

# Adaptive Languages

## An Information-Theoretic Account of Linguistic Diversity

Christian Bentz  
*University of Tübingen*

January 17, 2017



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

**EVO LAEMP**  
LANGUAGE EVOLUTION: THE EMPIRICAL TURN

# OVERVIEW

Information and Language

Applications to Typology

Explanations of Diversity

Conclusions

# INTUITIVE TERMINOLOGY

- ▶ order  $\leftrightarrow$  disorder
- ▶ regularity  $\leftrightarrow$  irregularity
- ▶ predictability  $\leftrightarrow$  unpredictability
- ▶ certainty  $\leftrightarrow$  uncertainty

} Entropy

“*Entropy as possibility* is my favorite short description of entropy because (...) unlike *uncertainty* and *missing information*, it has positive connotation.”

“*Entropy as possibility* is my favorite short description of entropy because (...) unlike *uncertainty* and *missing information*, it has positive connotation.”

“Entropy is an additive measure of the  
**number of possibilities**  
available to a system.”

Lemons, 2013

# ENTROPY AS CHOICE

## Minimal

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaaa

# ENTROPY AS CHOICE

## Minimal

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

aaaaaaaaaaaaaaaaaaaaaaaa

## Maximal

fcbihspm hkgiwlelbj

sdmkfuuf cvkym cfcsqdvcc

trdgjmpnkjhujril

unnapsfmgbk ggqvntxprl

kfkmpsgjetn

grycfjuxxcusejlexfhkfrmh

jknecxjgg isonkqcwmxr

ymwwuieumi brlromnqyq

yclvlkmtgfd fcmvulfkyawa

# ENTROPY AS CHOICE

## Minimal

aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa  
aaaaaaaaaaaaaaaaaaaaaaaaa

## Maximal

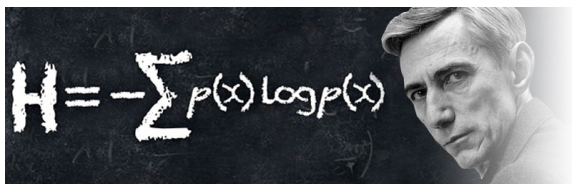
fcbihspm hkgiwlelbj  
sdmkuuf cvkym cfcsqdvcc  
trdgjmpnkjhujril  
unnapsfmgbk ggqvntxprl  
kfkmpsgjetn  
grycfjuxxcusejlexfhkfrmh  
jknecxjgg isonkqcmwxr  
ymwuiemi brlromnqyq  
yclvlkmtgfd fcmvulfkyawa

## Intermediate

all human beings are born  
free and equal in dignity  
and rights they are endowed  
with reason and conscience  
and should act towards one  
another in a spirit of  
brotherhood everyone is  
entitled to all the rights  
and freedoms



# SHANNON ENTROPY



Shannon & Weaver (1949) The mathematical theory of communication

# SHANNON ENTROPY

$$H(X) = - \sum_{i=1}^N p(x_i) \log_2 p(x_i) \quad (1)$$

- $-\log_2 p(x_i)$  is the **information content** of a unit  $x_i$  (e.g. word type).

# SHANNON ENTROPY

$$H(X) = - \sum_{i=1}^N p(x_i) \log_2 p(x_i) \quad (1)$$

- ▶  $-\log_2 p(x_i)$  is the **information content** of a unit  $x_i$  (e.g. word type).
- ▶ For example:
  - “human” in the UDHR:  $-\log_2(\frac{13}{2000}) \sim 7.27$
  - “the” in the UDHR:  $-\log_2(\frac{121}{2000}) \sim 4.05$

# SHANNON ENTROPY

$$H(X) = - \sum_{i=1}^N p(x_i) \log_2 p(x_i) \quad (2)$$

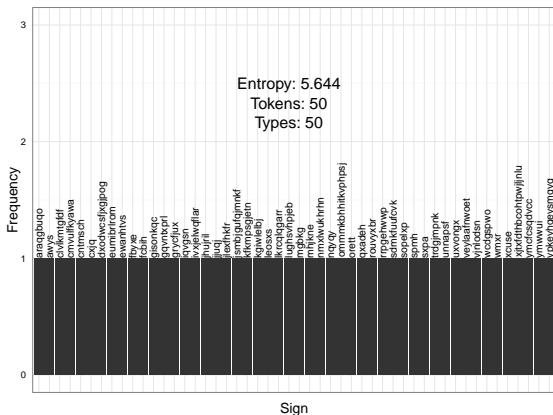
- The entropy is the **average information content** of information encoding units:  $H(X) = [0, \infty[$

# MINIMAL $H(X)$

aaa  
aaa  
aaa  
aaa

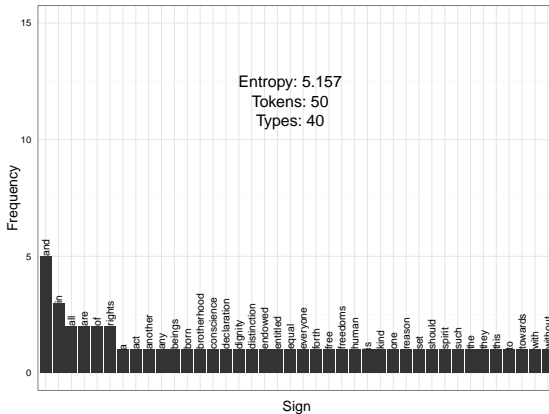
$$H(X) = 1 \times \log_2(1) = 0 \quad (3)$$

# RANDOM TEXT



$$H(X) = -\left(\frac{1}{50} \log_2\left(\frac{1}{50}\right) + \dots + \frac{1}{50} \log_2\left(\frac{1}{50}\right)\right) = 5.644 \quad (4)$$

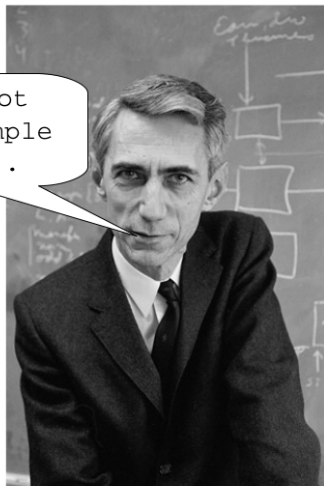
# ENGLISH UDHR



$$H(X) = -\left(\frac{5}{50} \log_2\left(\frac{5}{50}\right) + \frac{3}{50} \log_2\left(\frac{3}{50}\right) + \dots + \frac{1}{50} \log_2\left(\frac{1}{50}\right)\right) = 5.157 \quad (5)$$

# METHODOLOGICAL ISSUES

But it's not  
quite as simple  
as that...





