

Learning pressures simplify morphology: Corpus, computational and experimental evidence

Christian Bentz

University of Tübingen

Aleksandrs Berdicevskis

The Arctic University of Norway

December 8, 2016



WORDS BONES GENES TOOLS
Tracking Linguistic, Cultural, and Biological Trajectories of the Human Past

EVOLAEMP
LANGUAGE EVOLUTION: THE EMPIRICAL TURN

OVERVIEW

INTRODUCTION

CORPUS ANALYSES

- Quantitative Measures

- Corpora

- Entropy Estimation

LEMMATIZATION

- Method

- Corpora

- Results

ITERATED LEARNING

- Experiments

- Results

MORPHOLOGICAL DIVERSITY

- (1) Hawaiian (Austronesian)

A ua olelo aku o Ioane ia ia [...]

Then PERF say.to SUBJ.Johan he.DAT [...]

"Then Johan said to him [...]"

- (2) Iñupiatun (Eskimo-Aleut)

Aglaan Jesus-ngum **itnagnigai** [...]

But Jesus-ERG this.say.report.3S.to.3PL

"But Jesus said to them [...]"

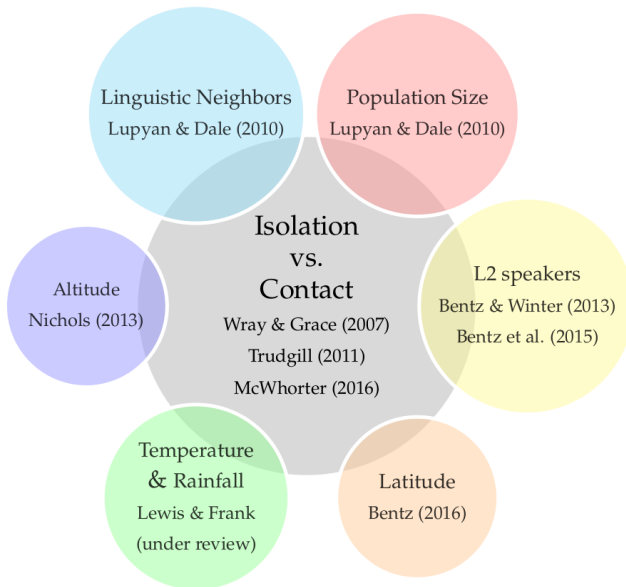


TYPOLICAL QUESTIONS

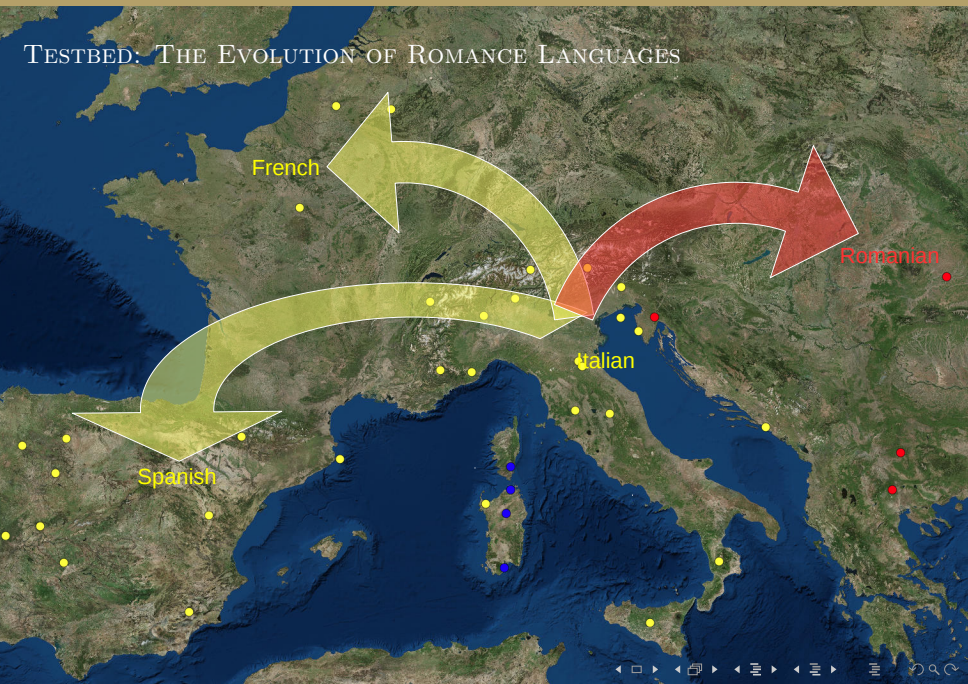
- ▶ How do we measure these differences in complexity?
- ▶ Are there systematic explanations for the patterns we find across the world?




