## The analysis of typological data

Michael Cysouw Leipzig Spring School on Linguistics Diversity

## Survey of this course

#### I. Collecting data

- Choosing Languages
- Establishing Types

#### 2. Implicational Universals (and the like)

- The typological tradition
- Statistical view of things
- Dryer's test
- Summing up parameters (don't !)
- 3. Semantic Maps (and other graphs)
  - The typological tradition
  - Taking frequencies into account
- 4. Relationships between languages

#### I. Collecting data

- Investigate worldwide linguistic diversity
- Sample the world's languages
- Classify languages into types
- Any results are statements about actual and not possible human language !
- By sampling, only major types are captured



## Choosing languages

- Tradition: sample from linguistic families
- Indeed: don't take 20 Indo-European languages and 5 other
- Watch out for large areal consistencies !

#### Diversity Sample (from Ö. Dahl, forthcoming)



## Choosing languages

- Tradition: sample from linguistic families
- Indeed: don't take 20 Indo-European languages and 5 other
- Watch out for large areal consistencies !

## Choosing languages

- Tradition: sample from linguistic families
- Indeed: don't take 20 Indo-European languages and 5 other
- Watch out for large areal consistencies !
- Watch out for internal variation in families !

## Establishing types

- Don't group the dissimilar !
- Specify internal structure of types
  - Based on definitional structure of types
  - Based on empirical measures of similarity

#### Undifferentiated Typology



#### Including Similarities



#### Undifferentiated Typology

	Τı	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>	T <sub>6</sub>	<b>T</b> <sub>7</sub>
T							
T <sub>2</sub>		I					
T <sub>3</sub>			I				
<b>T</b> 4				I			
<b>T</b> 5					I		
T <sub>6</sub>						I	
<b>T</b> <sub>7</sub>							I

### Undifferentiated Typology

	Tı	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>	T <sub>6</sub>	<b>T</b> <sub>7</sub>
Tı	I	0	0	0	0	0	0
<b>T</b> <sub>2</sub>	0		0	0	0	0	0
T <sub>3</sub>	0	0	I	0	0	0	0
<b>T</b> 4	0	0	0		0	0	0
<b>T</b> 5	0	0	0	0	I	0	0
T <sub>6</sub>	0	0	0	0	0		0
<b>T</b> <sub>7</sub>	0	0	0	0	0	0	Ι

#### Specifying similarities

	Τı	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>	T <sub>6</sub>	<b>T</b> 7
Tı	I						
T <sub>2</sub>		I					
T <sub>3</sub>			I				
<b>T</b> 4				I			
<b>T</b> 5					I		
T <sub>6</sub>						Ι	
<b>T</b> <sub>7</sub>							I

#### Specifying similarities

	T	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>	T <sub>6</sub>	T <sub>7</sub>
Tı	I	0.37	0.28	0.27	0.51	0.76	0.66
T <sub>2</sub>	0.37		0.58	0.78	0.35	0.51	0.65
T <sub>3</sub>	0.28	0.58		0.6	0.55	0.28	0.67
<b>T</b> 4	0.27	0.78	0.6	I	0	0.58	0.70
<b>T</b> 5	0.51	0.35	0.55	0	I	0	0.68
T <sub>6</sub>	0.76	0.51	0.28	0.58	0	Ι	0.51
<b>T</b> <sub>7</sub>	0.66	0.65	0.67	0.70	0.68	0.51	Ι





## 'Deconstructing' Typology

	Lı	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>	• • •
Lı									
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>									
•••									

	Lı	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>	• • •
Lı									
$L_2$		I							
L <sub>3</sub>			I						
L <sub>4</sub>				I					
$L_5$					I				
$L_6$						I			
L <sub>7</sub>							I		
L <sub>8</sub>								I	
•••									

	Lı	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	I								
$L_2$		Ι							
L <sub>3</sub>			I						
L <sub>4</sub>				I					
$L_5$					I				
$L_6$						I			
L <sub>7</sub>									
L <sub>8</sub>								I	
•••									

	Lı	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	I	I	I						
$L_2$	I	Ι	I						
L <sub>3</sub>	I	I	I						
L <sub>4</sub>				I	I				
L <sub>5</sub>				I	I				
$L_6$						I		I	
L <sub>7</sub>						I		I	
L <sub>8</sub>						I		I	
•••									

