

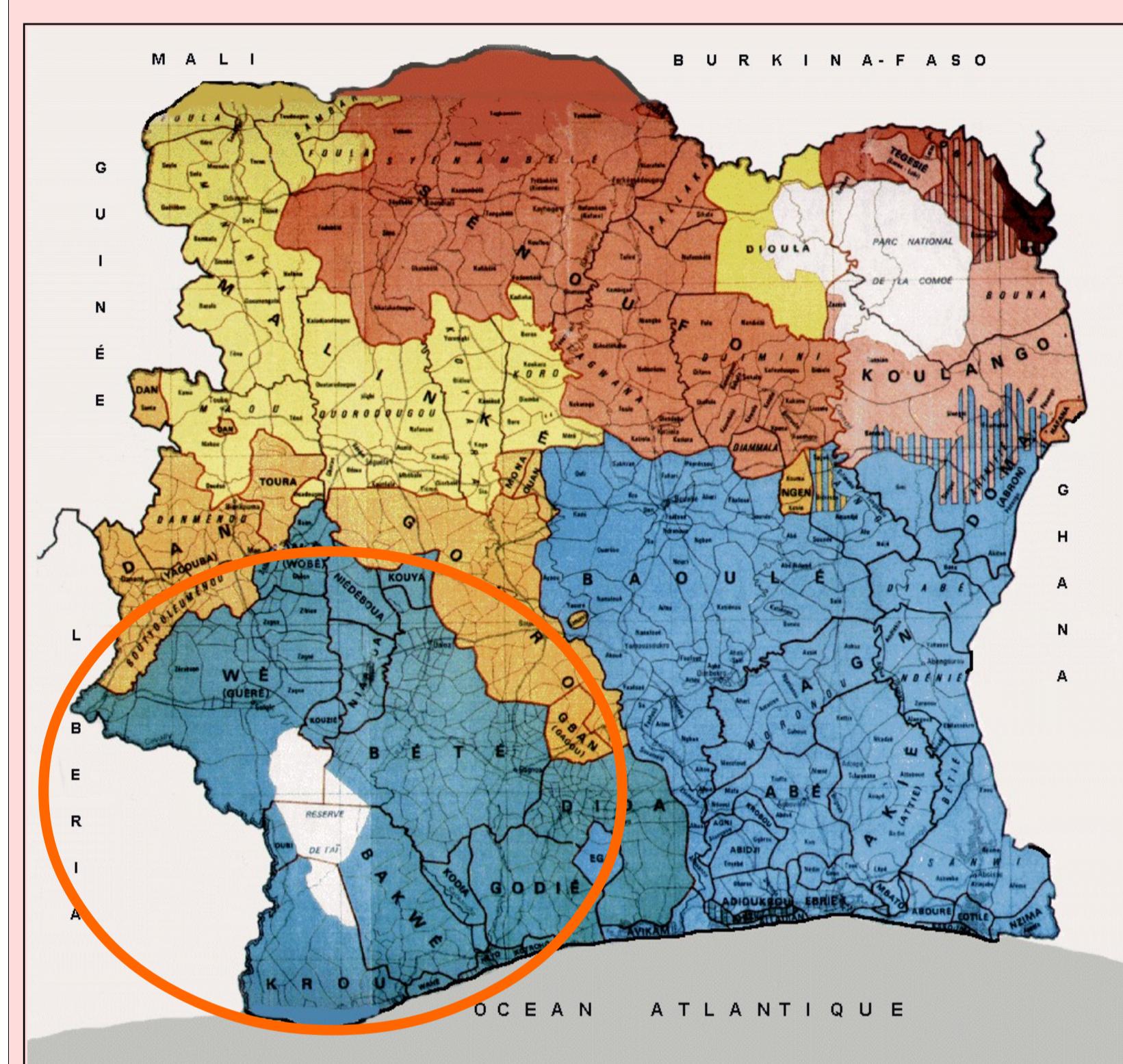
# Legacy Language Atlas Data Mining: Mapping Kru Languages

Dafydd Gibbon (Universität Bielefeld, Germany)

Resource type: data + online tool 'DistGraph'

- \* Classification of Kru languages in the context of a 'Language in Context' documentation project of Côte d'Ivoire languages: Stavros Skopeteas (Bielefeld), Firmin Ahoua (Abidjan), Dafydd Gibbon (Bielefeld), in cooperation with François Kipré Blé (Abidjan)
- \* Data revival: re-use of >30 year old Language Atlas:
  - . digital reproduction of data: scanning, retyping, redrawing
  - . crosscheck of historical/typological classifications from Language Atlas as basis for new atlas