# METHODOLOGICAL ISSUES

- ▶  $H(X)$  depends on **text size** (number of tokens)

# METHODOLOGICAL ISSUES

- ▶  $H(X)$  depends on **text size** (number of tokens)
- ▶ the **probability** of words is not a simple function of their **frequency**, they depend on **co-text**

# METHODOLOGICAL ISSUES

- ▶  $H(X)$  depends on **text size** (number of tokens)
- ▶ the **probability** of words is not a simple function of their **frequency**, they depend on **co-text**
- ▶ what are “**words**” anyways?

Haspelmath (2011) The indeterminacy of word segmentation

Wray (2014) Why are we so sure we know what a word is?

# METHODOLOGICAL ISSUES

*Article*

## The entropy of words – Estimations across more than 1000 languages

Christian Bentz <sup>1,2\*</sup>, Dimitrios Alikaniotis <sup>3</sup>, Michael Cysouw <sup>4</sup> and Ramon Ferrer-i-Cancho <sup>5</sup>

<sup>1</sup> DFG Center for Advanced Studies, University of Tübingen, Rümelinstraße 23, D-72070 Tübingen, Germany; [chris@christianbentz.de](mailto:chris@christianbentz.de)

<sup>2</sup> Department of General Linguistics, University of Tübingen, Wilhelmstraße 19-23, D-72074 Tübingen, Germany

<sup>3</sup> Department of Theoretical and Applied Linguistics, University of Cambridge, Cambridge, United Kingdom

<sup>4</sup> Forschungszentrum Deutscher Sprachatlas, Philipps-Universität Marburg, Deutschhausstrasse 3, 35037 Marburg

<sup>5</sup> Complexity and Quantitative Linguistics Lab, LARCA Research Group, Departament de Ciències de la Computació, Universitat Politècnica de Catalunya, Barcelona, Catalonia, Spain

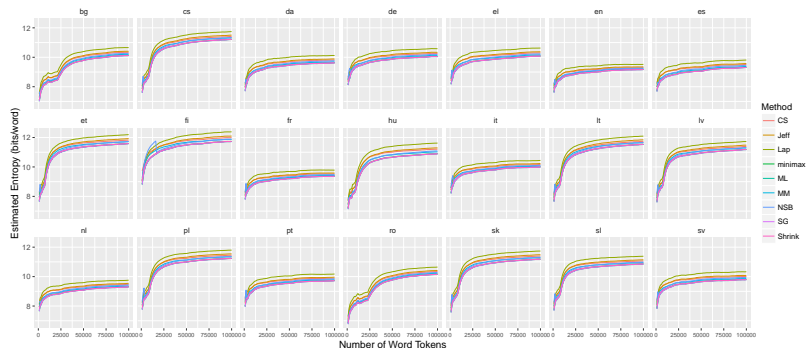
\* Correspondence: [chris@christianbentz.de](mailto:chris@christianbentz.de)

Academic Editor: name

Version January 9, 2017 submitted to Entropy; Typeset by L<sup>A</sup>T<sub>E</sub>X using class file mdpi.cls

- 1 **Abstract:** The uncertainty associated with words is a fundamental property of natural languages.
- 2 It lies at the heart of quantitative linguistics, computational linguistics, and language sciences
- 3 more generally. Information-theory gives us tools at hand to measure precisely this uncertainty
- 4 – the word entropy. We here use three parallel corpora – encompassing ca. 450 million words in

# TEXT SIZE DEPENDENCE



Bentz, Alikaniotis, Cysouw & Ferrer-i-Cancho (forthcoming)

# MISCONCEPTION

“(...) as many critics have since noted, and as Shannon was well aware, this model is not appropriate as a model of human communication, because “information” in Shannon’s technical sense is not equivalent to “meaning” in *any* sense.”

Fitch (2011) The Evolution of Language

# THE MATHEMATICAL THEORY OF COMMUNICATION

**LEVEL A.** How accurately can the symbols of communication be transmitted? (The technical problem.)

**LEVEL B.** How precisely do the transmitted symbols convey the desired meaning? (The semantic problem.)

**LEVEL C.** How effectively does the received meaning affect conduct in the desired way? (The effectiveness problem.)

Shannon & Weaver (1949), p.4

“The mathematical theory of the engineering aspect of communication (...) admittedly applies in the first instance only to problem A (...)”

However



“The mathematical theory of the engineering aspect of communication (...) admittedly applies in the first instance only to problem A (...)”

However

“(...) levels B and C, above, can make use only of those signal accuracies which turn out to be possible when analyzed at Level A. Thus any limitations discovered in the theory at Level A necessarily apply to levels B and C.”

Shannon & Weaver (1949), p.6

“(...) the analysis at Level A discloses that this level overlaps the other levels more than one could possibly naively suspect. Thus the theory of Level A is, at least to a significant degree, also a theory of levels B and C.”

Shannon & Weaver (1949), p.6

Entropy is a **necessary** but **not sufficient** condition for communication

# APPLICATIONS TO TYPOLOGY

- ▶ What is the **information encoding potential** (entropy share) of different linguistic features?

