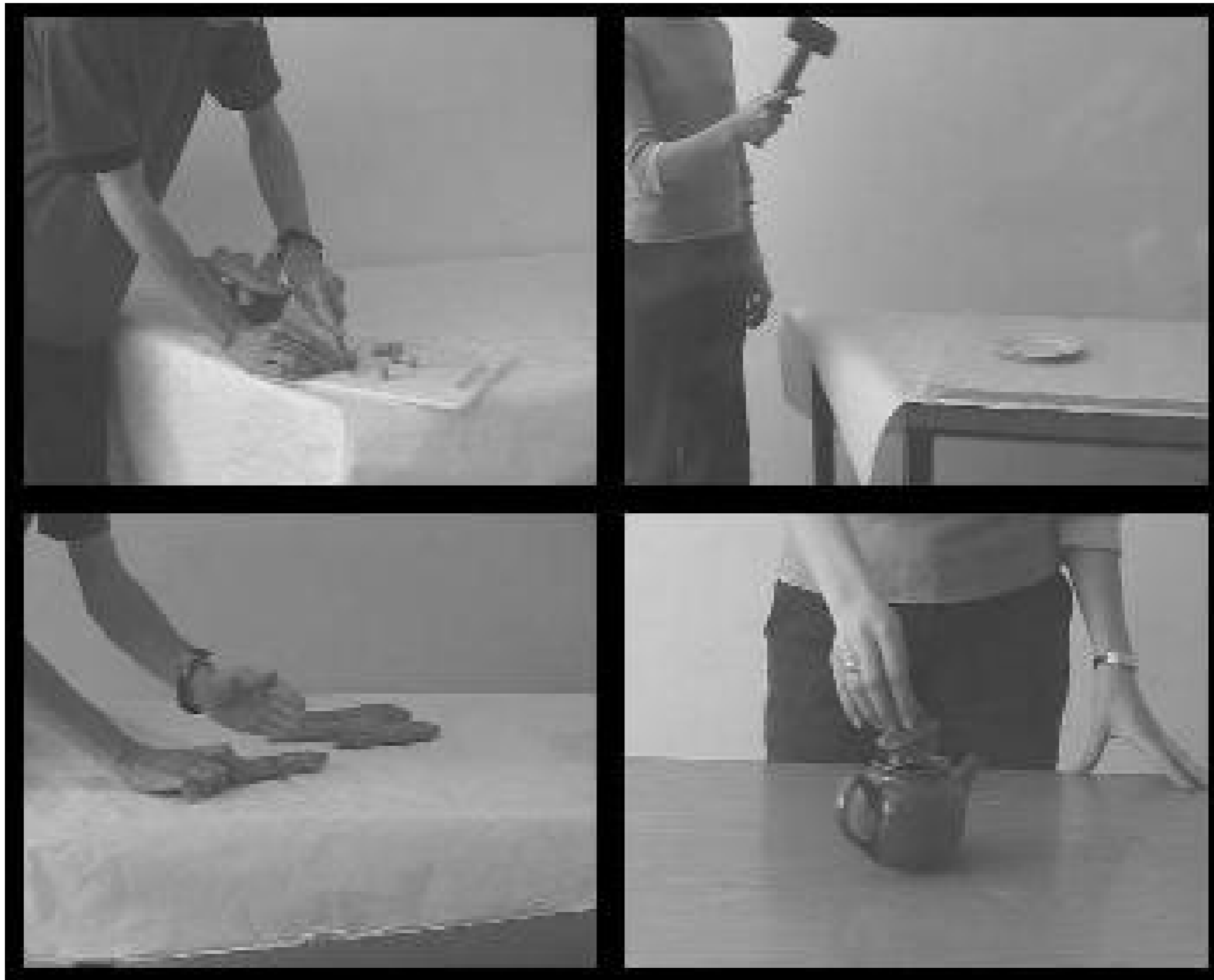
# The TMA questionnaire

Context indications are given within square brackets. Words within parentheses are not to be translated.

### Part A – sentences

- (1) [Standing in front of a house] The house BE BIG
- (2) [Talking about the house in which the speaker lives (the house is out of sight)] The house BE BIG
- (3) [Talking about a house in which the speaker used to live but which has now been torn down] The house BE BIG
- (4) [Talking about a house which the speaker saw for the first time yesterday and doesn't see now:] The house BE BIG
- (5) [Q: What your brother DO right now? (=What activity is he engaged in?) A by someone who can see him] He WRITE letters





Majid, Asifa *et al.* (2004) Event categorization: A crosslinguistic perspective. *Proceedings of AMCSS*, pp. 885-890.

