

Syntactic measuring of linguistic distances

Giuseppe Longobardi
longbard@units.it

(Laboratorio di Linguistica & Antropologia cognitiva, Trieste)

Vienna, October 19th 2010

Darwin's challenge

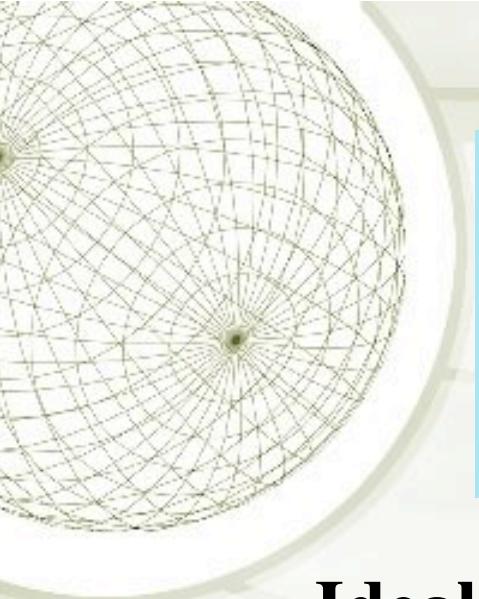
[The Origin of Species, ch. 14]

“If we possessed a perfect pedigree of mankind, a genealogical arrangement of the races of man would afford the best classification of the various languages now spoken throughout the world; and if all extinct languages, and all intermediate and slowly changing dialects, were to be included, such an arrangement would be the only possible one”

Cavalli Sforza, Menozzi, and Piazza (1994)

The History and Geography of Human Genes, p.18

**“We believe that the major breakthrough in the
study of human variation has been the introduction of
genetic markers, which are strictly inherited and
basically immune to the problem of rapid changes
induced by the environment”**



Phylogenetic reconstruction and comparative linguistics

Ideal properties of taxonomic characters for linguistic phylogenies:

Universality, discreteness (useful for clustering in relative taxonomies),

improbable similarity (to prove absolute relatedness against chance similarity):

Classical Comparative Method

Based on **lexical** arbitrariness

comparanda = patently beyond chance probability

(e.g. ‘regular/recurrent’ sound correspondences)

Solid phylogenetic conclusions and surprising etymologies (e.g.
full/pieno)

BUT

so improbable similarities = inevitably rare across languages

OK for Italian/English,

not

Italian/Japanese (long-range comparison)

Mass (multilateral) comparison

[Greenberg's life-long attempt toward a global classification]

Also based on **lexical** arbitrariness

**Relies on universal meaning lists
(e.g. Swadesh lists)**

**Hence, in principle applicable to any set
of languages**

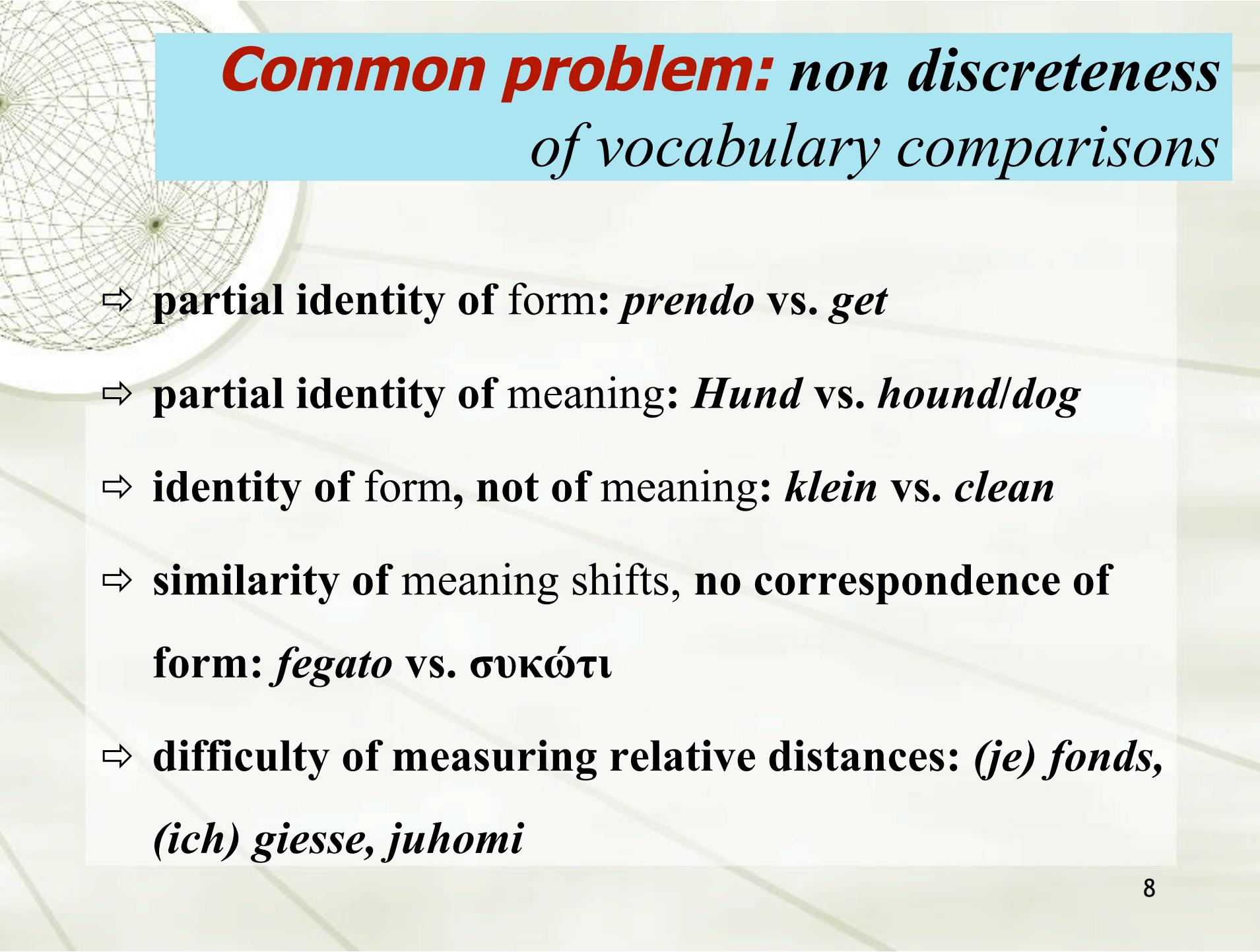
but...