# APPLICATIONS TO TYPOLOGY

- ▶ What is the **information encoding potential** (entropy share) of different linguistic features?
- ▶ How do **information encoding strategies** differ across languages?

# PARALLEL CORPORA

```
# language_name:      English
# closest ISO 639-3:  eng
# year_short:         1890
# year_long:          Not available
# title:              The Bible in English, Darby Translation
# URL:                http://unbound.biola.edu/index.cfm?method=downloa
# copyright_short:    © Public Domain
# copyright_long:     First published in 1890 by John Nelson Darby, an
Darby also published translations of the Bible in French and German.
01001001      In the beginning God created the heavens and the earth .
01001002      And the earth was waste and empty , and darkness was on t
01001003      And God said , Let there be light . And there was light .
01001004      And God saw the light that it was good ; and God divided
01001005      And God called the light Day , and the darkness he called
01001006      And God said , Let there be an expanse in the midst of th
01001007      And God made the expanse , and divided between the waters
01001008      And God called the expanse Heavens . And there was evenin
01001009      And God said , Let the waters under the heavens be gather
01001010      And God called the dry [ land ] Earth , and the gathering
```

# PARALLEL CORPORA

```
# language_name:      Amharic
# closest ISO 639-3:  amh
# year_short:         1994
# year_long:          E-Text in transliterated ASCII format by Lapsley/Brooks
                        (www.nt-text.net). Revised Amharic Bible in XML ( 2003 ).
# title:              The New Testament in Amharic
# URL:                http://unbound.biola.edu/index.cfm?method=downloads.show
# copyright_short:    © Printed Version by United Bible Societies 1962
# copyright_long:     Not available
```

40001001 የዳዊት ልጅ የአብርሃም ልጅ የኢየሱስ ክርስቶስ ትውልድ መጽሐፍ ።

40001002 አብርሃም ይስሐቅን ወለደ ፤ ይስሐቅም ያዕቆብን ወለደ ፤ ያዕቆብም ይሁዳንና ወንድሞቹን ወ፤

40001003 ይሁዳም ከትዕማር ፋሬስንና ዛሬን ወለደ ፤ ፋሬስም ኢስሮምን ወለደ ፤

40001004 ኢስሮምም አራምን ወለደ ፤ አራምም ለሚናዳብን ወለደ ፤ ለሚናዳብም ነአሶንን ወለደ ፤ ነአሶ

40001005 ሰልሞንም ከራኩብ ቦሌዝን ወለደ ፤ ቦሌዝም ከሩት ኢየቤድን ወለደ ፤ ኢየቤድም ለሴይን ወለ

40001006 ለሴይም ንጉሥ ዳዊትን ወለደ ።

40001007 ሰሎሞንም ሮብዓምን ወለደ ፤ ሮብዓምም አቢያን ወለደ ፤ አቢያም ለሣፍን ወለደ ፤

40001008 ለሣፍም ኢዮሣፍጥን ወለደ ፤ ኢዮሣፍጥም ኢዮራምን ወለደ ፤ ኢዮራምም ያዝያንን ወለደ ፤

40001009 ያዝያንም ኢዮአታምን ወለደ ፤ ኢዮአታምም ለካዝን ወለደ ፤

40001010 ለካዝም ሕዝቅያስን ወለደ ፤ ሕዝቅያስም ምናሲን ወለደ ፤ ምናሲም ለሞዕን ወለደ ፤

40001011 ለሞዕም ኢዮስያስን ወለደ ፤ ኢዮስያስም በባቢሎን ምርኮ ጊዜ ኢኮንያንንና ወንድሞቹን ወለደ

# PARALLEL CORPORA

```
# language_name: ភាសាខ្មែរ
# closest ISO 639-3:  khm
# year_short:      2011
# year_long:       Not available
# title:           Khmer Christian Bible<br>The New Testament in Khmer
# URL:             https://www.bible.com/de/bible/315/mat.1.kcb
# copyright_short:  © Words of Life Ministries 2011
# copyright_long:   Khmer Christian Bible<br>Copyright © Holy Bible, Khmer Christian Bible
#                  copyright 2011 by Words of Life Ministries, P.O. Box 2581, Phnom Penh,
#                  3, Cambodia. All rights reserved.
```

40001001 កំណត់ត្រាវង្សត្រកូលរបស់ព្រះយេស៊ូគ្រីស្ទដែលជាពូជពង្សរបស់ស្តេចដាវីឌ និងលោកអំប្រាហាំ :

40001002 លោកអំប្រាហាំបង្កើតលោកអ៊ីសាក លោកអ៊ីសាកបង្កើតលោកយ៉ាកុប លោកយ៉ាកុបបង្កើតលោកយូដា និង

40001003 លោកយូដា និងនាងតាម៉ារបង្កើតលោកពេរេស និងលោកសេវ៉ាស ឯលោកពេរេសបង្កើតលោកហេស្រ្តុន

40001004 លោករ៉ាមបង្កើតលោកអ៊ីមីណាដាប់ លោកអ៊ីមីណាដាប់បង្កើតលោកណាសូន លោកណាសូនបង្កើតលោក

40001005 លោកសាលម៉ូន និងនាងរ៉ាហាប់បង្កើតលោកបូអូស ហើយលោកបូអូស និងនាងរស់បង្កើតលោកអូប៊ិឌ ។

40001006 លោកអ៊ីសាយបង្កើតស្តេចដាវីឌ ស្តេចដាវីឌ និងប្រពន្ធលោកអ៊ូរីបង្កើតស្តេចសាឡូម៉ូន

40001007 ស្តេចសាឡូម៉ូនបង្កើតស្តេចអេហ្វោម ស្តេចអេហ្វោមបង្កើតស្តេចអ៊ីប៊ីយ៉ា ស្តេចអ៊ីប៊ីយ៉ាបង្កើតស្តេចអេសេ