“EXTERNAL” FACTORS



TESTBED: THE EVOLUTION OF ROMANCE LANGUAGES



TESTBED: THE EVOLUTION OF ROMANCE LANGUAGES

- 
- The background is a map of Europe. Large, semi-transparent arrows indicate the spread of languages. A yellow arrow points from the center of Europe (France/Italy area) towards the northwest (Spain/France) and southwest (Iberian Peninsula). A red arrow points from the center of Europe towards the northeast (Russia/Balkans area). Yellow dots are scattered across Western and Central Europe, while red dots are in the East. Labels for 'French', 'Spanish', 'Italian', and 'Romanian' are placed near their respective regions.
- ▶ As **Vulgar Latin** varieties spread throughout Europe, their **morphological complexity** was reduced (Herman & Wright 2000)
 - ▶ This is argued to be (partly) due to **L2 contact** (Herman & Wright 2000, Bentz & Christiansen 2013)

SIMPLE EXAMPLE: WORD FOR “BROTHER” IN THE BIBLE

► Latin

01004008 Dixitque Cain ad Abel **fratrem** suum [...]

01004009 Ubi est Abel **frater** tuus?

01004011 [...] suscepit sanguinem **fratris** tui de manu tua!

SIMPLE EXAMPLE: WORD FOR “BROTHER” IN THE BIBLE

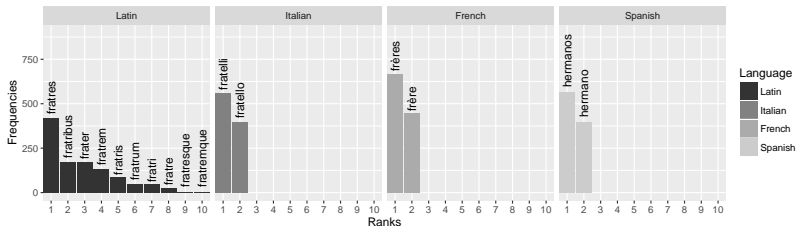
► Latin

01004008 Dixitque Cain ad Abel **fratrem** suum [...]
01004009 Ubi est Abel **frater** tuus?
01004011 [...] suscepit sanguinem **fratris** tui de manu tua!

► Italian

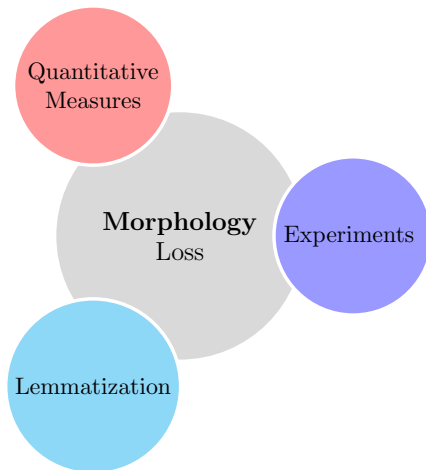
01004008 Caino disse al **fratello** Abele [...]
01004009 Dov'è Abele , tuo **fratello**?
01004011 [...] ha bevuto il sangue di tuo **fratello**!

SIMPLE EXAMPLE: WORD FOR “BROTHER” IN THE BIBLE



CONVERGING EVIDENCE

Bentz et al. (2015)
Bentz & Alikaniotis
(2016)
Ferrer-i-Cancho &
Bentz (forthcoming)



Bentz, Alikaniotis,
Samardžić & Buttery
(forthcoming)

CONVERGING EVIDENCE

Bentz et al. (2015)
Bentz & Alikaniotis
(2016)
Ferrer-i-Cancho &
Bentz (forthcoming)

Quantitative
Measures

Morphology
Loss

Experiments

Lemmatization

Bentz, Alikaniotis,
Samardžić & Buttery
(forthcoming)

Kirby et al. (2008,
2015)

Berdicevskis (2012)

Berdicevskis &
Semenuks
(forthcoming)