Mass (multilateral) comparison

**supposed cognates are ‘identified’
on mere word resemblance,
not on any precise and improbable phenomena**



not safe from chance similarity



Common problem: non discreteness of vocabulary comparisons

- ⇒ partial identity of form: *prendo* vs. *get*
- ⇒ partial identity of meaning: *Hund* vs. *hound/dog*
- ⇒ identity of form, not of meaning: *klein* vs. *clean*
- ⇒ similarity of meaning shifts, no correspondence of form: *fegato* vs. *συκώτι*
- ⇒ difficulty of measuring relative distances: *(je) fonds*,
(ich) giesse, juhomi

Population Genetics

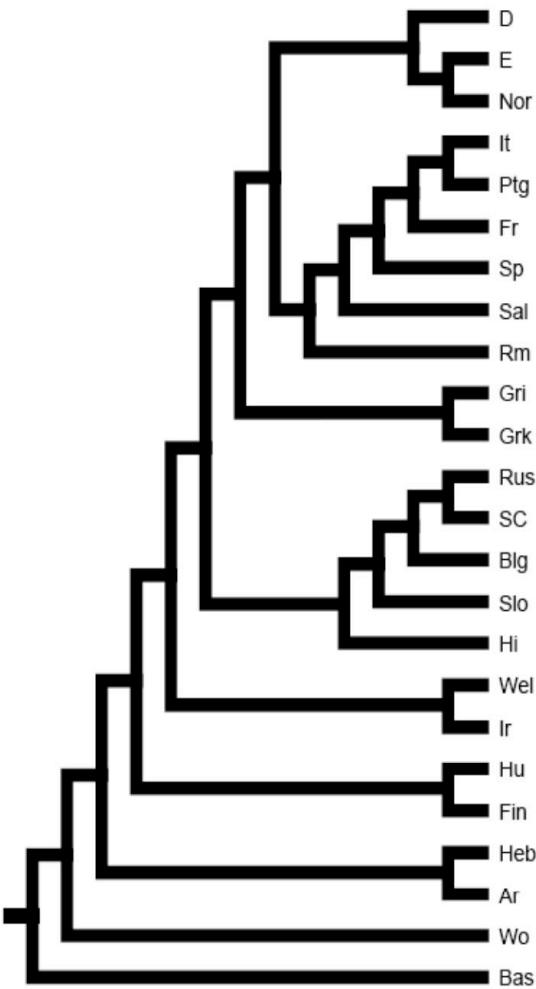
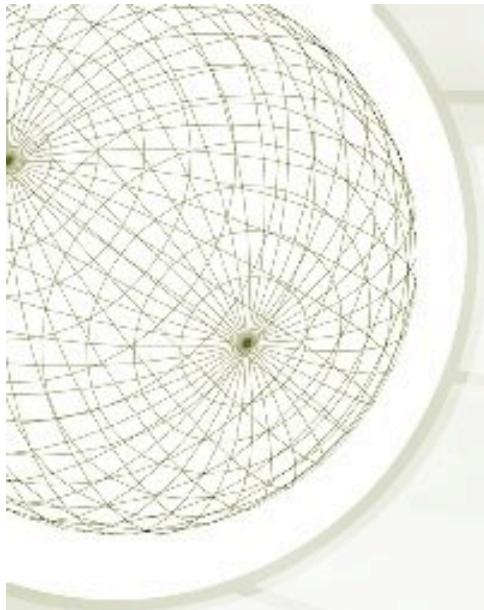
- ⇒ *Comparanda are drawn from a finite universal list of discrete biological options* (alleles of genetic polymorphisms)
- ⇒ **Comparison is safe w.r.t. chance and universal possible** (measuring of taxonomic distances possible between any two populations)

Parametric Comparison Method (PCM)

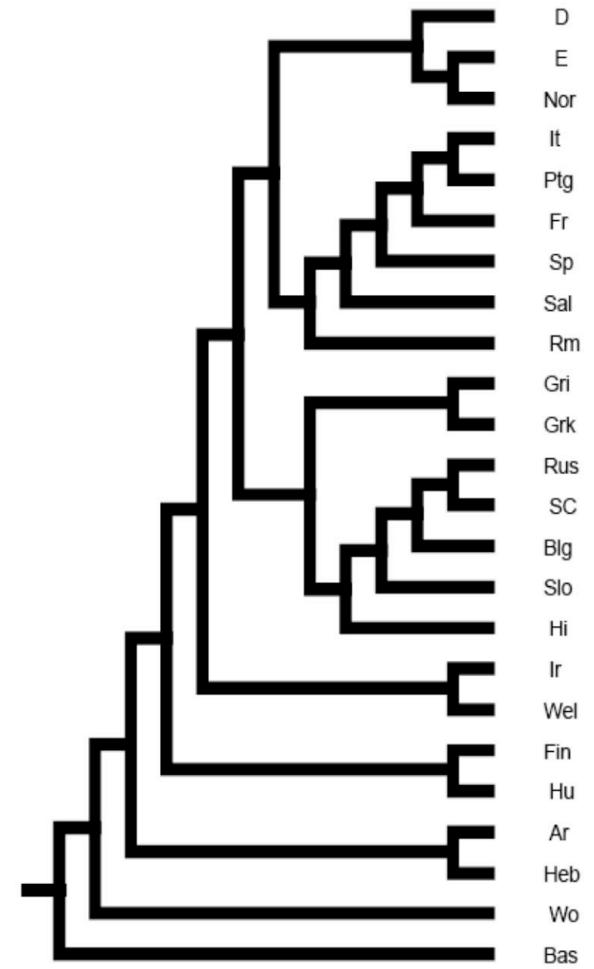
Longobardi (2003), Longobardi & Guardiano (2009):