	Lı	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				0	0	0	0	0	
L <sub>2</sub>			I	0	0	0	0	0	
L <sub>3</sub>			I	0	0	0	0	0	
L <sub>4</sub>	0	0	0	I	I	0	0	0	
L <sub>5</sub>	0	0	0	I	I	0	0	0	
L <sub>6</sub>	0	0	0	0	0	I	I	I	
L <sub>7</sub>	0	0	0	0	0	I	I	I	
L <sub>8</sub>	0	0	0	0	0	I	I	I	
•••									

#### Undifferentiated Typology

	Lı	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>	$L_5$	L <sub>6</sub>	L <sub>7</sub>	L <sub>8</sub>	• • •
L		I	I	0.37	0.37	0.28	0.28	0.28	
L <sub>2</sub>	I	I	I	0.37	0.37	0.28	0.28	0.28	
L <sub>3</sub>		I	I	0.37	0.37	0.28	0.28	0.28	
L <sub>4</sub>	0.37	0.37	0.37	I		0.58	0.58	0.58	
L <sub>5</sub>	0.37	0.37	0.37	I	I	0.58	0.58	0.58	
$L_6$	0.28	0.28	0.28	0.58	0.58	I		I	
L <sub>7</sub>	0.28	0.28	0.28	0.58	0.58	I		I	
L <sub>8</sub>	0.28	0.28	0.28	0.58	0.58	Ι		I	
• • •									

#### Inter-type similarities

	Lı	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>	•••
Lı									
$L_2$		I							
L <sub>3</sub>			I						
L <sub>4</sub>				I					
L <sub>5</sub>					I				
$L_6$						I			
L <sub>7</sub>							I		
L <sub>8</sub>								I	
•••									

	Lı	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>	L <sub>5</sub>	$L_6$	L <sub>7</sub>	L <sub>8</sub>	• • •
Lı		0.55	0.72	0.31	0.70	0.61	0.50	0.58	
L <sub>2</sub>	0.55	I	0.55	0.31	0.40	0.44	0.31	0.48	
L <sub>3</sub>	0.72	0.55	I	0.29	0.53	0.5 I	0.48	0.60	
L <sub>4</sub>	0.31	0.31	0.29	I	0.38	0.36	0.26	0.27	
L <sub>5</sub>	0.70	0.40	0.53	0.38		0.64	0.5 I	0.46	
L <sub>6</sub>	0.61	0.44	0.51	0.36	0.64	-	0.57	0.43	
L <sub>7</sub>	0.50	0.31	0.48	0.26	0.51	0.57		0.47	
L <sub>8</sub>	0.58	0.48	0.60	0.27	0.46	0.43	0.47		
•••									

#### 'Deconstructed' Typology

## B.Wälchli's data on motion events

- 72 languages
- 335 clauses for each language from Bible
- clauses describing motion events
- here, only the lexical verb used is included
- contextually situated exemplars

	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











dimension 1



# 2. Implicational Universals (and the like)

- The typological tradition
- Statistical view of things
- Dryer's test (with variations and critique)
- Summing up parameters (don't !)

## The typological tradition

- Implicational Universal
- Bidirectional Universal (Equivalence)
- Implicational Hierarchy
- Nested Implicational Universal
# Greenberg (1963)

- Universal 3: Languages with dominant VSO order are always prepositional
- Universal 2: In languages with prepositions, the genitive almost always follows the governing noun, while in languages with postpositions it almost always precedes

	А	В	С	D	
1	+	+	+	+	26
2	_	+	+	+	<b>78</b>
3	_	_	+	+	99
4	_	_	_	+	20
5	_	_	_	_	21
6	+	_	+	+	3
7	_	+	_	+	12
8	_	_	+	_	4
9	+	_	_	+	1
10	_	+	+	_	0
11	+	+	_	+	0
12	+	_	+	_	0
13	_	+	_	_	0
14	+	+	+	_	1
15	+	+	_	_	0
16	+	_	_	_	0
Total +	31	117	211	239	

