# At the boundaries of syntactic prehistory: metric and non-metric distances

Andrea Sgarro
sgarro@units.it

Dept. of Mathematics and Geosciences, University of Trieste (I)
Human Language Technologies Research Center, Bucharest (Ro)

RADH 2021
Bucharest, October 2021

Andrea Ceolin, Cristina Guardiano, Monica Alexandrina Irimia,
**Giuseppe Longobardi**, Luca Bortolussi, Andrea Sgarro

## *At the boundaries of syntactic prehistory*

Andrea Ceolin, Cristina Guardiano, Monica Alexandrina Irimia, **Giuseppe Longobardi**, Luca Bortolussi, Andrea Sgarro

**At the boundaries of syntactic prehistory**

Philosophical Transactions B, Royal Society (2021)

Laura Franzoi, Andrea Sgarro, Anca Dinu, Liviu P. Dinu

**Random Steinhaus distances for robust syntax-based classification of partially inconsistent linguistic data**

IPMU 2020, Lisbon (Pt)

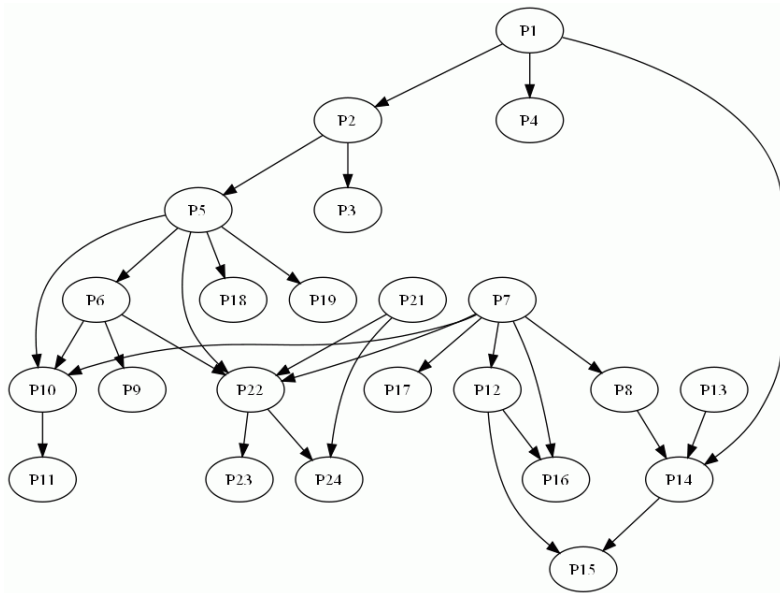Parametric Comparison Method PCM

94 syntactic parameters, 58 languages from the Old World

possible languages

94 parameters as before, 5000 possible languages

## results

controversial clusters such as Altaic (Japanese, Korean, Mongolian, ...)
or Uralo-Altaic were signifcantly supported, while other possible
macro-groupings as Indo-Uralic or Basque-Caucasian were not

## Longobardi distances, Hamming-like and Jaccard-like

L = 0 | 1 | * | 1 | 0 | 1
Λ = 0 | 0 | * | * | 1 | 1

$$\text{dist}_H(\Lambda, L) = \frac{\#\text{ bit differences}}{\text{"sound" bit length}} = \frac{2}{4} = \frac{1}{2}$$

## Longobardi distances, Hamming-like and Jaccard-like

L = 0 | 1 | * | 1 | 0 | 1
Λ = 0 | 0 | * | * | 1 | 1

$$\text{dist}_H(\Lambda, L) = \frac{\# \text{ bit differences}}{\text{"sound" bit length}} = \frac{2}{4} = \frac{1}{2}$$