Lexically blind method of **syntactic** comparison

- ⇒ **universally applicable**
- a. mathematically **discrete** results
- c. choice of comparanda: **safe from chance**
- d. parameter settings: virtually **immune** from natural selection and conscious cultural choice



Kitsch



UPGMA

Pairs	Distance
SC-Rus	0,0286
It-Ptg	0,0377
It-Fr	0,0392
It-Sal	0,06
Ir-Wel	0,0698
Sp-Ptg	0,0755
Fr-Ptg	0,0784
Sal-Ptg	0,08
Blg-SC	0,0882
Rus-Slo	0,0882
E-Nor	0,093
Sal-Fr	0,104
It-Rm	0,106
Ptg-Rm	0,106
D-Nor	0,109
It-Sp	0,113
E-D	0,114
Sp-Fr	0,118
Blg-Rus	0,118
SC-Slo	0,118
Gri-SC	0,121
Sal-Gri	0,122
It-Gri	0,128
Sal-Rm	0,128
Sp-Rm	0,128
Ir-Slo	0,129
It-D	0,13
Sal-D	0,13
Ptg-D	0,13
Rm-Gri	0,13
Gri-Grk	0,13
Rm-Nor	0,133
It-Nor	0,136

Sal-Nor	0,136
Ptg-Nor	0,136
Grk-SC	0,147
Blg-Slo	0,147
Rus-Hi	0,148
Ptg-Gri	0,149
Sal-Slo	0,152
Gri-Rus	0,152
Fr-Rm	0,156
Gri-D	0,156
It-E	0,159
Fr-D	0,159
Ptg-E	0,159
Gri-Ir	0,159
Heb-Ar	0,159
Sal-Sp	0,16
Gri-Blg	0,163
Fr-Nor	0,167
Hu-Fin	0,167
D-Slo	0,171
It-Blg	0,174
Sp-D	0,174
Ptg-Blg	0,174
Grk-Rus	0,176
Sal-Blg	0,178
Fr-Gri	0,178
Nor-Blg	0,178
Hi-Slo	0,179
Sp-Nor	0,182
Gri-Slo	0,182
Grk-Slo	0,182
SC-Hi	0,185
Sal-E	0,186
Fr-E	0,186
Sal-Rus	0,188

Sp-Ir	0,19
Gri-Nor	0,19
D-Ir	0,19
Blg-Heb	0,195
It-Hi	0,2
Sal-Hi	0,2
Ptg-Hi	0,2
Rm-D	0,2
Rm-Blg	0,2
Gri-Hu	0,2
D-Hi	0,2
Nor-Hi	0,2
Rus-Ir	0,2
Sal-Ir	0,205
Sp-E	0,205
Rm-Grk	0,205
Grk-Blg	0,205
E-Heb	0,205
Nor-Heb	0,205
D-Rus	0,206
Sal-Fin	0,207
Rm-Hi	0,207
Blg-Hi	0,207
It-Heb	0,209
It-Rus	0,212
It-Slo	0,212
Sp-Slo	0,212
Ptg-Rus	0,212
Ptg-Slo	0,212
Rm-Rus	0,212
Nor-Slo	0,212
It-Grk	0,213
Sp-Gri	0,213
Sp-Grk	0,213
It-Ir	0,214

Ptg-Ir	0,214
Rm-E	0,214
Gri-Hi	0,214
Sp-Blg	0,217
Sal-SC	0,219
Gri-Wo	0,219
Gri-E	0,22
D-Wel	0,22
Fr-Ir	0,225
Grk-Ir	0,225
Fr-Blg	0,227
D-Blg	0,227
Rm-Heb	0,231
Rm-Hu	0,231
E-Wel	0,231
Nor-Ir	0,231
Ir-Heb	0,231
Fin-Hi	0,231
Ptg-Heb	0,233
E-Blg	0,233
Nor-Fin	0,233
SC-Ir	0,233
Ptg-Grk	0,234
D-SC	0,235
Grk-Ar	0,238
Ir-Hi	0,24
It-Fin	0,241
Ptg-Fin	0,241
Rm-Fin	0,241
Gri-Fin	0,241
It-SC	0,242
Ptg-SC	0,242
Rm-SC	0,242
Rm-Slo	0,242
E-Ir	0,243