## Data: the languages



South-West Ivory Coast

### Ethnologue:

Niger-Congo (1537)  
Atlantic-Congo (1440)  
Volta-Congo (1367)

### Kru (39)

- + Aizi (3)
- Eastern (11)
  - + Bakwe (2)
  - + Bete (5)
  - + Dida (3)
  - + Kwadia (1)
  - + Kuwaa (1)
- + Seme (1)
- Western (23)
  - + Bassa (3)
  - + Grebo (9)
  - + Kla (2)
  - + Wee (9)

The 39 Ethnologue entries have ISO 639-3 codes, but Cedepo, Dewoin, Koyo and Niaboua are not listed. In some cases more than one language variety is listed. Dida and some Dida varieties are listed, but Atlas varieties named Dida de Lozoua and Dida F are not.

Few Kru languages have ISO 639-3 codes

## Data: the language atlas

- \* Marchese, Lynell. 1984. *Atlas linguistique kru*. Agence de coopération culturelle et technique, Université d'Abidjan, 3ème éd.
- \* Contents: language sketch tables & maps for 19 languages
- \* Selection: consonant tables for 19 languages, 44 different consonants
- \* Why consonants and not lexical items?
  - . Lexical items are highly heterogeneous, easily borrowed
  - . Consonant systems are relatively stable, slow changing
  - . Consonant change laws are well-established for many language families (cf. Grimm's Law, Verner's Law, High German Sound Shift)

## Method ('BLARK' for language typology?)

1. Input:
  - 19 ordered consonant sets x 44 features (consonants)
2. Outputs:
  - pairwise difference matrix (Hamming distance)
  - feature ranking list (variance)
  - distance distribution histogram
  - table of average distance/isolation
  - table of specific pairwise differences

## Implementation

- \* Server-side web application:
  - . HTML → CGI → HTML+graphics
  - . Linux, Windows (public & localhost)
  - . Python 2.7
  - . GraphViz neato engine (line drawings)
  - . SciPy + Matplotlib (dendrogram)
- \* Client:
  - . (almost) any browser
  - . resource demo:
    - localhost tablet & laptop
    - internet (see address on footer)

This study is dedicated to the memory of our late colleague and Symposium host, Henrike Grohs, Director of Abidjan Goethe Institut, cruelly murdered by terrorists in Grand Bassam, Côte d'Ivoire  
13<sup>th</sup> March 2016.

## Data flow



Godié de Dakpado et Legiko (Marchese, 1975)	Koyo (Kakora, 1976, p. 23)
p t c k kp kw	p t c k kp cj
b d ž g gb gm	b d ž g gb
f s	f s
v z	v z
b l j y w	b i j y(2) w
m n n q gn	m n n q

### Parameter settings (+ CSV input field for consonant table)

<b>IO parameters</b>	<b>Output type:</b>
Input table CSV separator: semicolon ▾	<input checked="" type="radio"/> parametrised LED graph <small>(properties of same attributes in same field position)</small>
Graphics format: GIF bitmap graphics (smallest files) ▾	<input checked="" type="radio"/> parametrised SIRD graph <small>(use only if properties in different fields are different, i.e. sets)</small>
	<input checked="" type="radio"/> CSV ▾ HTML ▾ XML formattted input data
	<input checked="" type="radio"/> CSV ▾ HTML ▾ XML output of LED distance matrix
	<input checked="" type="radio"/> CSV ▾ HTML ▾ XML output of LED distance triples

### Graph parameters

- Graph engines (from AT&T GraphViz package):
- neato spring model
  - dot undirected graph model
  - twopi centred circle model
  - circo circle model

Numerical parameters:  
range of distances to be processed: 0 ... 6 (check distance matrix for full data range)  
random seed for neato spring model (trial and error): 6  
minimal scaling ▾ % graph width (percent of window): 90

<i>Language Atlas mining:</i><br>Consonant sets of Kru Languages.<br>Data source: Marchese, Lynell. 1984. <i>Atlas Linguistique des Langues Kru.</i><br>Lynell, Marchese, 1984. <br>atlas linguistique kru NOUVELLE EDITION