40001008 ស្តេចអេសេបង្កើតស្តេចយ៉ូសាផាត ស្តេចយ៉ូសាផាតបង្កើតស្តេចយ៉ូរាម ស្តេចយ៉ូរាមបង្កើតស្តេចអូសៀស

40001009 ស្តេចអូសៀសបង្កើតស្តេចយ៉ូថាម ស្តេចយ៉ូថាមបង្កើតស្តេចអេហាស ស្តេចអេហាសបង្កើតស្តេចអេសេគ

40001010 ស្តេចអេសេគាសបង្កើតស្តេចម៉ាណាសេ ស្តេចម៉ាណាសេបង្កើតស្តេចអាំម៉ូន ស្តេចអាំម៉ូនបង្កើតស្តេចយ៉ូ



# PARALLEL CORPORA

- ▶ *Parallel Bible Corpus* (PBC): ca. **1200** lang.  
Mayer & Cysouw (2014)
- ▶ *Universal Declaration of Human Rights* (UDHR): ca. **400** lang.  
[www.unicode.org/udhr](http://www.unicode.org/udhr)
- ▶ *Open Subtitles Corpus* (OSC): ca. **100** lang.  
Tiedemann (2012)
- ▶ *Europarl Parallel Corpus* (EPC): **21** lang.  
Koehn (2005)

# APPLICATIONS TO TYPOLOGY

JOURNAL OF QUANTITATIVE LINGUISTICS, 2016  
<http://dx.doi.org/10.1080/09296174.2016.1265792>

 **Routledge**  
Taylor & Francis Group

## Variation in Word Frequency Distributions: Definitions, Measures and Implications for a Corpus- Based Language Typology

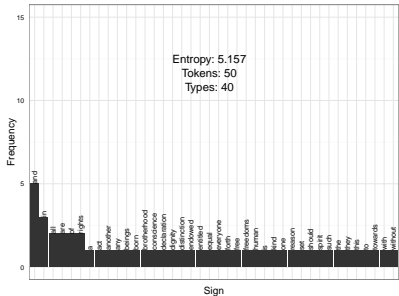
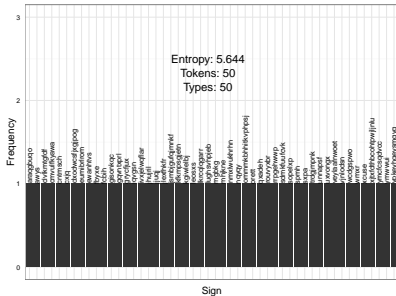
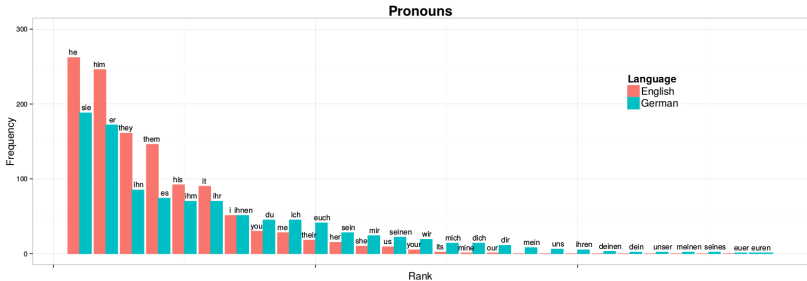
Christian Bentz<sup>a,b</sup>, Dimitrios Alikaniotis<sup>a</sup>, Tanja Samardžić<sup>c</sup> and  
Paula Buttery<sup>a</sup>

<sup>a</sup>Department of Theoretical and Applied Linguistics, University of Cambridge, Cambridge, United Kingdom; <sup>b</sup>Department of General Linguistics, University of Tübingen, Tübingen, Germany; <sup>c</sup>URPP on Language and Space, University of Zürich, Zürich, Switzerland

### ABSTRACT

Word frequencies are central to linguistic studies investigating processing difficulty, learnability, age of acquisition, diachronic transmission and the relative weight given to a concept in society. However, there are few cross-linguistic studies on entire distributions of word frequencies, and even less on systematic changes within them. Here, we first define and test an exact measure for the relative

## ENTROPY DIFFERENCES



# ANALYSIS 1

- ▶ Languages: English and German
- ▶ Corpora: OSC, Bible, UDHR
- ▶ **manually neutralize** inflections, derivations, compounds, clitics/contractions

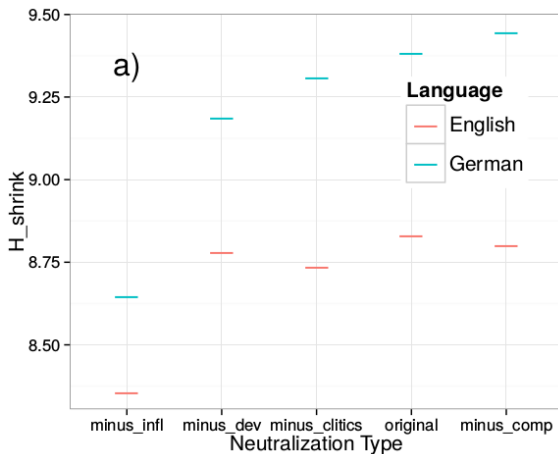
# ANALYSIS 1

## Open Subtitles Corpus

Go back to Oslo and meet up with your  
boss, report what you've seen here, and  
ask him to send you back.  
It's about Hellfjord.  
There is someone in this room who has  
read too much Donald Duck.  
It's quiet without Salmander.  
Oh my god.  
Hey, what are you doing?  
What is this?  
This is no fish eye.  
Chop off an arm and a leg, and blame  
the sea serpent.