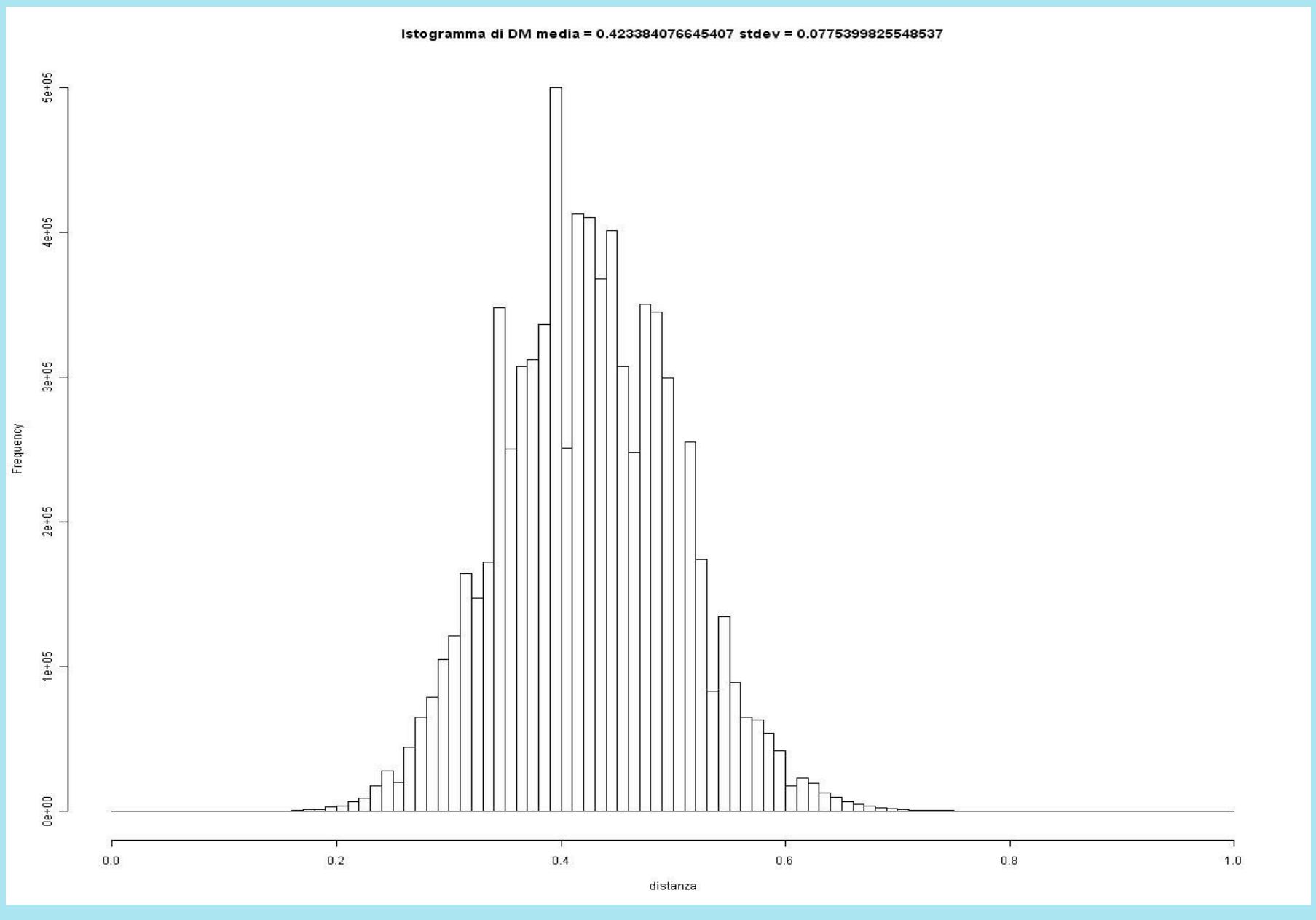
It-Hu	0,244
Sal-Grk	0,244
Sal-Hu	0,244
Rm-Ir	0,244
Gri-Heb	0,244
Sal-Heb	0,25
Fr-Hi	0,25
Grk-Hi	0,25
E-Hi	0,25
D-Heb	0,25
D-Hu	0,25
Nor-Rus	0,25
Sp-Heb	0,256
Gri-Wel	0,256
Fr-Rus	0,258
Fr-Slo	0,258
D-Fin	0,258
Wel-Slo	0,258
Fr-Fin	0,259
Hu-Hi	0,259
E-Hu	0,263
Nor-Hu	0,263
Sp-Hi	0,267
Fr-Grk	0,267
Grk-D	0,267
Rus-Heb	0,267
Fr-Heb	0,268
Ptg-Hu	0,268
Grk-Heb	0,268
Heb-Hu	0,27
Sp-Rus	0,273
Sp-Ar	0,273
Grk-Hu	0,275
Blg-Hu	0,275
Heb-Hi	0,28

E-Rus	0,281
Nor-SC	0,281
Fr-Hu	0,282
Blg-Ir	0,282
Sp-Wel	0,286
Ir-Hu	0,286
Wel-Hu	0,286
Wo-Hu	0,286
Nor-Wel	0,289
Ir-Ar	0,289
Fr-SC	0,29
Fin-Slo	0,29
Fin-Bas	0,292
Sp-Hu	0,293
It-Ar	0,295
Sal-Wel	0,295
Ir-Fin	0,296
Heb-Fin	0,296
Blg-Fin	0,3
SC-Heb	0,3
SC-Fin	0,3
Rus-Fin	0,3
Heb-Slo	0,3
Sp-SC	0,303
It-Wel	0,31
Sp-Fin	0,31
Ptg-Wel	0,31
Grk-Nor	0,31
E-SC	0,312
E-Slo	0,312
Fr-Wel	0,317
Rm-Wel	0,317
Ptg-Ar	0,318
SC-Hu	0,323
Rus-Hu	0,323

Rm-Ar	0,333
Grk-E	0,333
Grk-Wel	0,333
Grk-Wo	0,333
Nor-Ar	0,333
SC-Wo	0,333
Rus-Wel	0,333
Ar-Slo	0,333
Blg-Ar	0,341
Wel-Heb	0,342
E-Fin	0,345
Hi-Bas	0,348
Ar-Hu	0,351
Nor-Bas	0,355
Hu-Slo	0,355
Fr-Ar	0,357
E-Ar	0,359
Bas-Slo	0,36
Sal-Ar	0,364
Sal-Wo	0,364
Gri-Ar	0,366
Rm-Wo	0,367
SC-Wel	0,367
Wel-Fin	0,37
D-Ar	0,375
Rus-Wo	0,375
Wel-Ar	0,378
Grk-Fin	0,379
Blg-Wel	0,385
Wel-Hi	0,385
Ar-Hi	0,385
Rm-Bas	0,387
E-Bas	0,387
D-Wo	0,387
Blg-Wo	0,387

It-Wo	0,394
Ptg-Wo	0,394
Rus-Ar	0,4
Wo-Fin	0,4
Wo-Hi	0,4
Wo-Bas	0,4
Ar-Fin	0,407
Sal-Bas	0,412
Ir-Wo	0,414
D-Bas	0,419
Sp-Wo	0,424
Sp-Bas	0,424
Fr-Wo	0,424
Ptg-Bas	0,424
Hu-Bas	0,429
Gri-Bas	0,433
SC-Ar	0,433
Wo-Slo	0,435
Blg-Bas	0,438
It-Bas	0,455
Rus-Bas	0,458
E-Wo	0,467
Heb-Wo	0,469
Nor-Wo	0,483
Wel-Wo	0,483
Fr-Bas	0,485
SC-Bas	0,5
Heb-Bas	0,5
Ar-Wo	0,5
Wel-Bas	0,519
Ir-Bas	0,52
Grk-Bas	0,533
Ar-Bas	0,533

From Longobardi, Bortolussi, Guardiano, Sgarro: in prep.