	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



	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

## **Contextually Situated Exemplar**

	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

**Languoid**

	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

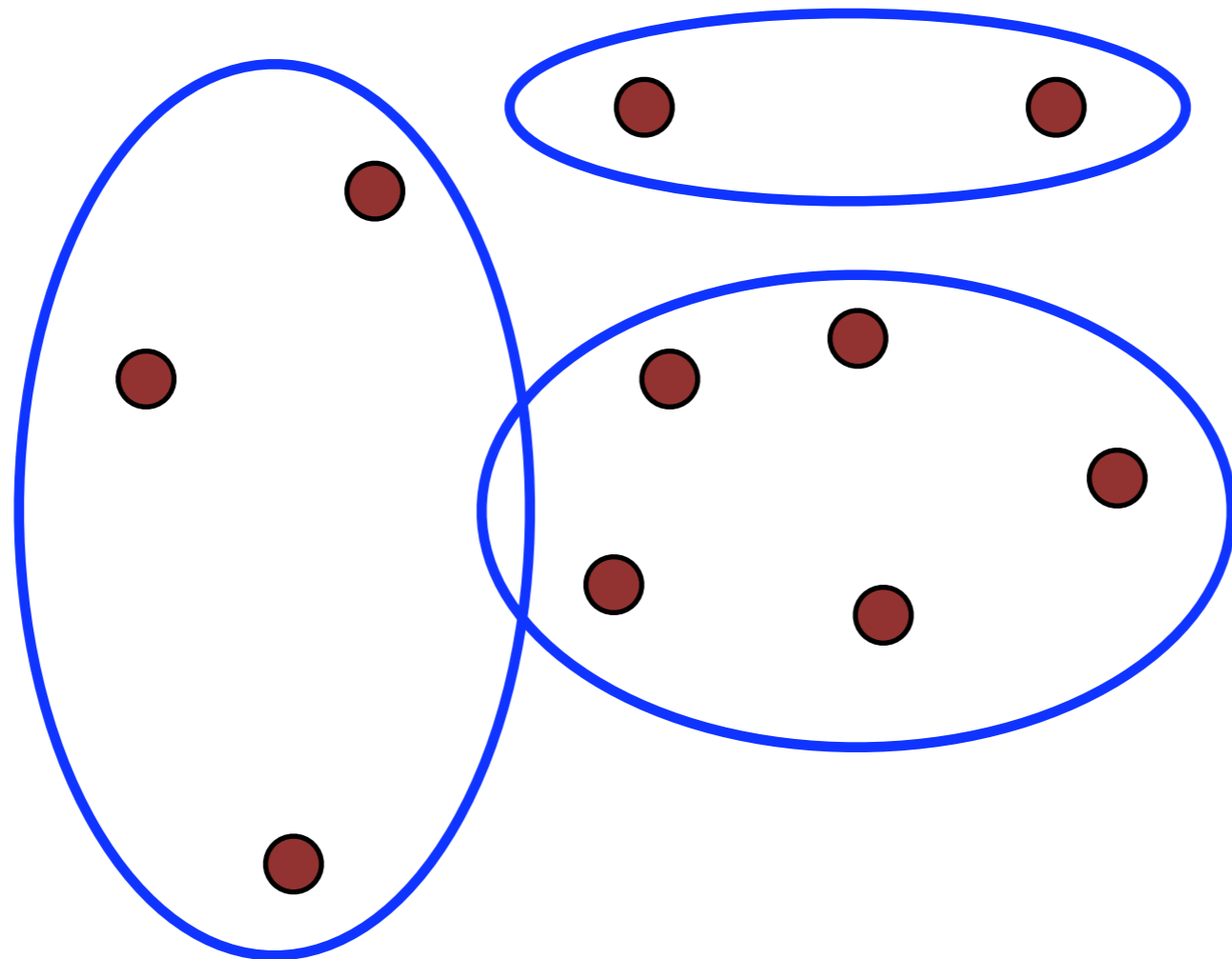
## Language Specific Category

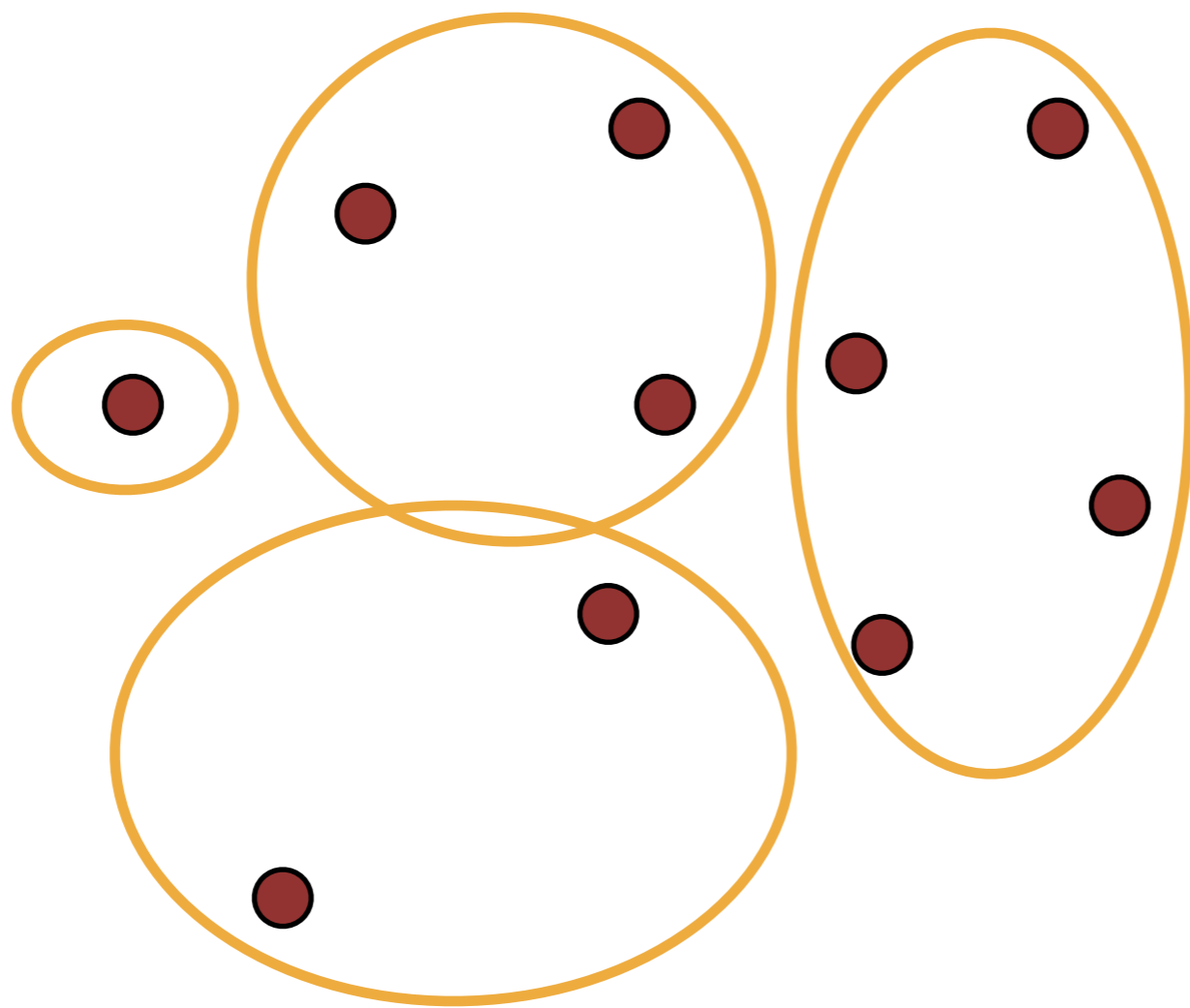
# Using Contextually Situated Exemplars

- A. Compare Exemplars (cf. Semantic Map)
- B. Compare Categories (The Real Thing!)
- C. Compare Linguoids (cf. Typology)