QUANTITATIVE MEASURES

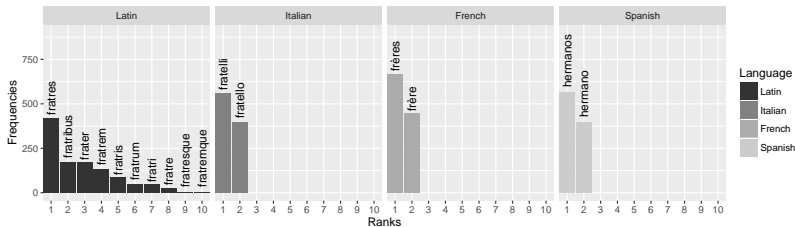
- ▶ sample of more than **500 languages** of 101 language families
- ▶ **4 corpus-based measures** compared to a measure based on **typological data (WALS)**
- ▶ strong Spearman correlations (up to 0.9) between all of them

Bentz, Ruzsics, Koplenig & Samardžić (2016)

SHANNON ENTROPY

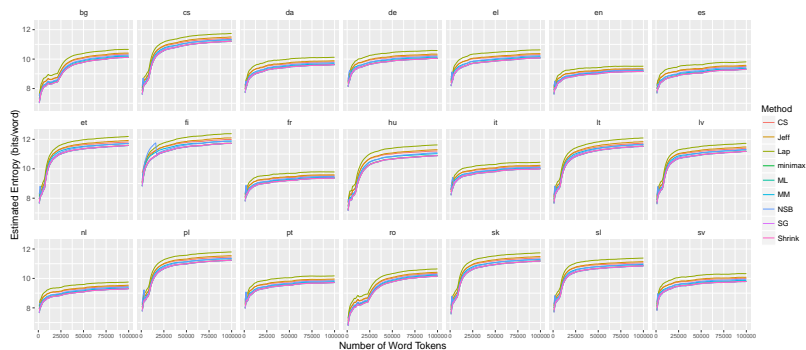
- Measure the skewness of the word form distribution via the **entropy H** according to Shannon (1949):

$$H(T) = - \sum_{i=1}^V p(w_i) \log_2(p(w_i)). \quad (1)$$



CONCEPTUAL PROBLEM

- ▶ word entropy depends on text size
- ▶ $H(T)$ converges onto stable value at ca. 50K tokens



Bentz et al. (forthcoming)

CORPUS ANALYSES

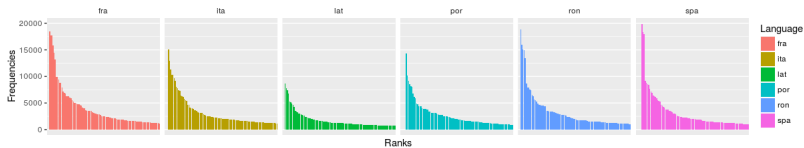
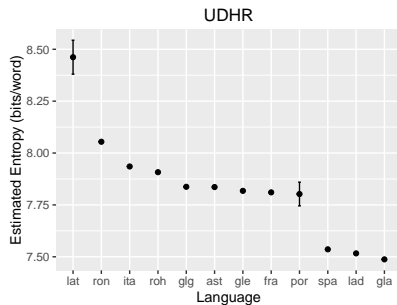
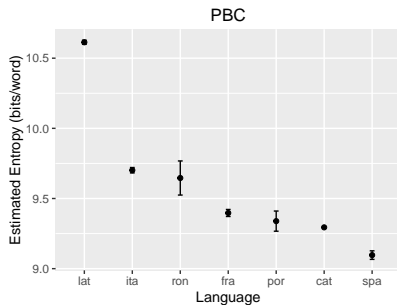
First Analysis

Measure the **entropy (bits/word) change** for **Latin towards Romance languages** using the so-called *James-Stein shrinkage estimator* (Hausser & Strimmer, 2014)

CORPUS ANALYSES

Select all the **Romance languages** (+ Latin) from the *Parallel Bible Corpus* (7 languages), and the *Universal Declaration of Human Rights* (10 languages).