'Anti-Babelic principle'

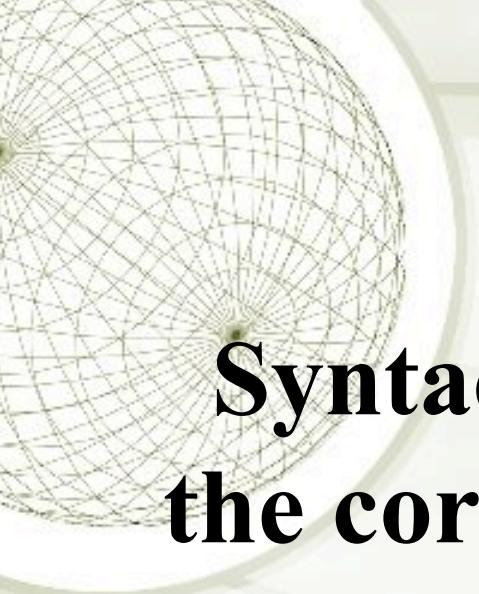
(Guardiano&Longobardi 2005)

*Similarities among languages can be due
either to **historical cause**
(relatedness/convergence) or to **chance**,
differences can only be due to **chance**
(no one ever made grammars diverge on purpose!)*

SYNTACTIC DISTANCES

	Bulgarian	English	French	German	Greek	Hindi	Irish	Italian	Norwegian	Portuguese	Rumanian	Russian	Serbocroat	Slovenian	Spanish	Welsh	
Bulgarian	0	0,233	0,209	0,227	0,205	0,207	0,282	0,156	0,178	0,156	0,2	0,118	0,0882	0,147	0,222	0,385	Bulgarian
English	0,772	0	0,186	0,114	0,333	0,25	0,243	0,159	0,093	0,159	0,214	0,281	0,312	0,312	0,205	0,231	English
French	0,791	0,764	0	0,159	0,267	0,25	0,225	0,04	0,167	0,08	0,156	0,258	0,29	0,258	0,1	0,317	French
German	0,769	0,422	0,756	0	0,267	0,2	0,19	0,13	0,109	0,13	0,2	0,206	0,235	0,171	0,174	0,22	German
Greek	0,811	0,838	0,843	0,812	0	0,25	0,225	0,213	0,31	0,234	0,205	0,176	0,147	0,182	0,213	0,333	Greek
Hindi	0,801	0,854	0,824	0,853	0,874	0	0,24	0,2	0,2	0,2	0,207	0,148	0,185	0,179	0,267	0,385	Hindi
Irish	0,818	0,817	0,812	0,806	0,859	0,878	0	0,214	0,231	0,214	0,244	0,2	0,233	0,129	0,19	0,0698	Irish
Italian	0,769	0,753	0,197	0,735	0,822	0,818	0,8	0	0,136	0,0385	0,106	0,212	0,242	0,212	0,0962	0,31	Italian
Norwegian	0,773	0,452	0,77	0,367	0,821	0,852	0,836	0,754	0	0,136	0,133	0,25	0,281	0,212	0,182	0,289	Norwegian
Portuguese	0,781	0,76	0,291	0,753	0,833	0,813	0,817	0,227	0,761	0	0,106	0,212	0,242	0,212	0,0577	0,31	Portuguese
Rumanian	0,798	0,773	0,421	0,751	0,843	0,827	0,837	0,34	0,786	0,371	0	0,212	0,242	0,242	0,128	0,317	Rumanian
Russian	0,365	0,758	0,778	0,755	0,832	0,8	0,782	0,761	0,758	0,773	0,781	0	0,0286	0,0882	0,273	0,333	Russian
Serbocroat	0,291	0,766	0,772	0,764	0,821	0,805	0,796	0,755	0,772	0,766	0,778	0,325	0	0,118	0,303	0,367	Serbocroat
Slovenian	0,385	0,751	0,782	0,733	0,821	0,8	0,809	0,76	0,762	0,781	0,79	0,386	0,316	0	0,212	0,258	Slovenian
Spanish	0,782	0,76	0,266	0,747	0,833	0,819	0,805	0,212	0,761	0,126	0,406	0,769	0,768	0,772	0	0,286	Spanish
Welsh	0,838	0,841	0,81	0,82	0,867	0,876	0,645	0,793	0,849	0,804	0,812	0,818	0,821	0,838	0,813	0	Welsh
	Bulgarian	English	French	German	Greek	Hindi	Irish	Italian	Norwegian	Portuguese	Rumanian	Russian	Serbocroat	Slovenian	Spanish	Welsh	

