

Advancing cognate reflex prediction

Incorporating expert evaluations and multi-tiered representations

Tiago **Tresoldi**¹, Freja **Lindgren**¹, Oscar **Billing**¹, John **Huisman**¹ & Fabrício Ferraz **Gerardi**²

1 Uppsala Universitet

2 Eberhard Karls Universität Tübingen

56th SLE meeting, Athens, 29.08.2023



UPPSALA
UNIVERSITET

Background

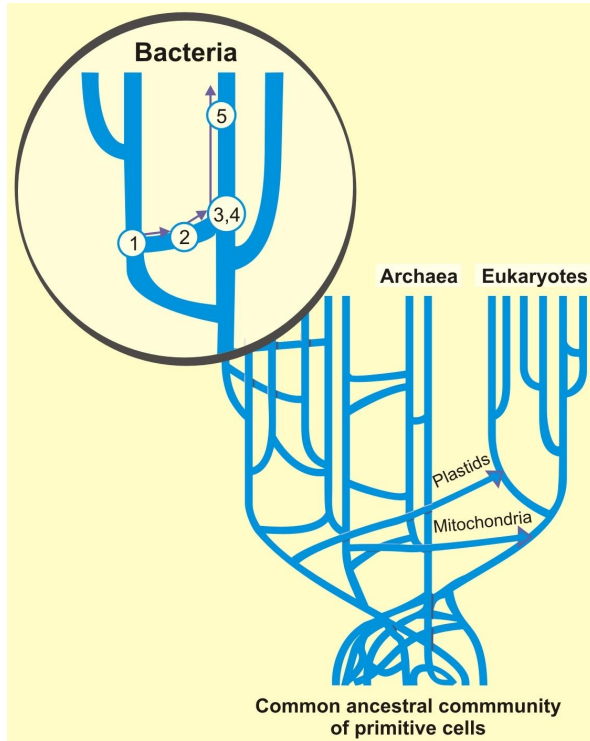
- Different but similar approaches
 - “computational historical linguistics” (Jäger, 2019)
 - “computational phylogenetics” (Bown, 2018)
 - “phylolinguistics” (Greenhill et al., 2020)
 - “computer-assisted language comparison” (List, 2017)
 - Moscow School of Comparative Linguistics
- Automated steps yielding faster, reproducible, unbiased(?) analyses; however:
 - Focus on tree-building from lexical replacement
 - Minimal commitment to the automation, enhancement, and support of routine tasks in historical linguistics (like inference of sound changes)



DALL-E 2 for “futuristic Tower of Babel”

The Challenge

- Cognate reflex prediction
 - Identifying semantic shifts
 - Cognacy probability and borrowings
 - Imputation of missing reflexes
 - Assistance with reconstructions
 - Building step towards sound change inference
- Different approach
 - To explain and not to predict
 - “interpretable” machine learning



“Horizontal gene transfer” (Wikipedia Commons)

Previous work

Cognate Set	German	English	Dutch
ASH	a ʃ ε	æʃ	ɑ s
BITE	b ai s ə n	b ai t	b ei t ə
BELLY	b au x	-	b œi k

- Notable earlier attempts:
 - character-based (Beinborn et al., 2013)
 - encoder-decoder (Meloni et al., 2021)
 - low-resource machine-translation (Fourrier et al. (2021)
 - RNN (Dekker and Zuidema, 2021)
 - SYGTYP2022 shared task
- Two contributions from Hammarström et al. (2019)
 - Focus on predicting rather than explaining (Shmueli, 2010)
 - “the most conspicuous **difference** between Computational Historical Linguistics [...] and the classical comparative method is the recognition of **sound changes and their role** in reconstruction and subgrouping. Sound changes are especially powerful for subgrouping as they are typically directional”