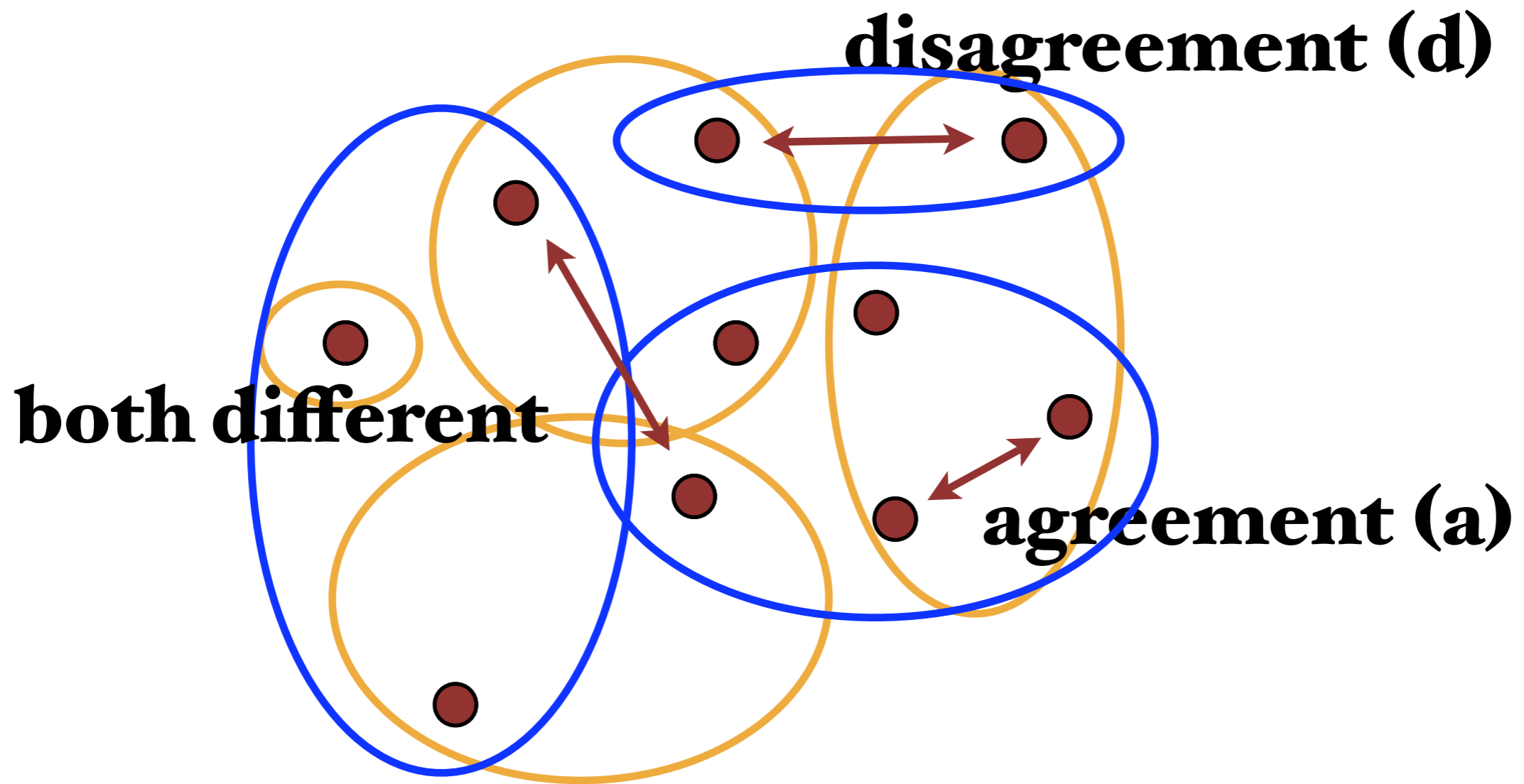
	MRD	LIT	ENG	FRE
1050	sams	eiti	go	aller
1070	sams	eiti	come	venir
1090	sams	eiti	come	venir
1104	lisems	kopti	come	sortir
1105	valgoms	zengti	descend	descendre
1114	–	–	come	se faire entendre
1120	vetjams	varyti	drive	pousser
1140	sams	eiti	come	se rendre
1160	jutams	eiti	walk	marcher