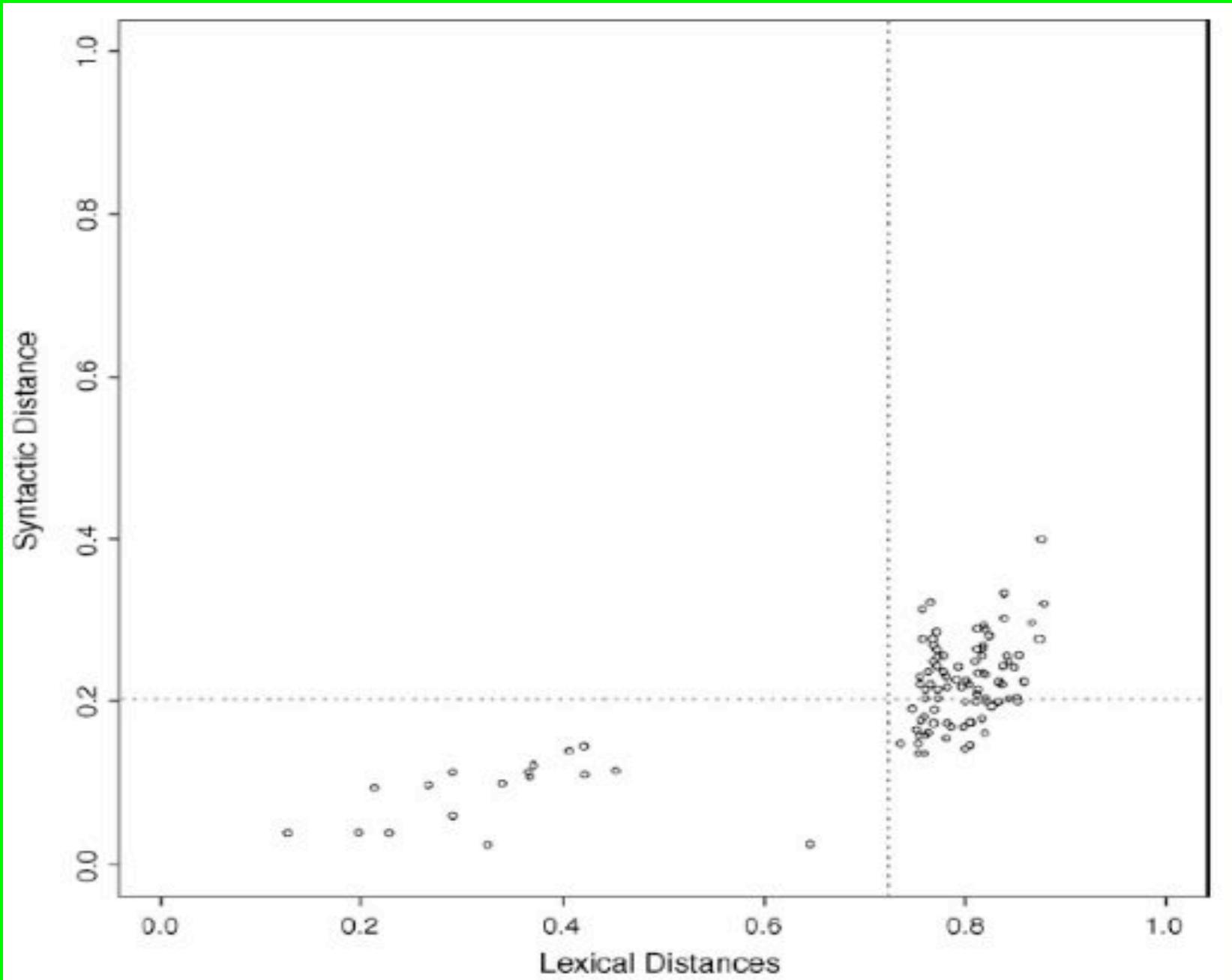
LEXICAL DISTANCES



**Syntactic distances are about 1/4
the corresponding lexical distances
(attended ratio 2/3)**

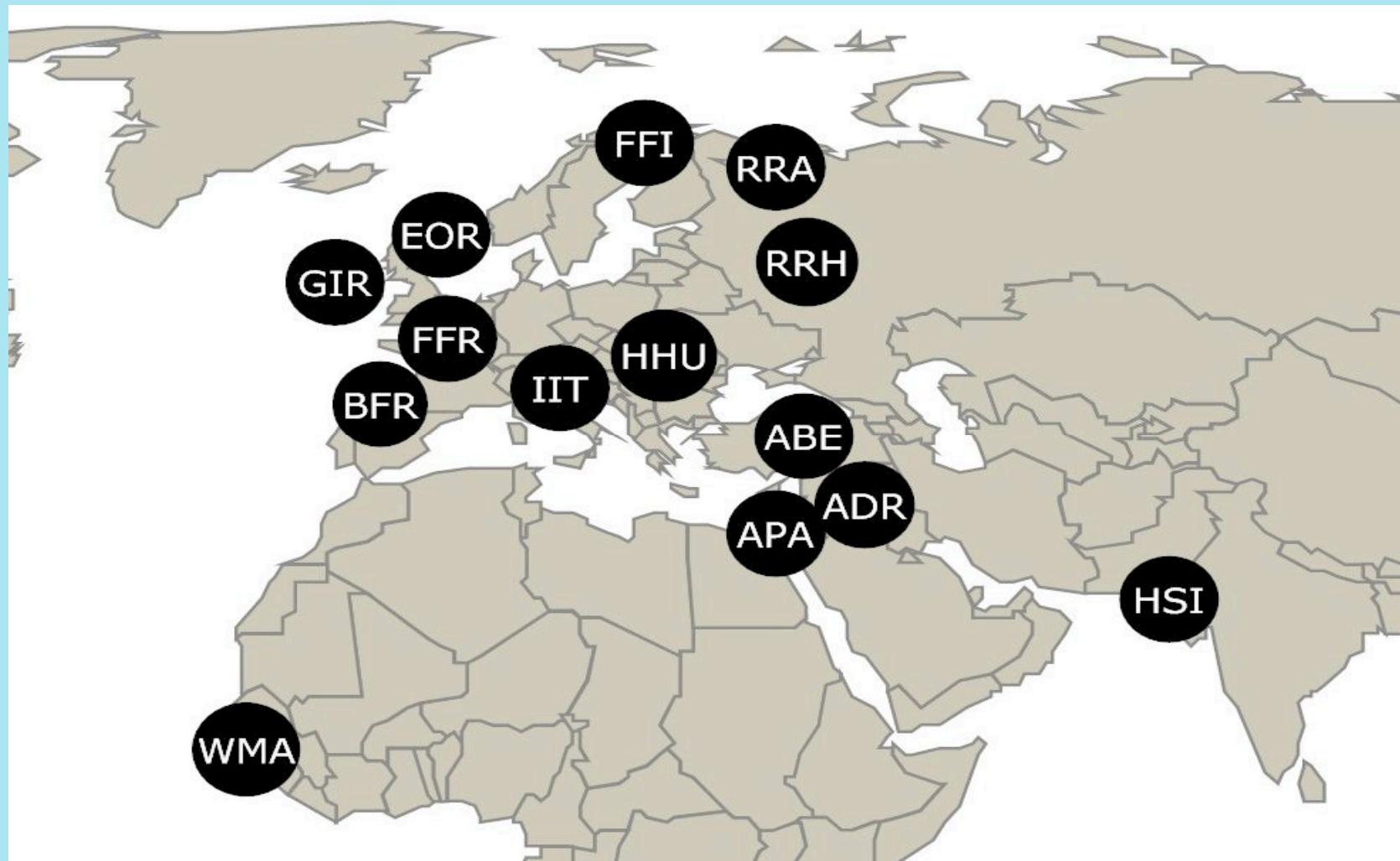
**Indirect argument for some version of an Inertial
Theory (Keenan 1994) of diachronic syntax**