Fahr nach Oslo, triff dich mit deinem  
Chef und erstatte ihm Bericht, - - was  
du hier gesehen hast, und bitte ihn,  
dich zurückzuschicken.  
Es geht um Hellfjord.  
In diesem Raum befindet sich jemand,  
der zu viel Donald Duck gelesen hat.  
Es ist so still hier ohne Salmander.  
Oh mein Gott.  
Hey, was machst du in meiner Küche?  
Was ist das?  
Das sind keine Fischeaugen.  
Hackt ihr einen Arm und ein Bein ab und  
beschuldigt das Seeungeheuer.

# ANALYSIS 1



# ANALYSIS 1: CONCLUSIONS

- **Inflectional marking** has the biggest **entropy share** in both languages, i.e. ca. 0.8 bits/word (9%) in German, and ca. 0.5 bits/word (6%) in English

# ANALYSIS 1: CONCLUSIONS

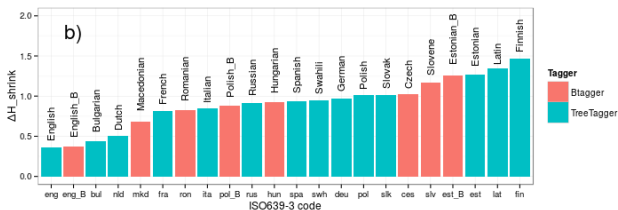
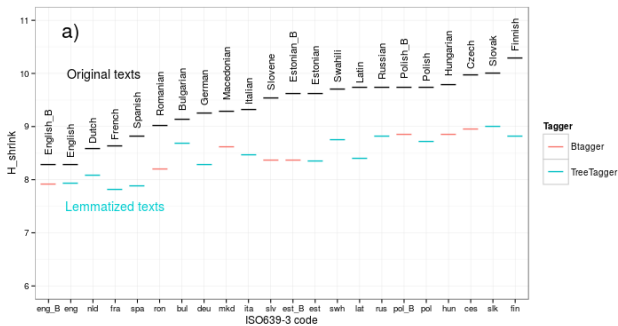
- ▶ **Inflectional marking** has the biggest **entropy share** in both languages, i.e. ca. 0.8 bits/word (9%) in German, and ca. 0.5 bits/word (6%) in English
- ▶  $\Delta H$  between **English and German** is due to inflections (48% decrease), derivations (27% decrease), compounds (14% increase), contractions/clitics (2% increase)



# ANALYSIS 2

- ▶ Languages: 19
- ▶ Corpora: PBC, UDHR
- ▶ **automatically neutralize** inflections via lemmatization

# ANALYSIS 2



## ANALYSIS 2: CONCLUSIONS

- ▶ The **entropy share** of inflections across 19 languages ranges from 0.4 bits/word (5%) in English to 1.5 bits/word (15%) in Finnish.

## ANALYSIS 2: CONCLUSIONS

- ▶ The **entropy share** of inflections across 19 languages ranges from 0.4 bits/word (5%) in English to 1.5 bits/word (15%) in Finnish.
- ▶  $\Delta H$  between **19 languages** (Indo-European, Uralic, Atlantic-Congo) is explained to **55% by inflectional marking** differences