### Aligned Language x Feature (=consonant) table

Bete	p t c k kp kw	— b d C	g gb	f s	v z	— — — — —	B	l	j	x w m n J N Nw	— — — — —
Godie	p t c k kp kw	— b d C	g gb	gw f s	v z	— — — — —	B	l	j	x w m n J N Nw	— — — — —
Koyo	p t c k kp kw	kj b d C	g gb	f s	v z	— — — — —	B	l	j	x w m n J N Nw	— — — — —
Neyo	p t c k kp kw	b d C	g gb	f s	v z	— — — — —	B	l	j	x w m n J N Nw	— — — — —
DidaDeLozoua	p t c k kp kw	b d C	g gb	gw f s	v z	— — — — —	B	l	j	x w m n J N Nw	— — — — —
DidaF	p t c k kp kw	b d C	g gb	gw f s	v z	— — — — —	B	l	j	x w m n J N Nw	— — — — —
Wobe	p t c k kp kw	b d C	g gb	gw f s	v z	— — — — —	B	l	j	x w m n J N Nw	Nm
Guere	p t c k kp kw	b d C	g gb	gw f s	v z	— — — — —	B	l	j	x w m n J N Nw	Nm km
Krahn	p t c k kp kw	b d C	g gb	f s	— — — — —	B	l	j	x w m n J N Nw	Nm km	
Cedepo	p t c k kp kw	b d C	g gb	f s	— — — — —	B	l	j	x w m n J N Nw	Nm km	
Kla	p t c k kp kw	b d C	g gb	f s	— — — — —	B	l	j	x w m n J N Nw	Nm km	
Niaboua	p t c k kp kw	b d C	g gb	gw f s	v z	— — — — —	B	l	j	x w m n J N Nw	Nm km
Dewoin	p t c k kp kw	b d C	g gb	gw f s	v z	— — — — —	B	l	j	x w m n J N Nw	Nm km
Bassa	p t c k kp kw	b d C	dj g gb	f s	v z	h hw	B	l	j	x w m n J N Nw	Nm km
Grebo	p t c k kp kw	b d C	g gb	f s	v z	h hw	B	l	j	x w m n J N Nw	Nm km
Teppo	p t c k kp kw	b d C	g gb	f s	v z	h hw	B	l	j	x w m n J N Nw	Nm km
KuwaaLiberia	p t c k kp kw	b d C	g gb	f s	v z	— — — — —	B	l	j	x w m n J N Nw	Nm km
SemeHauteVolta	p t c k kp kw	b d C	g gb	f s S v	— — — — —	B	l	j	x w m n J N Nw	Nm km	
AiziCdl	p t c k kp kw	b d C	g gb	f s S v z	— — — — —	B	l	j	x w m n J N Nw	Nm km	