Scatter plot: syntax/lexicon

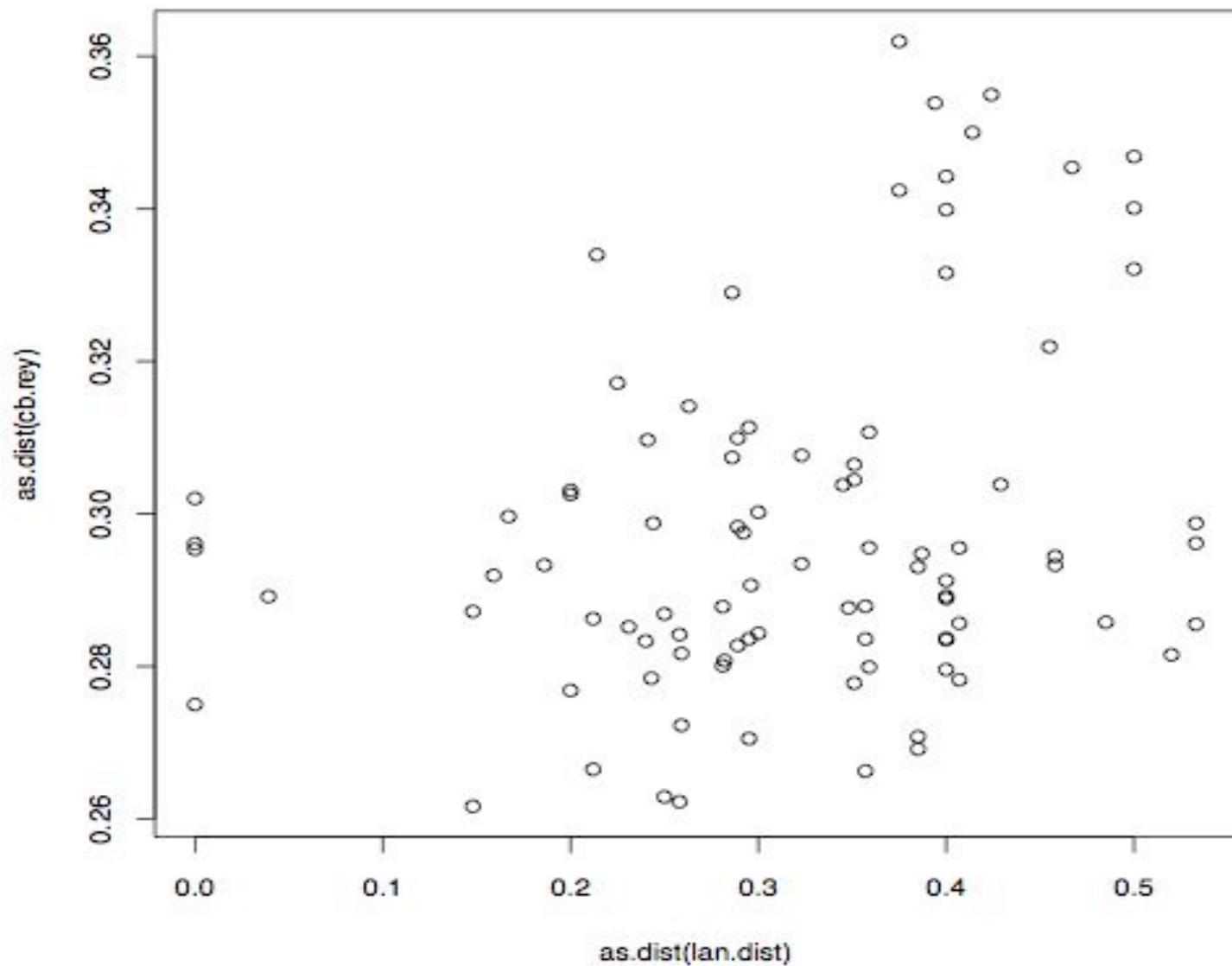


20

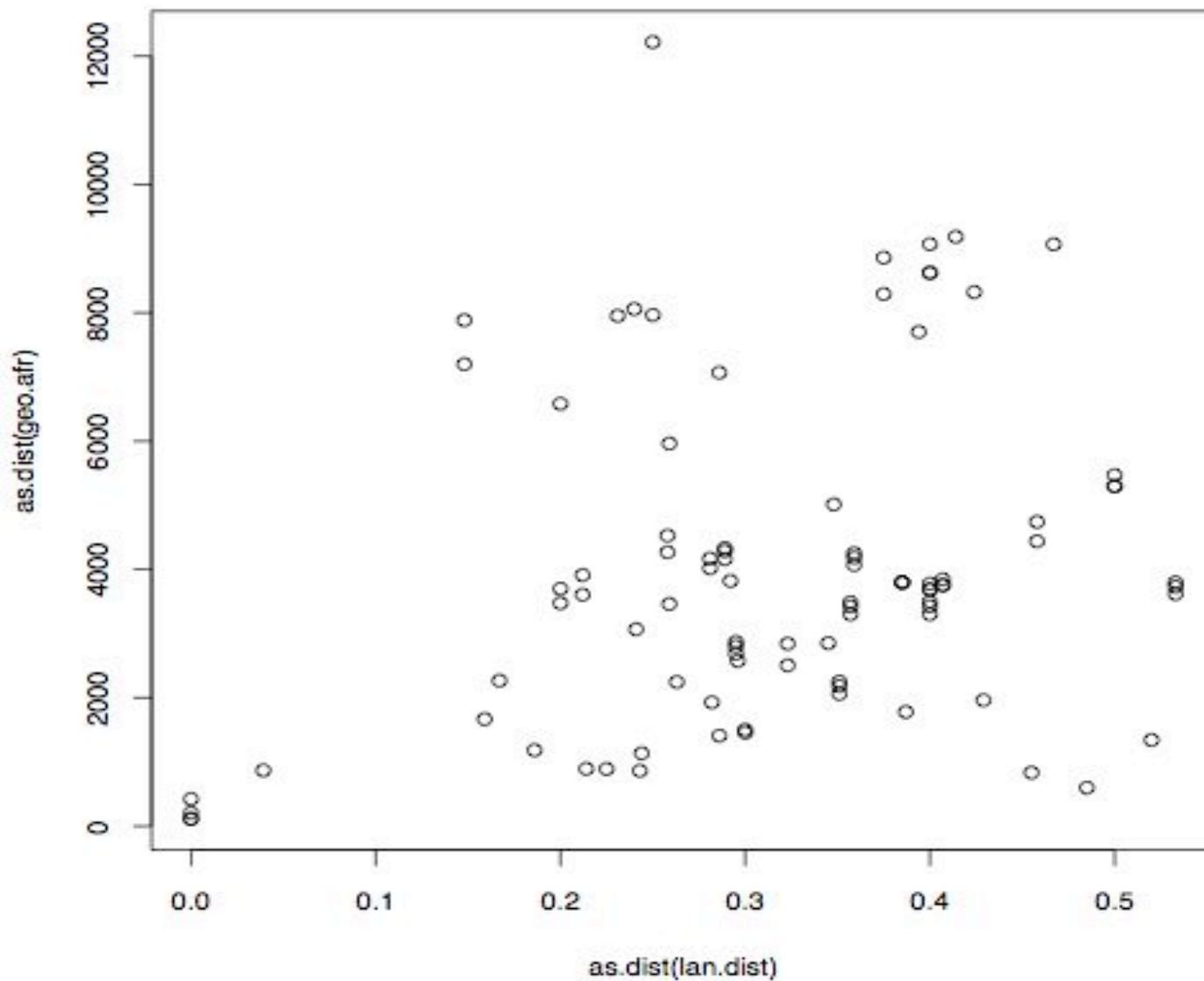
From Colonna, Boattini, Dall'Ara, Guardiano,
Pettener, Longobardi, Barbujani: in *Human
Heredity*, 2010