RESULTS



CONCLUSION (CORPUS ANALYSES)

- ▶ Entropy is reduced from **ca. 11 (bits/word) in Latin** to **ca. 9.75-9.0 (bits/word) in 6 Romance languages** of the PBC, i.e. by around **10-15%**
- ▶ A similar pattern is found for the UDHR though with overall lower entropy values (due to differences in text size)

LEMMATIZATION

Second Analysis

Neutralize **inflectional marking** in Latin and the Romance languages to measure the effect of **inflectional differences** (Bentz, Alikaniotis, Samardžić & Buttery, in print)

TREETAGGER (SCHMIDT 1994, 1995)

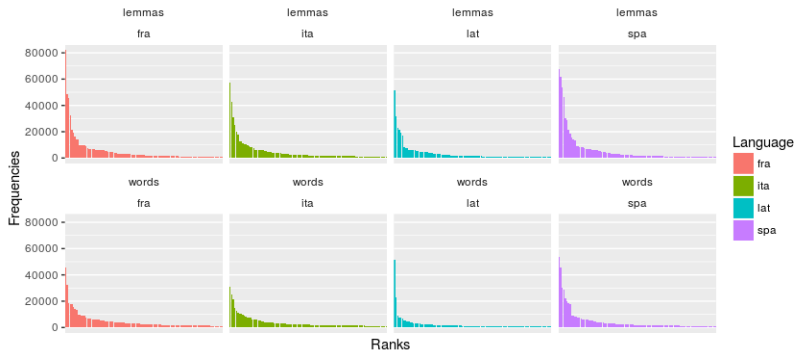
Input	→	Tag	Lemma
fratrem	→	N:acc	frater
fratris	→	N:gen	frater
fratribus	→	N:dat	frater
vivit	→	V:IND	vivo
movetur	→	V:IND	moveo
humanum	→	ADJ	humanus

...

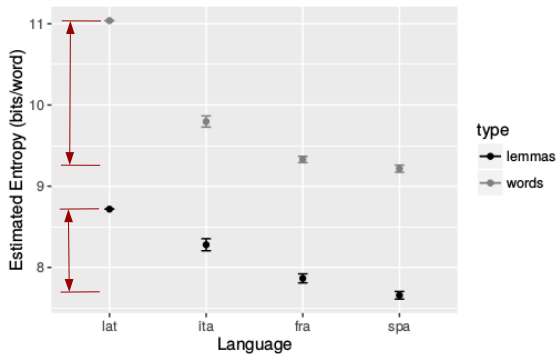
CORPORA

Select the **Romance languages** (+ Latin) - which can be **lemmatized with the TreeTagger** - from the *Parallel Bible Corpus* (4 languages: Latin, French, Spanish, Italian)

RESULTS (LEMMATIZATION)



RESULTS (LEMMATIZATION)



CONCLUSIONS (LEMMATIZATION)

- ▶ In **Latin, Italian, Spanish and French** entropy (bits/word) is reduced by **ca. 15-20% through lemmatization**, i.e. when inflectional marking is neutralized.
- ▶ The **entropy difference to Latin** is reduced via lemmatization by **ca. 50%**

ITERATED LEARNING EXPERIMENTS

Third Analysis

Illustrate via **iterated learning experiments** how inflectional marking is lost through **learning pressures (non-native, i.e. L2)** over several generations

EXPERIMENTAL DESIGN (BERDICEVSKIS & SEMENUKS, FORTHCOMING)

Setup

- ▶ **artificial language learning** task
- ▶ Overall **300 participants**
(3 types of chains × 10 generations × 10 subjects)
- ▶ native speakers of **Russian**
- ▶ experiment on **webpage** (*jsPsych* javascript)

ARTIFICIAL LANGUAGE: EPSILON



segn



















segn bv-OAGR



segn-lpl bv-OAGR

(Berdicevskis & Semenuks, forthcoming)

ARTIFICIAL LANGUAGE: EPSILON

		event: none	event: fall apart	event: grow antlers	event: fly
agent: round animal	number: singular	 segn	 segn mv-OAGR	 segn tv-OAGR	 segn bv-OAGR
	number: plural	 segn-lpl	 segn-lpl mv-OAGR	 segn-lpl tv-OAGR	 segn-lpl bv-OAGR
agent: square animal	number: singular	 fuvN	 fuvN mv-iAGR	 fuvN tv-iAGR	 fuvN bv-iAGR
	number: plural	 fuvN-lpl	 fuvN-lpl mv-iAGR	 fuvN-lpl tv-iAGR	 fuvN-lpl bv-iAGR

(Berdicevskis & Semenuks, forthcoming)