$$\text{dist}_J(\Lambda, L) = \frac{\# \text{ bit differences}}{\text{sound length} - \# \text{ "irrelevant" positions}} = \frac{2}{4 - 1} = \frac{2}{3}$$
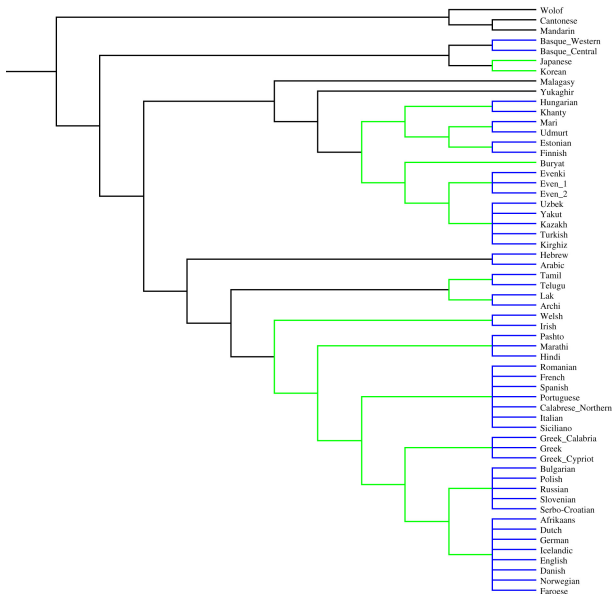
## Longobardi distances, Hamming-like and Jaccard-like

L = 0 | 1 | * | 1 | 0 | 1
Λ = 0 | 0 | * | * | 1 | 1

$$\text{dist}_H(\Lambda, L) \ = \ \frac{\# \text{ bit differences}}{\text{"sound" bit length}} \ = \ \frac{2}{4} = \frac{1}{2}$$

$$\text{dist}_J(\Lambda, L) \ = \ \frac{\# \text{ bit differences}}{\text{sound length} - \# \text{ "irrelevant" positions}} \ = \ \frac{2}{4 - 1} = \frac{2}{3}$$

both might violate the triangle inequality

# what should a distance be?

## at least...

- $d(x, y) \geq 0$
- $d(x, x) \leq \min\left[d(x, y), d(y, x)\right]$

## (ordered) triangle inequality

$$d(x, y) \leq d(x, z) + d(z, y)$$

Steinhaus transform or biotope transform of the distance $d$:

$$S_d(x, y) \doteq \frac{2d(x, y)}{d(x, y) + d(x, z) + d(y, z)}$$

where:

- $x, y, \ldots$ are objects (possibly strings)
- $d(x, y)$ is their distance
- $z$ is a fixed object called the pivot $z$

We'll have to generalize to several pivots
$S_d(x, y)$ preserves metricity

Steinhaus transform or biotope transform of the distance $d$:

$$S_d(x, y) \doteq \frac{2d(x, y)}{d(x, y) + d(x, z) + d(y, z)}$$

where:

- $x, y, \ldots$ are objects (possibly strings)
- $d(x, y)$ is their distance
- $z$ is a fixed object called the pivot $z$

We'll have to generalize to several pivots
$S_d(x, y)$ preserves metricity

From (normalized) Hamming to Jaccard:
the objects are $n$-lenght strings,
the pivot $z = \underline{z}$ is the all-0 string

$x, y$ strings of $n$ logical values

$$d(x, y) = \sum_i \left[x_i \text{ AND } \neg y_i\right] \text{ OR } \left[\neg x_i \text{ AND } y_i\right]$$

$x, y$ strings of $n$ logical values

$$d(x, y) = \sum_i \left[ x_i \ \text{AND} \ \neg y_i \right] \ \text{OR} \ \left[ \neg x_i \ \text{AND} \ y_i \right]$$

standard fuzzy logical operators, $\text{OR} = \max$, $\text{AND} = \min$

Solomon Marcus (1925-2016)

$x, y$ strings of $n$ logical values

$$d(x, y) = \sum_i \left[ x_i \ \text{AND} \ \neg y_i \right] \ \text{OR} \ \left[ \neg x_i \ \text{AND} \ y_i \right]$$

standard fuzzy logical operators, OR $=$ max , AND $=$ min

Solomon Marcus (1925-2016)

why do not start from the fuzzy Hamming distance?