### Hamming distance measure

For binary sequences of equal length:  
Feature coding is {1,0}

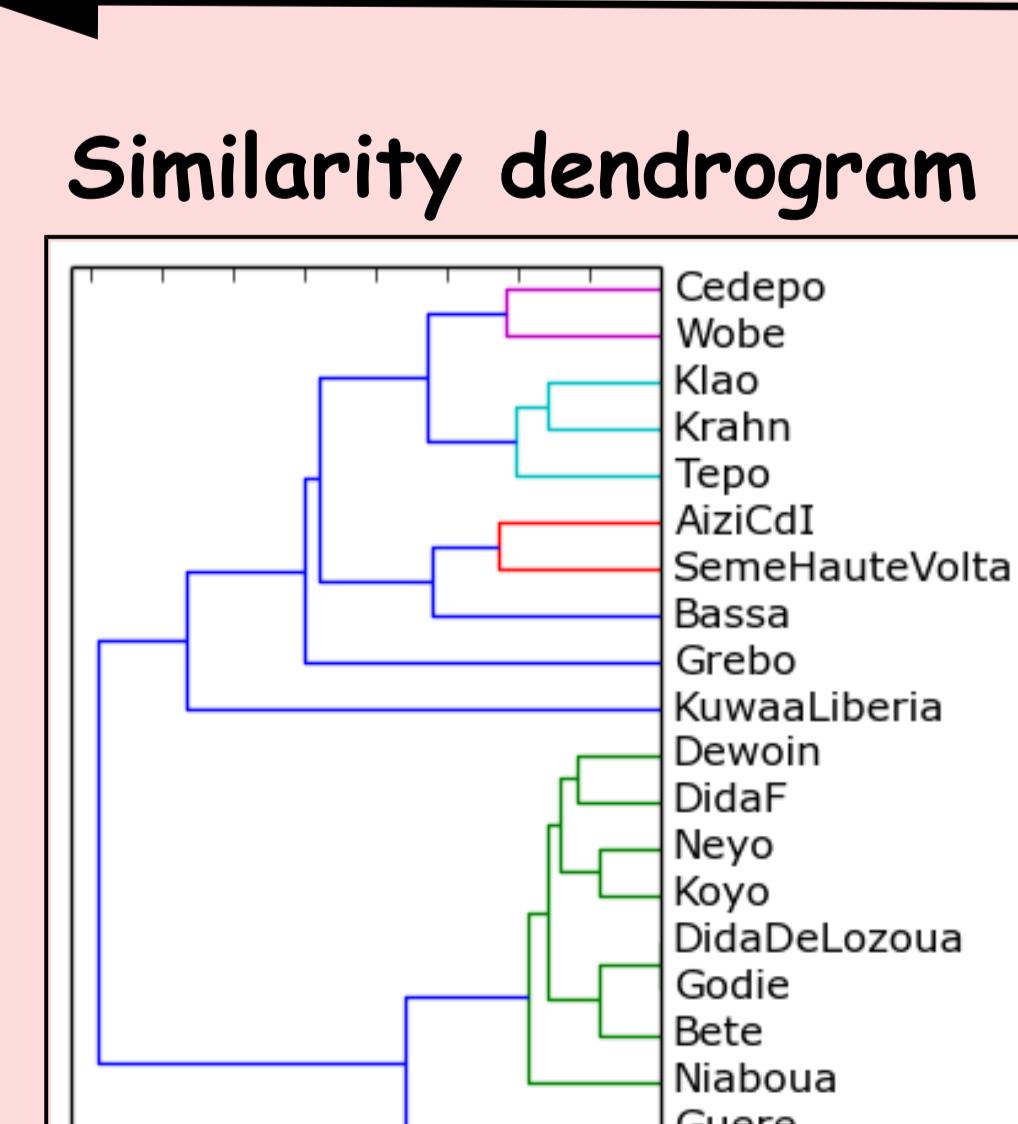
$$\sum_{i=1}^n |x_i - y_i|$$

Length normalisation not used:  
 $\frac{\sum_{i=1}^n |x_i - y_i|}{n}$

### Pairwise distance matrix (column headers = row headers)

Eastern	Bete	0 1 2 1 1 3 10 6 9 11 8 4 4 7 11 8 12 9 6
	Godie	1 0 3 2 0 2 11 5 10 12 9 3 3 8 12 9 13 10 7
	Koyo	2 3 0 1 3 12 8 9 11 8 4 4 9 13 8 12 9 13 10 6
	Neyo	1 2 1 0 2 1 7 8 10 7 3 3 8 12 9 13 10 7
	DidaDeLozoua	1 0 3 2 0 2 11 5 10 12 9 3 3 8 12 9 13 10 7
Western	DidaF	3 2 3 2 2 0 11 5 10 10 7 3 3 10 12 7 13 10 7
	Wobe	10 11 12 11 11 11 0 8 6 6 4 10 12 12 11 8 14 11 12
	Guere	6 5 8 7 5 8 0 11 11 8 4 6 9 13 10 18 11 10
	Krahn	9 10 9 8 10 10 6 11 0 4 3 7 9 10 12 5 11 8 9
	Cedepo	11 12 11 10 12 10 6 11 4 0 3 9 11 10 10 5 13 8 11
Isolates	Kla	8 9 8 7 9 7 4 8 3 3 0 6 8 11 9 4 10 7 8
	Niaboua	4 3 4 3 3 3 10 4 7 9 6 0 2 7 13 8 14 7 6
	Dewoin	4 3 4 3 3 3 12 6 9 11 8 2 0 9 13 8 12 9 6
	Bassa	7 8 9 8 8 10 12 9 10 11 7 9 0 10 11 19 8 9
	Grebo	11 12 13 12 12 11 13 12 10 9 13 13 10 0 7 17 10 11
	Teppo	8 9 8 7 9 7 8 10 5 4 8 8 11 7 0 12 7 8
	KuwaaLiberia	12 13 12 11 13 13 14 18 11 13 10 14 12 19 17 12 0 15 14
	SemeHauteVolta	9 10 9 8 10 10 11 8 8 7 7 9 8 10 7 15 0 5
	AiziCdl	6 7 6 5 7 7 12 10 9 11 8 6 6 9 11 8 14 5 0

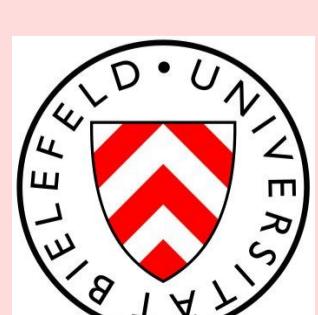
### Similarity dendrogram



### Distance (Difference) Map (force/spring map)

DIMENSION REDUCTION  
CLASSIFICATION  
VISUALISATION

### Typological Similarity Dendrogram (hierarchical clustering)



gibbon@uni-bielefeld.de

LREC 2016, Portorož, Slovenia

<http://wwwhomes.uni-bielefeld.de/gibbon/DistGraph/>