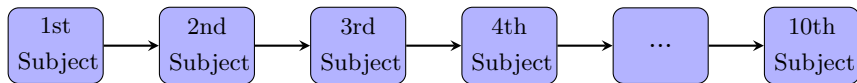
ARTIFICIAL LANGUAGE: EPSILON

“This system **resembles the more complex Russian morphosyntactic system** where nouns are marked for number, and adjectives and verbs agree with nouns in number and gender. The bottom line is that agreement is salient and pervasive in Russian morphosyntax, and thus the **mother tongue is not imposing pressure** on the participants to shed agreement.” (Berdicevskis & Semenuks, forthcoming)

EXPERIMENTAL DESIGN (BERDICEVSKIS & SEMENUKS, FORTHCOMING)

Condition I

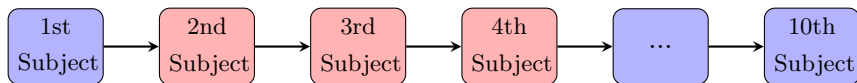
Set of **10 uninterrupted chains** of 10 “generations”, i.e. subjects



EXPERIMENTAL DESIGN (BERDICEVSKIS & SEMENUKS, FORTHCOMING)

Condition II

Set of **10 temporarily interrupted chains**

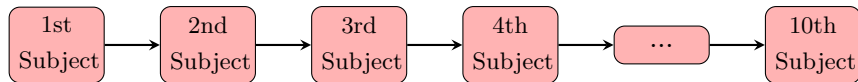


Interruption (i.e. L2 influence) is here introduced via **less exposure** to the target language in the learning phase (6 versus 3 learning blocks).

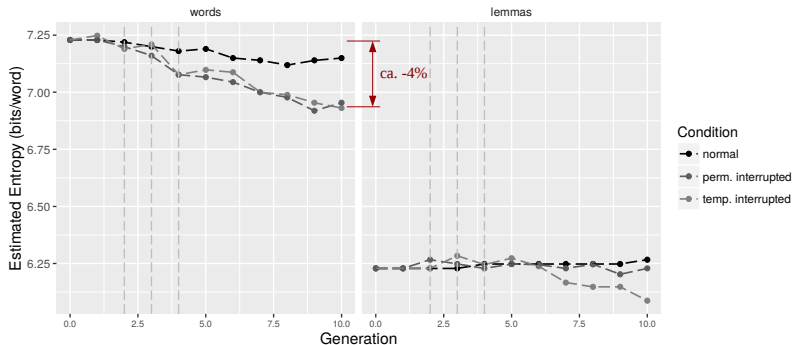
EXPERIMENTAL DESIGN (BERDICEVSKIS & SEMENUKS, FORTHCOMING)

Condition III

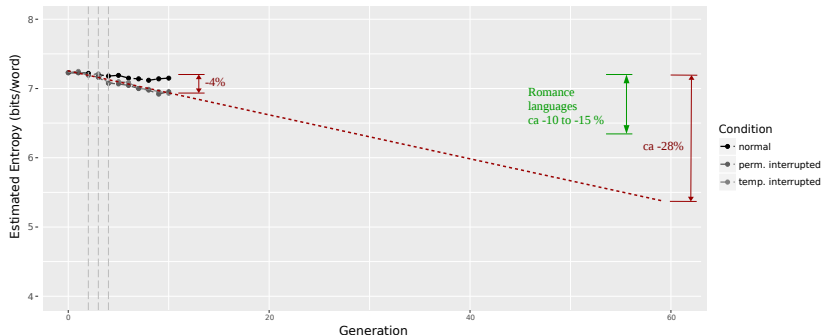
Set of **10 permanently interrupted chains**



RESULTS: ENTROPY IN EPSILON



RESULTS: “EXTRAPOLATION”



60 generations \sim 1800 years

CONCLUSIONS (ITERATED LEARNING)

- ▶ in the interrupted conditions word **entropy is reduced by around 4%** in 10 generations
- ▶ this reduction is mostly due to **loss of inflection** rather than **base vocabulary**

CONCLUSIONS

Corpus
Measure

Diachronic reduction of word
entropy (Romance languages)

Lemmatization

Reduction is (largely) due to
loss of inflectional marking

Experiments

Artificial languages shed
inflectional morphology
via learning pressure

THANK YOU!