### Statistical view of things

#### What do typologists say?

Smallest number	Kind of universal	Hypothetical distributions of a 100-language sample												
7	Exceptionless		33	34		26	48		20	60	14	72		
Zero	universal	-	0	33		0	26		0	20	 0	14	-	
	Strong		36	23		31	33		27	41	22	51		
Five tendency	-	5	36	. <u> </u>	5	31		5	27	 5	22	-		
Ten Statistical tendency	en Statistical tendency	Statistical		38	14		33	24		30	30	25	40	
		-	10	38	·	10	33		10	30	 10	25	-	
						25	1.5		21		20			
Fifteen	Maybe something					35 15	35		31 15	31	 28 15	29		
	C						20					~		
Nineteen	Nothing							_	31	19	 27	27	_	
									19	31	19	27		

#### What do statisticians say?

33	34	26	48	20	60	14	72
0	33	0	26	0	20	0	14
	1		I		1		1
36	23	31	33	27	41	22	5
5	36	5	31	5	27	5	22
38	14	33	24	30	30	25	4
10	38	10	33	10	30	10	2
		35	15	31	23	28	2
		15	35	15	31	15	2
				31	19	27	2
				19	31	19	2

Hypothetical distributions of a 100-language sample

Kind of interaction	Very strongly significant	Strongly significant	Significant	No interaction
Fisher's Exact two-tailed	<i>p</i> < 0.000001	<i>p</i> < 0.001	<i>p</i> < 0.05	<i>p</i> > 0.2

	The total number of cases= 97 The smallest value= 10 The smallest marginal= 33
Fisher	Use of * (starred) statistics is advised
Exact	The p for exactly this table= 0.000449
Table:   10 23   43 21	The p-value for the same or a stronger association = 0.000555* The p-value for a stronger association = 0.000105 The mid p-value = 0.000330 The p-value for the same or the reverse association = 0.999894
Calculate Help Fisher Exact This procedure by SISA, 1989,1997,2000.	Two sided p-values for p(O>=E O<=E) p-value= 0.0010737593* (the sum of small p's) p-value= 0.001110 (double the single sided p) Two sided p-value for p(O>E O <e) p-value= 0.0006242531 (the sum of small p's) Two sided p-value mid-p</e) 
	Clear



A windows version of this procedure is available here

## Dryer's test

# Summing up parameters (Don't !)



Sum of Head and Dependent marking: 'complexity':

'... the complexity (Dependent points plus Head points ...) has a roughly normal distribution. Neither zero complexity nor the theoretical maximum complexity of [18] points (9 Head points plus 9 Dependent points ...) occurs. the highest attested complexity is 15, found in only two languages. Figure 4 shows the complexity values attested in my sample. ... The normal distribution and preference for moderate complexity shown in the overall sample are echoed in most ... areas, with high complexity predominating in only two.' (Nichols 1992: 88-89)



Ratio of Dependent and Head points: indicating the relative strength of head or dependent marking in a language.



'... computing the ration of dependent to head marking ... gives us 35 different ratios among the 174 sample languages. Their distribution is shown in figure 1. It is bimodal, with the greatest peaks at the extremes of exclusive head marking (ration of zero since D = 0) and exclusive dependent marking (since H = 0, an actual ratio cannot be computed as it has a zero denominator). The other ratios, whose without zeroes, run from 0.14 (two languages) to 8.00 (one language). The highest frequencies are:

- 0.00 34 languages (radically head marking)
- 0.17 9 languages
- 0.50 8 languages
- 1.00 11 languages
- 2.00 12 languages
- H = 0 19 languages (radically dependent marking)

... The other three frequency peaks suggest that preferred patterns cluster at perceptually simple ratios: two to one, one to one, and one to two. Overall, then , we have a preference for neatness of some sort: polar types, two-to-one ratios and even splits.' (Nichols 1992: 72-73)

ing)

[should be '0.33', MC]

![](_page_47_Figure_0.jpeg)

## 3. Semantic Maps

- Traditional view
  - Extension of Implicational Hierarchy
- Multidimensional Scaling

## Extension of Hierarchy

Guarani:	$\left[ first/second \ prn \ - \ third \ prn \ - \ human \ N \ - \ animate \ N \ - \ inanimate \ N$
Usan:	$\left[ \text{first/second } prn \ - \ \text{third } prn \ - \ \text{human} \ N \ - \ \text{animate} \ N \ - \ \text{inanimate} \ N \right.$
Tiwi:	first/second prn - third prn - human N - animate N - inanimate N
Kharia:	$\left[ first/second \; prn \; - \; third \; prn \; - \; human \; N \; - \; animate \; N \; - \; inanimate \; N \;$
English:	first/second prn - third prn - human N - animate N - inanimate N