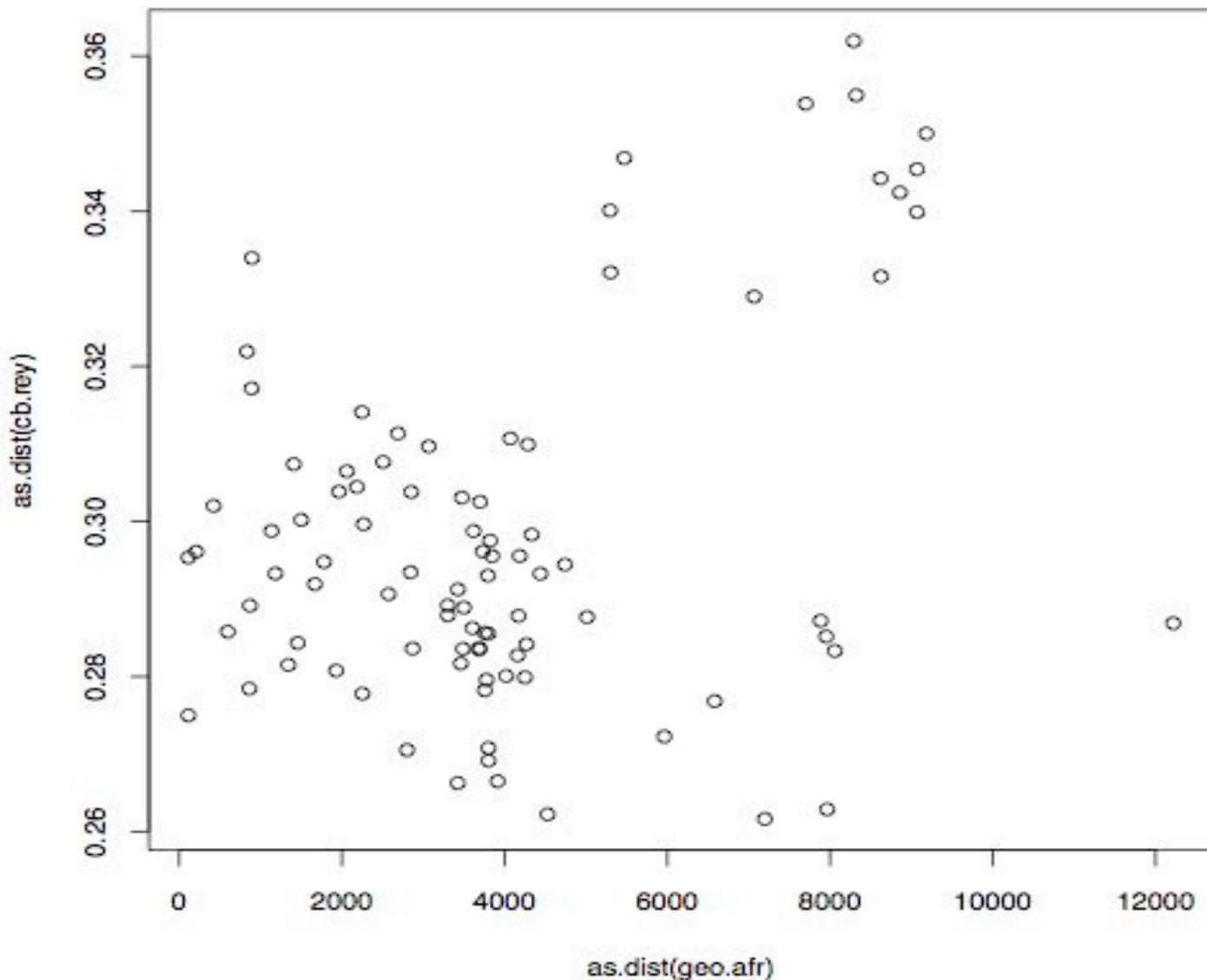
Languages and genes



Languages and geography



Genes and geography



Grammar as a population science: overcoming a XX century bias

In biology a sort of Galilean revolution (molecular analyses)

has affected both theoretical and evolutionary work leading to formal theories of biological diversity

In linguistics a Galilean revolution in the study of grammatical theory has been proposed by Chomsky but it has not yet been fully exploited for a mathematical approach to linguistic diversity