Our Approach

- Roots in Chacon & List (2015)
- Novel method combining
 - extended phonological alignments (“multitiers”)
 - expert evaluations
 - focus on explainability
 - ultimate focus on sound changes
- Linguistics groups under study
 - Dravidian
 - Anatolian
 - Japonic
 - Tupi-Guarani






“SumperJumbo black box” (Getty Images)

Handling of phonological data

- Properly modelling sounds/words and sounds changes
 - “properly” means **not** treating them as orthographic characters and strings
 - Understand that segments are abstractions and suprasegmentals are essential
- **maniphono** allows the “symbolic manipulation” of sounds, segments, and sequences, and includes tools for normalizing transcriptions
 - Both human- and machine- representation
 - Building on experience from **CLTS** (Anderson et al., 2018)
- **alteruphono** allows to apply sound changes in forward and backwards direction
 - Necessary for any non-trivial process, such as apophonies and metaphonies
 - Building on experience from **foma** (Hulden, 2009)

Alteruphono

802 lines (802 sloc) 40.1 KB					Raw	Blame			
Search this file...									
1	ID	RULE	WEIGHT	TEST_ANTE	TEST_POST				
2	0	V s > @1 z @1 / # p b r _ t d	1	presto	prezeto				
3	1	C N > @1 / _ #	1	agrogŋ	agrog				
4	2	C > :null: / _ #	1	adja:d	adja:				
5	3	C > :null: / r _	1	ik?erja	ik?era				
6	4	s C > :null: / _ #	1	gegsisk	gegsi				
7	5	s k C > @1 / _ #	1	akankmiks	akankmik				
8	6	K > ? / _ s	2	akankmiks	akankmi?s				
9	7	L > d / # _	1	labjoplɔl	dabjoplɔl				
10	8	N S > h @2	3	rimbsu	rihbsu				
11	9	N > m / _ #	2	akθun	akθum				
12	10	N > n / _ #	5	apmbirom	apmbiron				
13	11	N V > n @2	1	akankmiks	akankniks				
14	12	N d l > n @2	1	imlimko	inlimko				
15	13	N v k s > n @2	1	imlimko	imlinko				

Alteruphono

802 lines (802 sloc) 40.1 KB

Raw Blame

Search the

1 ID

2 0

3 1

4 2

5 3

6 4

7 5

8 6

9 7

10 8

11 9

12 10

13 11


14 12 N d|l > n @2 1 imlimko inlimko

15 13 N v|k|s > n @2 1 imlimko imlinko

/parba/

p -> b / _ v

/barba/



Alteruphono

802 lines (802 sloc) 40.1 KB

Raw Blame

Search the

1 ID

2 0

3 1

4 2

5 3

6 4

7 5

8 6

9 7

10 8

11 9

12 10

13 11

14 12 N d|l > n @2 1 imlimko inlimko

15 13 N v|k|s > n @2 1 imlimko imlinko

/barba/, /barpV/, /pVrba/, /pVrpV/

p -> b / _ v

/barba/

Multitiered representations

- The “sheet music” of alignments (List)
- Not aligned linear components, but dataset records
- The dataset becomes a 2D-matrix for machine learning
- Designed for improved explainability
 - Identify which tiers (“features”) contribute positively or negatively towards an outcome, how much they contribute, and their interaction effect

Aufführungsrecht vorbehalten.

Der 100. Psalm.

Max Reger, Op.106.

Maestoso (J.-72).
(Animato.)

2 große Flöten.
2 Oboen.
2 Klarinetten in A.
2 Fagotte.
2 Trompeten in C.
4 Hörner in F.
2 Tenorposaunen.
Baßposaune.
Baßtuba.
3 Pauken in A C D.
Große Trommel.
Becken.
Sopran.
Alt.
Tenor.
Baß.
Violinen I.
Violinen II.
Bratschen.
Violoncelli.
Kontrabässe.
Orgel.
Pedale.