# FURTHER ANALYSES: TONE MARKING?

## Usila Chinantec (Otomanguen)

40001001 I<sup>4</sup>la<sup>3</sup> ti<sup>2</sup>ton<sup>3</sup> la<sup>4</sup>jang<sup>34</sup> sa<sup>1</sup>jeun<sup>3</sup> quian<sup>1</sup> Jesucristo a<sup>3</sup>lang<sup>43</sup> jon<sup>43</sup>tyie<sup>1</sup> A<sup>3</sup>br  
 40001002 A<sup>3</sup>brang<sup>23</sup> lang<sup>43</sup> jmai<sup>3</sup> I<sup>3</sup>sa<sup>23</sup> . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> I<sup>3</sup>sa<sup>23</sup> a<sup>4</sup>hyon<sup>23</sup> Ja<sup>3</sup>co<sup>23</sup> . I<sup>2</sup>  
 40001003 I<sup>2</sup> i<sup>2</sup>con<sup>23</sup> Judá jian<sup>23</sup> A<sup>1</sup>ta<sup>3</sup>mar<sup>23</sup> ra<sup>5</sup>sian<sup>3</sup> Fares jian<sup>3</sup> Zara . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup>  
 40001004 I<sup>2</sup> i<sup>2</sup>con<sup>23</sup> Aram ja<sup>34</sup> Aminadab . I<sup>2</sup> Aminadab ja<sup>34</sup> Naasón i<sup>2</sup>con<sup>23</sup>i<sup>3</sup> . Jian<sup>3</sup> i  
 40001005 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Salmón jian<sup>23</sup> A<sup>1</sup>ra<sup>3</sup>hab<sup>23</sup> ra<sup>5</sup>sian<sup>3</sup> Booz . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Booz j  
 40001006 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Isai hain<sup>4</sup> ra<sup>5</sup>sian<sup>3</sup> re<sup>1</sup> Da<sup>3</sup>vei<sup>23</sup> . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> re<sup>1</sup> Da<sup>3</sup>vei<sup>23</sup>  
 40001007 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Salomón ra<sup>5</sup>sian<sup>3</sup> Roboam . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Roboam ra<sup>5</sup>sian<sup>3</sup> Abías  
 40001008 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Asa ra<sup>5</sup>sian<sup>3</sup> Josafat . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Josafat ja<sup>34</sup> Joram . Jia  
 40001009 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Uzias ja<sup>34</sup> Jotam . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Jotam ja<sup>34</sup> Acáz . Jian<sup>3</sup> i<sup>2</sup>co  
 40001010 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Ezequías ra<sup>5</sup>sian<sup>3</sup> Manasés . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Manasés ra<sup>5</sup>sian<sup>3</sup> A<sup>3</sup>  
 40001011 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Josías hain<sup>4</sup> cuan<sup>34</sup> Jeconías jian<sup>23</sup> si<sup>3</sup>reunh<sup>1</sup> ma<sup>2</sup>a<sup>4</sup>han<sup>5</sup> Israe  
 40001012 I<sup>2</sup> coh<sup>5</sup> a<sup>4</sup>húan<sup>3</sup>i<sup>3</sup> chion<sup>32</sup> jeu<sup>3</sup> Babilonia jon<sup>3</sup> , jon<sup>3</sup> ra<sup>5</sup>sian<sup>3</sup> a<sup>3</sup>jon<sup>43</sup> Jecon  
 40001013 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Zorobabel ra<sup>5</sup>sian<sup>3</sup> Abiud . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Abiud ja<sup>34</sup> Eliaquim  
 40001014 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Azor ja<sup>34</sup> Sadoc . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Sadoc ja<sup>34</sup> Aquim . Jian<sup>3</sup> i<sup>2</sup>co  
 40001015 Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Eliud ra<sup>5</sup>sian<sup>3</sup> Eleazar . Jian<sup>3</sup> i<sup>2</sup>con<sup>23</sup> Eleazar ra<sup>5</sup>sian<sup>3</sup> Matán  
 40001016 I<sup>2</sup> i<sup>2</sup>con<sup>23</sup> Ja<sup>3</sup>co<sup>23</sup> hain<sup>4</sup> ra<sup>5</sup>sian<sup>3</sup> Se<sup>1</sup> , a<sup>3</sup>hain<sup>4</sup> i<sup>3</sup>cúa<sup>3</sup> Ma<sup>3</sup>rei<sup>2</sup> a<sup>3</sup>cuan<sup>34</sup> Jes  
 40001017 I<sup>2</sup> la<sup>4</sup>ne<sup>3</sup> ja<sup>34</sup> quia<sup>5</sup>quin<sup>4</sup> sa<sup>1</sup>jeun<sup>3</sup> liah<sup>4</sup>ma<sup>2</sup>sian<sup>3</sup> A<sup>3</sup>brang<sup>23</sup> la<sup>4</sup>teg<sup>4</sup> Da<sup>3</sup>vei<sup>23</sup>  
 chion<sup>32</sup> ta<sup>5</sup> Babilonia ; conh<sup>4</sup>liah<sup>4</sup> i<sup>2</sup>quia<sup>5</sup>quin<sup>4</sup> liah<sup>4</sup>ma<sup>2</sup>tionh<sup>3</sup>i<sup>3</sup> chion<sup>32</sup> Babilonia la<sup>4</sup>teg<sup>4</sup>  
 40001018 La<sup>4</sup>la<sup>3</sup> ra<sup>5</sup>sian<sup>3</sup> Jesucristo : A<sup>4</sup>leg<sup>34</sup> re<sup>3</sup> i<sup>4</sup>sanh<sup>4</sup> Ma<sup>3</sup>rei<sup>2</sup> sie<sup>23</sup>i<sup>3</sup> jian<sup>23</sup> Se<sup>1</sup>

# FURTHER ANALYSES: WORD ORDER?

## Word Order Typology through Multilingual Word Alignment

**Robert Östling**

Department of Linguistics

Stockholm University

SE-106 91 Stockholm, Sweden

`robert@ling.su.se`

### Abstract

With massively parallel corpora of hundreds or thousands of translations of the same text, it is possible to automatically perform typological studies of language structure using very large language samples. We investigate the domain of word order using multilingual word alignment and high-precision annotation transfer in a corpus with 1144 translations in 986 languages of the New Testament. Results are

jection of part-of-speech (PoS) tags and dependency parse trees to investigate five different word order properties in 986 different languages, through a corpus of New Testament translations. The results are validated through comparison to relevant chapters in the World Atlas on Language Structures, WALS (Dryer and Haspelmath, 2013), and we find a very high level of agreement between this database and our method.

We identify two primary applications of this method. First, it provides a new tool for basic research in linguistic typology. Second, it has been

# APPLICATIONS TO TYPOLOGY: CONCLUSIONS

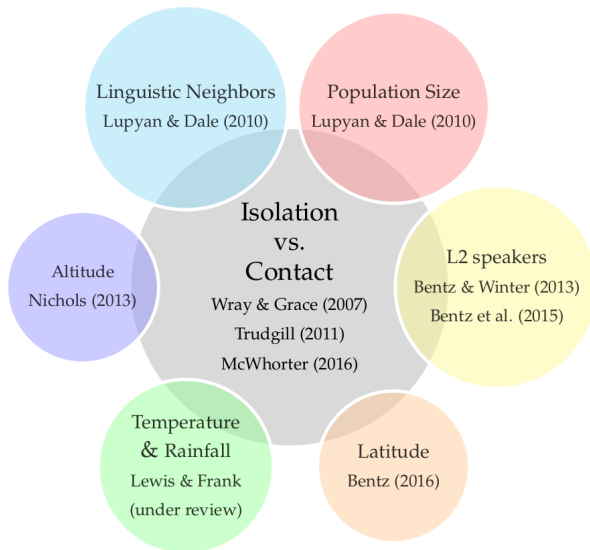
- ▶ The **information encoding potential** of *any linguistic feature* can be measured, if there is a systematic way of manipulating it in texts

# APPLICATIONS TO TYPOLOGY: CONCLUSIONS

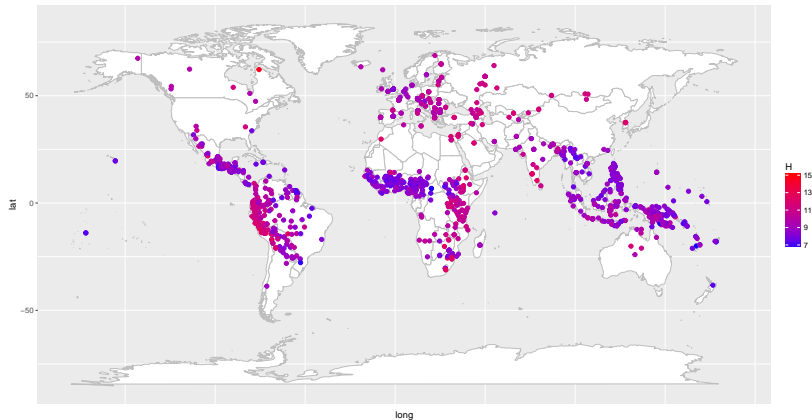
- ▶ The **information encoding potential** of *any linguistic feature* can be measured, if there is a systematic way of manipulating it in texts
- ▶ The analyses do **not have to depend** on **word types**, they could also use characters, morphemes, constructions, etc.