Figure 5.1 Semantic maps of plural inflection in various languages

![](_page_50_Figure_0.jpeg)

![](_page_50_Figure_1.jpeg)

## Examples of indefinite Pronouns

![](_page_51_Figure_1.jpeg)

![](_page_51_Figure_2.jpeg)

#### Evaluation

- Good: only 10 out of 45 (= 9x8/2) possible lines needed
- But: one line is indeterminated
- But: 105 groups predicted, though only 39 attested
- But: frequencies do not play a role

### Multidimensional scaling

![](_page_53_Figure_1.jpeg)

# Semantic Map of Person Marking

![](_page_54_Figure_1.jpeg)

### Frequencies included

![](_page_55_Picture_1.jpeg)

### Multidimensional scaling

![](_page_56_Figure_1.jpeg)

dimension 1

![](_page_57_Figure_0.jpeg)

dimensions

# B.Wälchli's data on motion events

- 72 languages
- 335 clauses for each language from Bible
- clauses describing motion events
- here, only the lexical verb used is included

	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

![](_page_60_Figure_0.jpeg)

dimension 1

venir

![](_page_61_Figure_1.jpeg)

dimension 1

aller

![](_page_62_Figure_1.jpeg)

dimension 1

#### screeplot

![](_page_63_Figure_1.jpeg)

dimensions

# 4. Relationships between languages

- Universality: linguistic diversity interpreted as a-historical generalisations
- Contingency: linguistic diversity interpreted as result of historical processes

#### Using Typological Data for Genealogical Investigations

![](_page_65_Figure_1.jpeg)

(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.

Kairiru

Dunn, Michael, Angela Terrill, Ger Reesink, Robert A. Foley & Stephen C. Levinson. 2005. Structural Phylogenetics and the Reconstruction of Ancient Language History. Science 309: 2072-2075.

![](_page_66_Figure_0.jpeg)

#### 

Dunn et al. tree based on typological

**,**•<sup>5</sup>

# Measuring typological stability

- Given an accepted tree, with many languages sampled from this tree
- how good does a typological feature predict this tree
- Energy-based measurement of fit between a dataset and a tree (work by Mihai Albu)
- Take a large set of random trees, and determine how good the "real" tree fits

#### Distribution of fits of all 125 features

(Too) many good fits!

![](_page_68_Figure_2.jpeg)

#### Oceania

![](_page_69_Picture_1.jpeg)

#### NNet of typological distances

-10.0

![](_page_70_Figure_2.jpeg)

#### NNet of typological distances

![](_page_71_Figure_1.jpeg)

![](_page_71_Figure_2.jpeg)






#### Typology/geography correlation



Mantel test p = .349

geographic distance (kilometers)

#### Correlation for selection only



Mantel test p = .001

geographic distance (kilometers)

#### When does correlation improve?

	Pearson's r
Nothing removed	.035

#### When does correlation improve?

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

#### Investigation typology/geography relation



geographic distance (kilometers)

#### Linguistically 'too similar'



#### Linguistically 'too similar'



# Summary

- Typology is correlated to genealogy
- but: typology is also correlated to geography
- When removing the (genealogically related) Austronesian languages, the typology/ geography correlation improves
- The language-pairs that are typologically more similar than expected from geography are genealogically related

# Towards an interpretation

- In longterm static (areal) interaction typological features diffuse individually, leading to regular geographical clines
- In relatively recent (genealogical) spread bundles of features 'move' together, leading to stronger similarities as expected from geography

## Eurasia





#### Typology/geography correlation







#### Remove 'worst-fitting' languages



geographic distance (kilometers)

#### Remove 'worst-fitting' languages



geographic distance (kilometers)





# Some interpretation

- Turkish and Hungarian are cases of relatively recent movement of whole languages
- But Lezgian (probably) not
- Link Hindi-Hungarian is unclear, and Burushaski-Basque is too cranky a speculation
- Chukchi, Georgian, Abkhaz simply unrelated, both genealogical and areal