	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





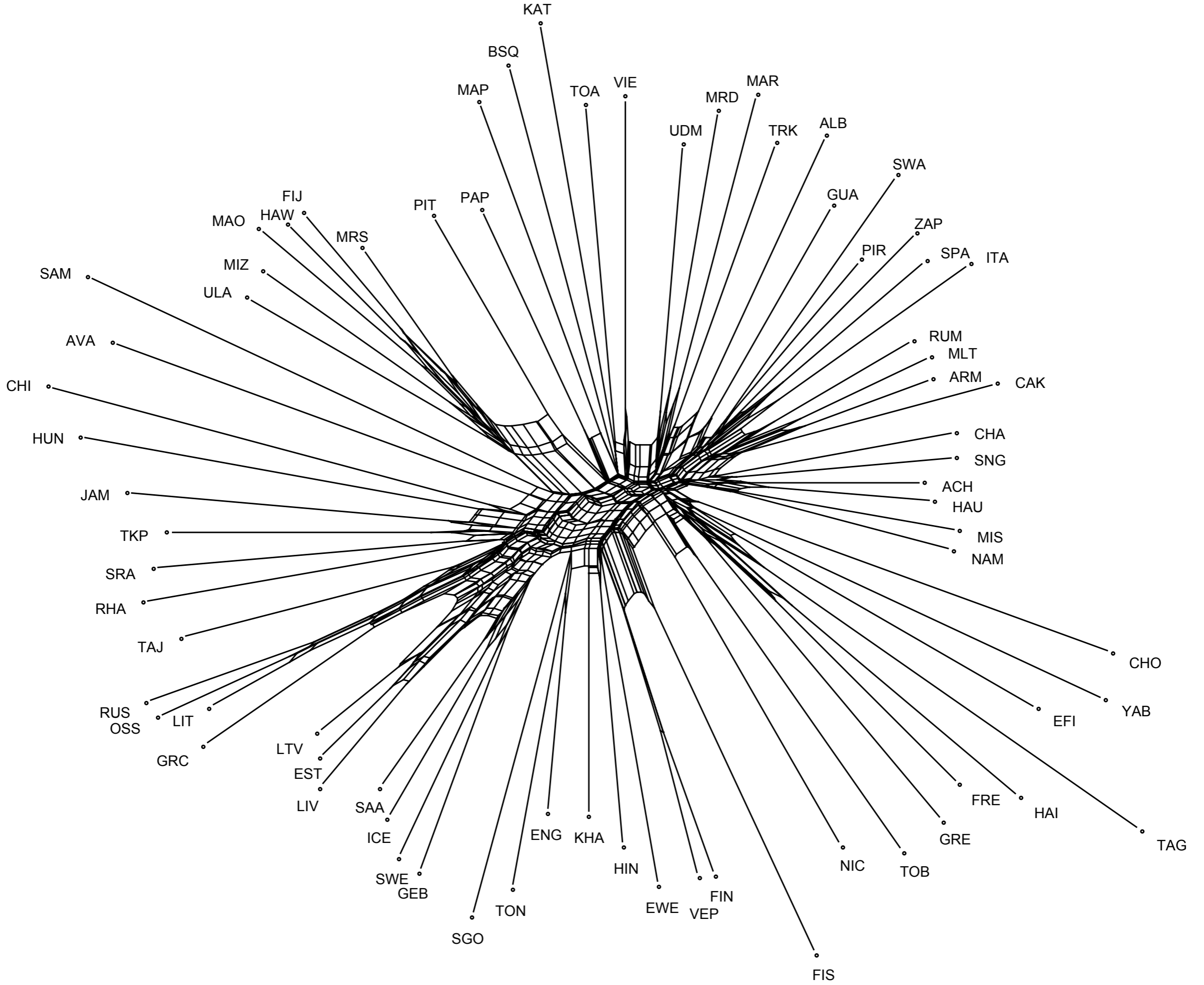




Jaccard similarity:

$$\frac{a}{a+d}$$





# Wordlists

INFORMATION

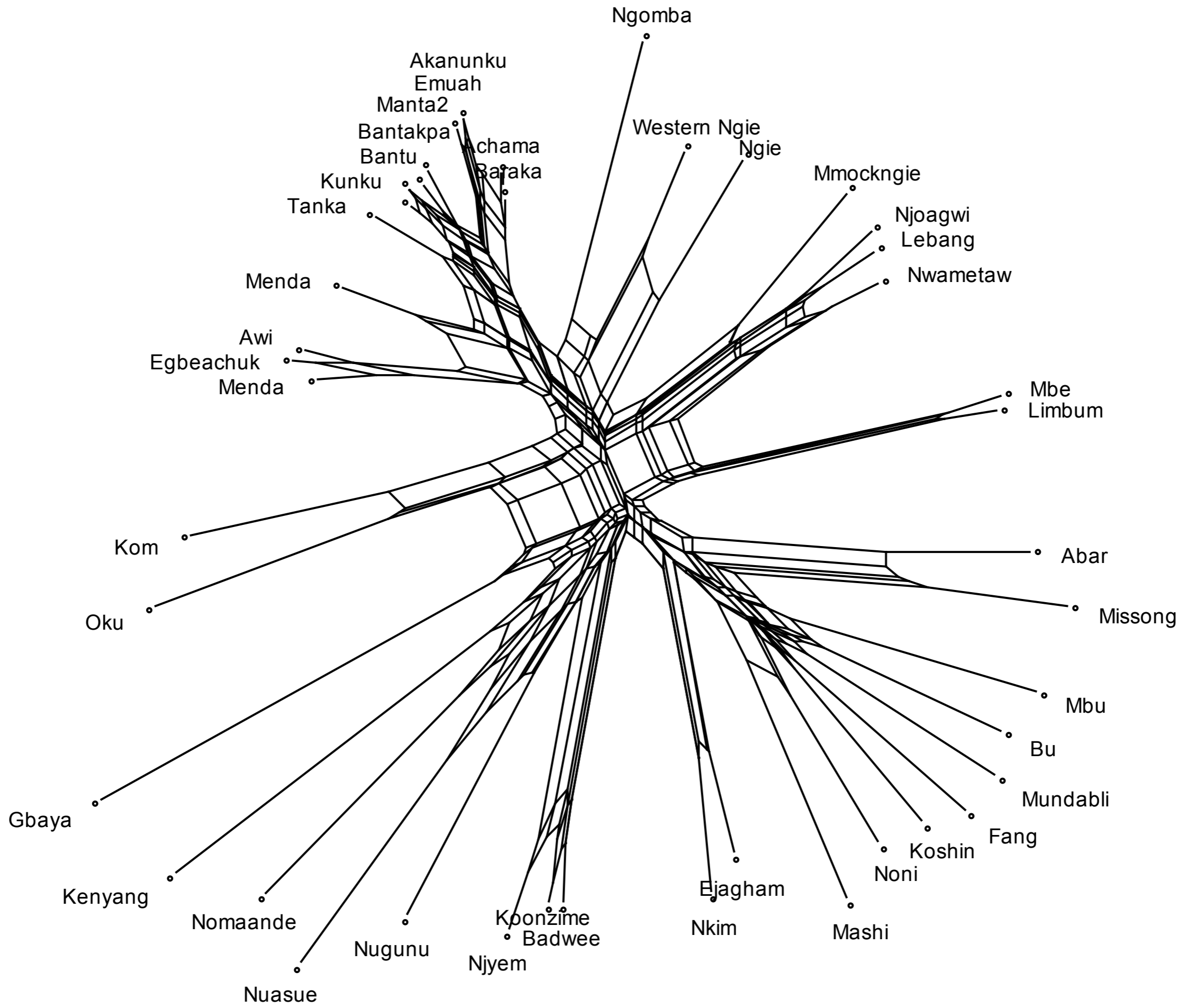
# Reminder!

- Wordlists are **not** only useful for investigating genealogical relations
- Wordlists are just ‘parallel dictionaries’ that allow for quantification of lexical similarity
- The meaning of such similarities is a research question

# Traditional Approach to Wordlist Analysis

- Presence vs. absence of cognates is used for computation of similarity
- Problems:
  - ▶ Judgement of cognacy presumes hypothesis about relationship
  - ▶ Very much of the available information is thrown away

Language	Group			'water'			
Abar	Beboid		a	nj	a		
Missong	Beboid		a	nj	ɛ		
Bu	Beboid			ŋg	ɨ	n	
Mundabli	Beboid			ŋg	i		
Koshin	Beboid			nd	i		
Fang	Beboid			ndz	ia	m	
Mbu'	Beboid			mg	iə	ŋ	
Mashi	Beboid			ngw	ɔ		
Noni	Beboid			j	oo		
Ejagham	Ekoid				á	b	
Nkim	Ekoid			l	ɨ	b	
Kendem	Mamfe		a	n	á	ʔ	
Kenyang	Mamfe	m	a	ny	ɛ	p	
Yambetta	Mbam A40	m	ə	n	í		
Nomaande	Mbam A40	m	e	ny	ïù	f	eà
Nuasue	Mbam A60	m	i	mb	i		
Badwee	Narrow Bantu, A80	m	o	d	ïù	b	eà
Koonzime	Narrow Bantu, A80	m	e	d	ïù	b	eà
Njyem	Narrow Bantu, A80	m	eè	d	ïù	b	oà







MAX-PLANCK-GESELLSCHAFT

**Max Planck Institute  
for Evolutionary Anthropology**