# EXPLANATIONS OF DIVERSITY



## Word Entropy across 1092 languages



Bentz, Alikaniotis, Cysouw & Ferrer-i-Cancho (forthcoming)

# NON-NATIVE LANGUAGE LEARNING



## RESEARCH ARTICLE

# Adaptive Communication: Languages with More Non-Native Speakers Tend to Have Fewer Word Forms

**Christian Bentz<sup>1\*</sup>, Annemarie Verkerk<sup>2</sup>, Douwe Kiela<sup>3</sup>, Felix Hill<sup>3</sup>, Paula Buttery<sup>1,3</sup>**

**1** Department of Theoretical and Applied Linguistics, University of Cambridge, Cambridge, United Kingdom, **2** Reading Evolutionary Biology Group, School of Biological Sciences, University of Reading, Reading, United Kingdom, **3** Computer Laboratory, University of Cambridge, Cambridge, United Kingdom

\* [chris@christianbentz.de](mailto:chris@christianbentz.de)



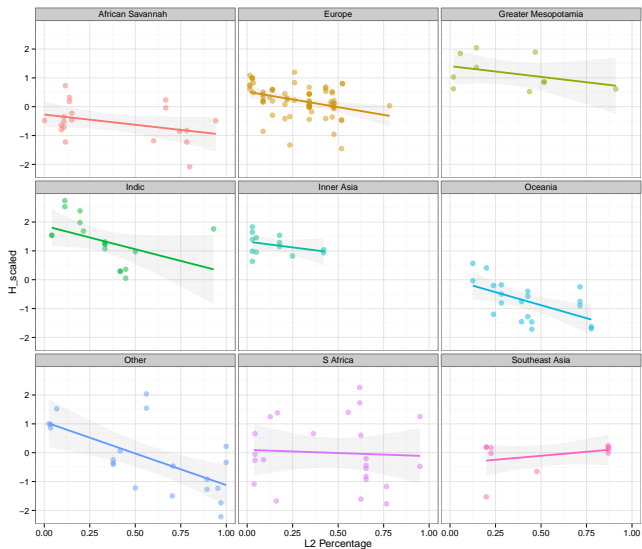
## Abstract

Explaining the diversity of languages across the world is one of the central aims of typological, historical, and evolutionary linguistics. We consider the effect of *language contact*—the number of non-native speakers a language has—on the way languages change and evolve. By analysing hundreds of languages within and across language families, regions, and text types, we show that languages with greater levels of contact typically employ fewer word

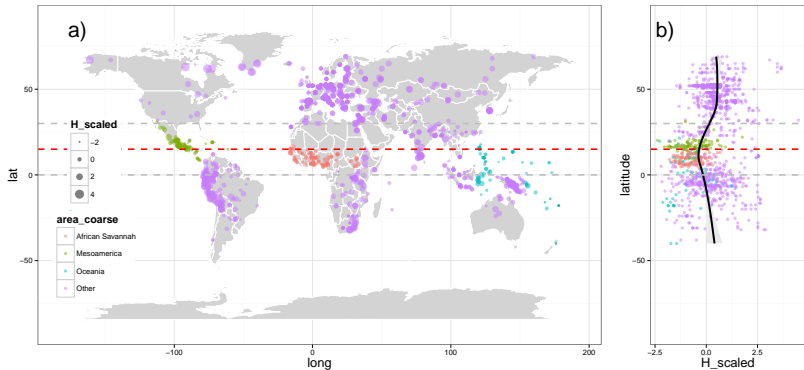
## OPEN ACCESS

**Citation:** Bentz C, Verkerk A, Kiela D, Hill F, Buttery P (2015) Adaptive Communication: Languages with More Non-Native Speakers Tend to Have Fewer