Jauch - zet, jauch - zet, jauch - zet, jauch - zet.

d. Zu diesem Hauptorchester tritt im letzten Satz ein Nebenorchester von wenigstens 4 Trompeten (in C) und 4 Tenorposaunen zur Durchführung des Choral.

Edition Peters.

10438 Copyright 1909 by C. F. Peters, Leipzig.

Max Reger, Der 100 Psalm
(Wikimedia Commons)

Language	Site #1	Site #2	Site #3	Site #4	Site #5	Site #6
English	f	ɑ:	ð	-	ə	ɹ
German	f	ɑ:	t ^h	-	ɐ	-
Italian	p	a	d	r	e	-
Spanish	p	a	ð	r	e	-

Tier	Site #1	Site #2	Site #3	Site #4	Site #5	Site #6
Segment_EN	f	ɑː	ð	-	ə	ɹ
Segment_DE	f	ɑː	tʰ	-	ɐ	-
Segment_IT	p	a	d	r	e	-
Segment_ES	p	a	ð	r	e	-

Tier	Site #1	Site #2	Site #3	Site #4	Site #5	Site #6
Segment_EN	f	ɑː	ð	-	ə	ɹ
Segment_DE	f	ɑː	tʰ	-	ɐ	-
Segment_IT	p	a	d	r	e	-
Segment_ES	p	a	ð	r	e	-

Tier	Site #1	Site #2	Site #3	Site #4	Site #5	Site #6
Segment_EN	f	ɑː	ð	-	ə	ɹ
SC_EN	B	A	D	-	A	R
Segment_DE	f	ɑː	tʰ	-	ɐ	-
SC_DE	B	A	T	-	E	-
Segment_IT	p	a	d	r	e	-
SC_IT	B	A	D	R	E	-
Segment_ES	p	a	ð	r	e	-
SC_ES	B	A	D	R	E	-

Tier	Site #1	Site #2	Site #3	Site #4	Site #5	Site #6
Segment_EN	f	ɑː	ð	-	ə	ɹ
Segment_EN_L1	-	f	ɑː	ð	-	ə
SC_EN	B	A	D	-	A	R
SC_EN_R1	A	D	-	A	R	-
Segment_DE	f	ɑː	tʰ	-	ɐ	-
Segment_DE_L1	-	f	ɑː	tʰ	-	ɐ
SC_DE	B	A	T	-	E	-
SC_DE_R1	-	B	A	T	-	E
Segment_IT	p	a	d	r	e	-
Segment_IT_L1	-	p	a	d	r	e
...

Methodology

- “Leave-one out”: each reflex is dropped from the dataset before training
 - One training round for each reflex: lots of computation!
- Different models are trained to obtain the word prediction
 - Predictions have probabilities
- Dual evaluation approach
 - Normal statistical evaluation of machine learning
 - Expert evaluations, later incorporated to improve explanations
- Building blocks for sound change inference
- ***Still on-going!***



ID	Dutch	English	German
*bainan	be:n	bəʊn	bain
*balgiz	balx	bɛlɪ	balk
*fadeer	va:dər	fɑ:ðer	fa:tər
*fasteenan	vastə	fɑ:st	fastən
*fuloon		fəʊl	fɔ:lən
*kalbaz	kalf	kɑ:f	kalp
...