$x, y$ strings of $n$ logical values

$$d(x, y) = \sum_i \left[ x_i \ \text{AND} \ \neg y_i \right] \ \text{OR} \ \left[ \neg x_i \ \text{AND} \ y_i \right]$$

standard fuzzy logical operators, OR $=$ max , AND $=$ min

Solomon Marcus (1925-2016)

why do not start from the fuzzy Hamming distance?

Łukasiewicz:   OR $=$ min $\left[ (x + y), 1 \right]$ , AND $=$ max $\left[ (x + y - 1), 0 \right]$

## taxicab or Minkowski or Łukasiewicz distance:

$$d(x, y) = \sum_i |x_i - y_i|$$

# taxicab or Minkowski or Łukasiewicz distance:

$d(x, y) = \sum_i |x_i - y_i|$

$$* \implies \frac{1}{2}$$

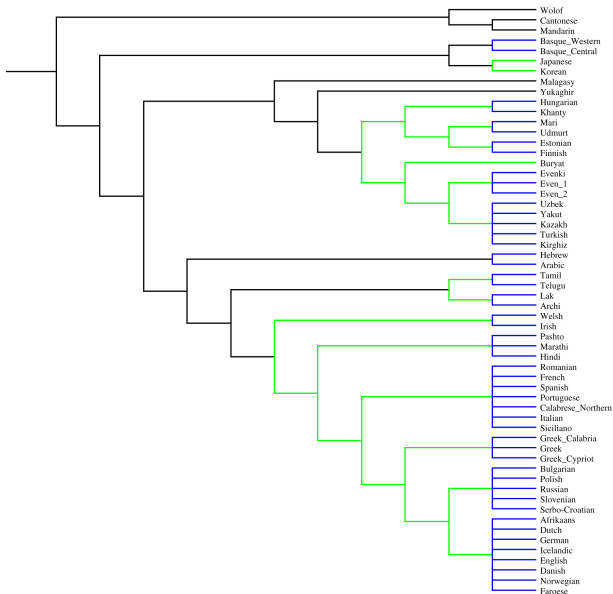$$d(bit, *) = d(*, bit) = \frac{1}{2}, \ d(*, *) = 0$$

pivot of the Steinhaus transform: the "totally unsound" all-∗ sequence

consistency $\chi(x)$ of the string $x$: its taxicab distance from the all-∗ string

$$S_d(x,y) \doteq \frac{2d(x,y)}{d(x,y) + \chi(x) + \chi(y)}$$

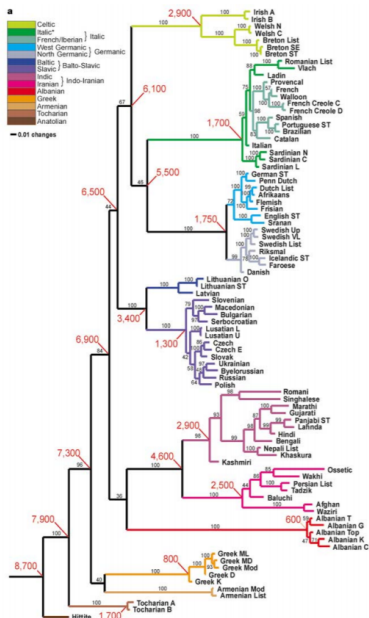pivot of the Steinhaus transform: the "totally unsound" all-$*$ sequence

consistency $\chi(x)$ of the string $x$: its taxicab distance from the all-$*$ string

$$S_d(x,y) \doteq \frac{2d(x,y)}{d(x,y) + \chi(x) + \chi(y)}$$

weight $w(x)$ of the string $x$: its taxicab distance from the all-0 string

$$S_d(x,y) \doteq \frac{2d(x,y)}{d(x,y) + \min\left[\chi(x) + \chi(y), w(x) + w(y)\right]}$$

thanks, mulțumesc, grazie