# NON-NATIVE LANGUAGE LEARNING

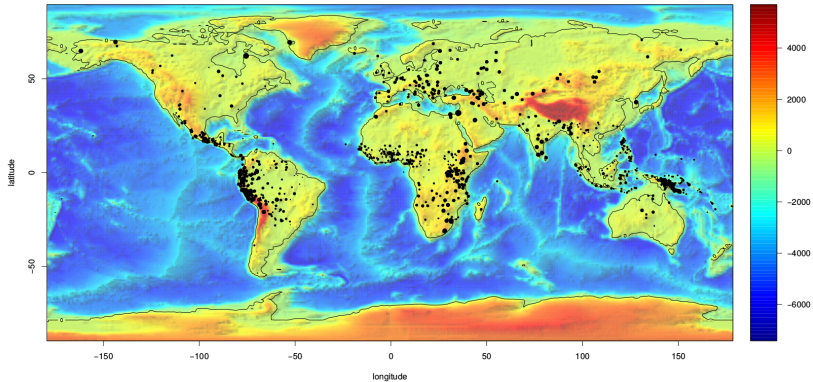


# LATITUDE



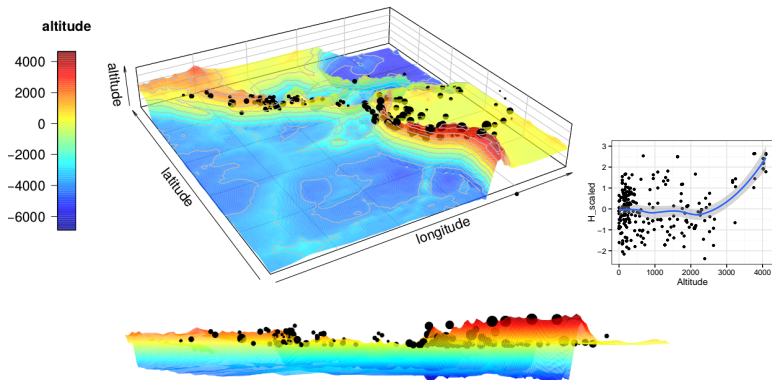
Bentz (2016) The Low-Complexity-Belt

# ALTITUDE



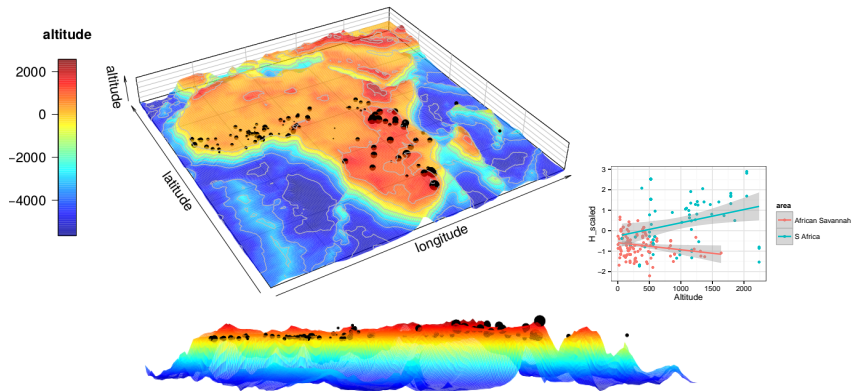
# MESOAMERICA AND THE ANDES

183 LANGUAGES, 90 FAMILIES, 3 AREAS



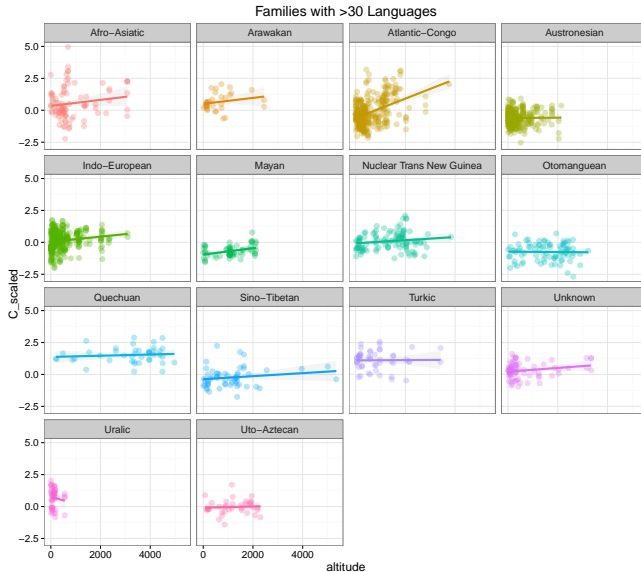
# AFRICAN SAVANNAH AND SOUTH AFRICA

127 LANGUAGES, 21 FAMILIES, 2 AREAS



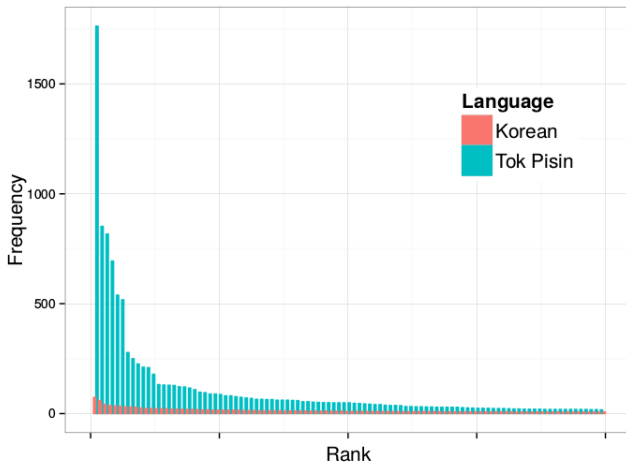


# ALTITUDE



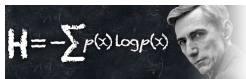
# DISCLAIMER

What's better?




# CONCLUSIONS

- Information theory *is* relevant to natural languages

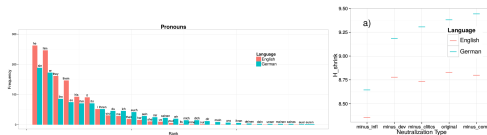

$$H = -\sum p(x) \log p(x)$$

# CONCLUSIONS

- Information theory *is* relevant to natural languages


$$H = -\sum p(x) \log p(x)$$


- We can determine the *entropy share* of linguistic features

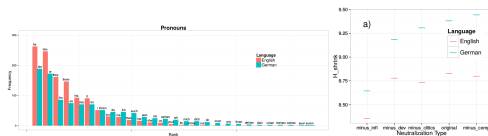


# CONCLUSIONS

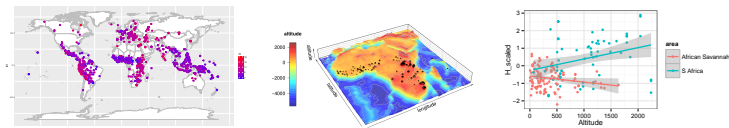
- Information theory *is* relevant to natural languages

$$H = -\sum p(x) \log p(x)$$


- We can determine the *entropy share* of linguistic features



- We start to understand *entropy diversity*



# THANK YOU