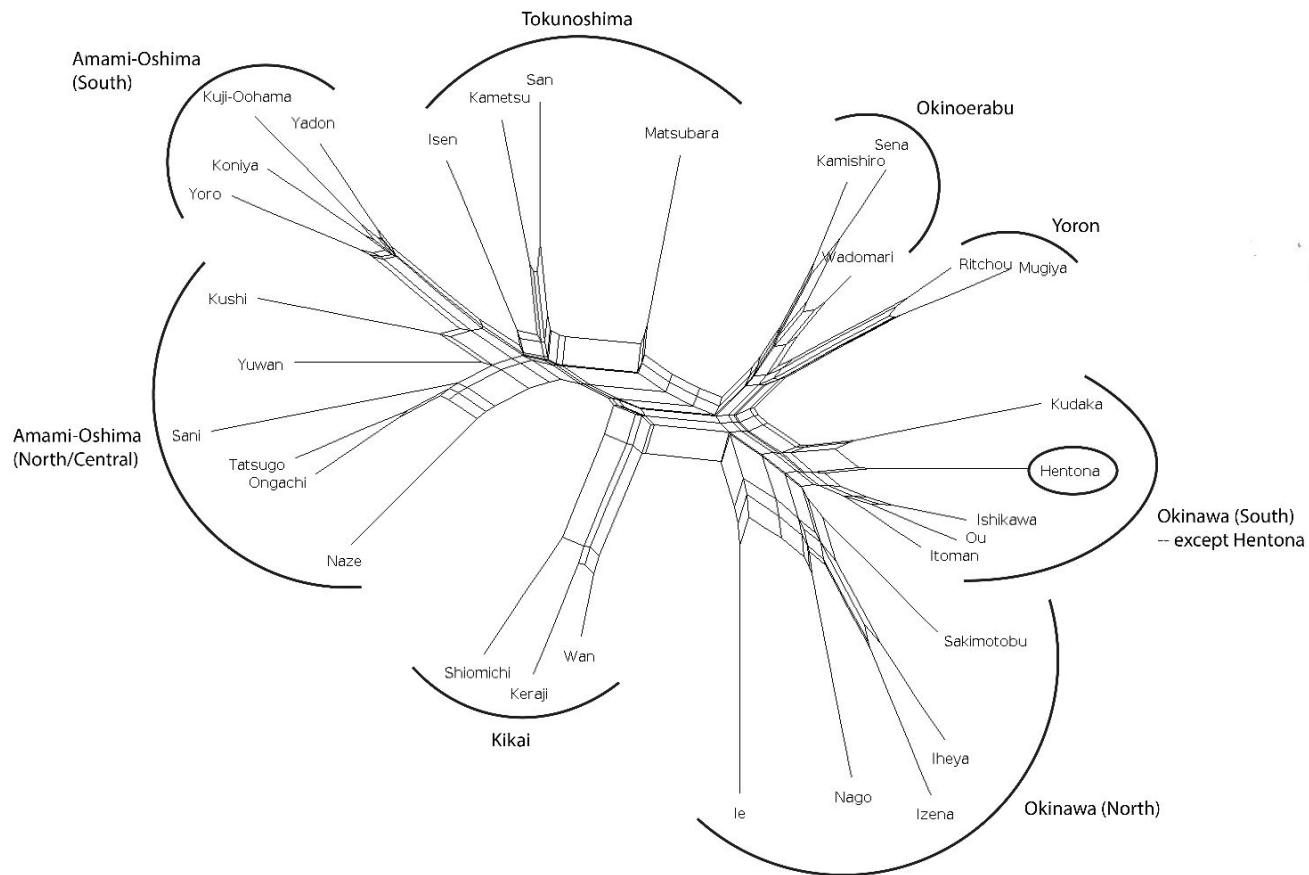
ID	Dutch	English	German
*bainan	be:n	bəʊn	bain
*balgiz	balx	bɛlɪ	balk
*fadeer	va:dər	fɑ:ðer	
*fasteenan	vastə	fɑ:st	fastən
*fuloon		fəʊl	fɔ:lən
*kalbaz	kalf	kɑ:f	kalp
...

	Site #1	Site #2	Site #3	Site #4	Site #5
Reference	f	a:	t	ə	r
Best	f (0.96)	a: (0.83)	t (0.88)	ə (0.92)	r (0.97)
Second best	v (0.03)	a (0.08)	d (0.10)	e (0.02)	- (0.02)
...



*fadeer	va:dər	fa:ðer	
*fasteenan	vastə	fa:st	fastən
*fuloon		fəʊl	fo:lən
*kalbaz	kalf	kɑ:f	kalp
...





Doculect	CONCEPT	IPA Alignment	Correspondences
Chabana	NOSE	p a n a	92 414 215 1881
Hentona	NOSE	ϕ a n a	92 414 215 1881
Izena	NOSE	ϕ a n a:	92 414 215 1881
Kametsu	NOSE	h a n a	92 414 215 1881
Chabana	WIND	h a d i	2971 414 298 299
Hentona	WIND	h a dʒ i	2971 414 298 299
Izena	WIND	h a dʒ i:	2971 414 298 299
Kametsu	WIND	k a d i	2971 414 298 299

	Chabana	p			
	Hateruma	p			
Docu	Iheya	h		IPA Alignment	Correspondences
Chaba	Ikema	h		p a n a	92 414 215 1881
Hento	Kudaka	p		ɸ a n a	92 414 215 1881
Izena	Kuroshima	p		ɸ a n a.	92 414 215 1881
Kame	Matsubara	h		h a n a	92 414 215 1881
Chaba		h a d i	2971 414 298 299
Hento		h a d ʒ i	2971 414 298 299
Izena	WIND			h a d ʒ i:	2971 414 298 299
Kametsu	WIND			k a d i	2971 414 298 299

Doculect		CONCEPT	IPA Alignment	Correspondences
Chabana		NOSE	p a n a	92 414 215 1881
Chabana	a		ɸ a n a	92 414 215 1881
Hateruma	a		ɸ a n a:	92 414 215 1881
Iheya	a		h a n a	92 414 215 1881
Ikema	a		h a d i	2971 414 298 299
Kudaka	a		h a dʒ i	2971 414 298 299
Kuroshima	a		h a dʒ i:	2971 414 298 299
Matsubara	a		k a d i	2971 414 298 299
...	...			

Docu	Chabana	d	IPA Alignment	Correspondences
Chaba	Hateruma	tʃ	p a n a	92 414 215 1881
Hentor	Iheya	dʒ	ɸ a n a	92 414 215 1881
Izena	Ikema	d	ɸ a n a:	92 414 215 1881
Kametsu	Kudaka	r	h a n a	92 414 215 1881
Chaba	Kuroshima	dʒ	h a d i	2971 414 298 299
Hentor	Matsubara	dʒ	h a dʒ i	2971 414 298 299
Izena	h a dʒ i:	2971 414 298 299
Kametsu	WIND		k a d i	2971 414 298 299

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

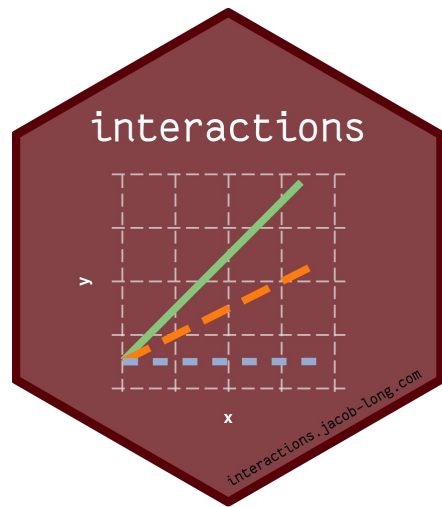
JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.



“Machine learning” (XKCD 1838)

Approach

- Prediction first (computers), explanation second (experts)
- Advanced statistical methods
 - Probabilities
 - Interaction effects
- A novel approach to cognate detection
- Ultimate result of reflex prediction: probabilities for segments and words
- Ultimate result of sound change inference: probabilities for sequences of sound changes, given in standard notation



Conclusions

- To explain and not to predict
- Not a black-box approach
- Agnostic about the machine learning solutions
- Encourage collaboration and a computer-assisted approach
- Paving the way for the Graal of Computational Historical Linguistics: inference of sound change history (relative chronology)

