Quantitative Comparison of Natural Languages and Other Sequences

Christian Bentz November 1, 2023

Department of General Linguistics, University of Tübingen

Acknowledgements

URPP Language and Space



DFG Center for Advanced Studies "Words, Bones, Genes, Tools"







Introduction

If a Martian scientist [...] received from Earth the broadcast of an extensive speech [...] what criteria would [...] determine whether the reception represented the effect of an animate process on Earth, or merely the latest thunderstorm on Earth? It seems that the only criteria would be the arrangement of occurrences of the elements [...]: the arrangement of the occurrences would be neither of **rigidly fixed regularity** [...] nor yet a completely random scattering of the same

Zipf (1936). The psycho-biology of language, p. 187.





• Is there a measure/algorithm that can distinguish between natural languages and other sequences?

- Is there a measure/algorithm that can distinguish between natural languages and other sequences?
- If yes why? What makes languages different?

Data

Corpus of Sequences (StringBase)

Hiwis: Tim Wientzek, Clara Garcia Baumgärtner

Number of Files: 138

```
Number of UTF-8
Characters:
ca. 10 to ca. 200 000
per file
```



Writing

Universal Declaration of Human Rights

 47 translations (different writing systems)



Eleanor Roosevelt

#type:: writing #spectiation: Val (val) #sortpictode: Vali #sororce: https://unicode.org/udhr/d/udhr_val.htnl (last access: 18.09.2019) #encoding: utf-8 #copyright: NA #comments: Full translation of the Universal Declaration of Human Rights into Val. The document starts with a preamble ("베바 법각도 아/봄"). Every subtiem within an article has its own line.

line_1> |||H H℃ ⊕ ☉ \#

</time 4> ካጜ ከ ዶ K ለቆ፲ ሂ ክ ፬ ዓ.፰ 3 ጠ ቴ ዾ ፬ ዓ.ም ቀ 8 ኳ ሺ 8 8 ሂ 8 0 ፕ 6 ⊫ ከ 6 ዜ ፁ ሚ ጠ 8 ኚ ኛ • ፕ ደ ዮ ጽ ጽ ሕ ፕ ኤ ዮ ሕ ሞ. ሂ ዶ ዩ K ጴ፱ ሂ ት ዎ ሮ ሞ ዝ ኳ ት መ 8 ኚ 8 መ 8

Sumerian, Akkadian, Elamite, Prakrit, Cretan Hieroglyphs

• 14 texts



#type: writing #specification: Cuneiform from the Lagash II period of the Sumerian language (sum) #scriptcode: Latn #source: https://cdli.ucla.edu/search/archival view.php?ObjectID=P232529 (last access: 28.03.2019); Full sign list: https:// cdli.ucla.edu/tools/SignLists/Rosengarten.pdf (last access: 28.03.2019) #encodina: utf-8 "Text in the pages of CDLI may be freely copied, aggregated and re-used according to common academic practice; we request, in #copyright: the case of re-use of considerable textual data, that mention be made of the source of such material, with reference to CDLI and its URL <http://cdli.ucla.edu>." From https://cdli.ucla.edu/?g=terms-of-use (last accessed 27.02.2020). #comments: Cuneiform of the Lagash II period (ca. 2200-2100 BC). A full sign list can be found on the above-mentioned link. Lines beginning with '#' are comments. List of transliteration conventions: https://cdli.ucla.edu/methods/ednotes.html

<obverse> <line_1> {d}nansze <line_2> nin uru16 <line_3> nin in-dub-ba <line_4> nin-a-ni

Palaeolithic Signs of the Aurignacian (ca. 35 000 BP)



• 20 sequences

#type: non-writing
#specification: palaeolithic signs (pal)
#scriptcode: Latn
#source: NA
#encoding: utf-8
#copyright: NA

#comments: These are encodings of geometric signs found on mobile objects of the palaeolithic. Each line represents a separate object. In "~>" the SignBase identifier for the object is given. White spaces indicate some visual distinction between rows and clusters of signs. Alphabet of this particular transcription: "z": zoomorph; "d": dot; "V": vulva; "m": motch; "o": cupule; "N": obmotch; "r": rhombus; "h": hatching; "L": line; "Z": zigagrow; "C": circummotch;

Nonwriting

Weather Symbols

1 sequence



#type: non-wrtting
#specifiation: weather symbols (wsy)
#scriptcode: Latn
#source: http://www.weather.com (last access: 03.03.2019)
#encoding: Ut-2
#encoding: Ut-2
#encoding: Ut-2
#encoding: the data was collected by looking at the 10-day f

#comments: The data was collected by looking at the 10-day forecast (which actually is a 14-days forecast) for different cities around the world. Following transcriptions of the weather symbols have been used: MostlyCloudy, Showers, Cloudy, PartlyCloudy, RainSnowShowers, SnowShowers, Thunderstorms, MostlySunny, IsolatedThunderstorms, Sunny, Mind, ScatteredThunderstorms, Rain

showers Wind Cloudy Cloudy Showers Show

Thunderstorms Thunderstorms MostlySunny NostlySunny MostlySunny IsolatedThunderstorms MostlySunny MostlySunny Thunderstorms MostlySunny Thunderstorms Thunderstorm

Cloudy IsolatedThunderstorms Thunderstorms Mostlys MostlyCloudy MostlyCloudy PartlyCloudy IsolatedThunderstorms

The Voynich Manuscript

 1 sequence (ca. 200 000 characters)



#type: unclassified #specification: Vovnich manuscript (vov) #scriptcode: Latn #source: http://www.voynich.com/pages/PagesH.txt (last accessed 05.02.2019) #encoding: utf-8 #copyright: NA #comments: This is a transcription of the Vovnich manuscript by Takeshi Takahashi in the so-called EVA alphabet. Details on the transcription and line identifiers can be found at http://www.vovnich.nu/transcr.html#n15. "." as word separators have been replaced by blank spaces. The character "*", which is used in the original EVA transcription system for unreadable characters, is here replaced by "?". Note that "?" has the meaning of "missing word" in the original EVA transcription system, but there are actually no "?" found in the transcription we used.

<fir> <fir.P1.1;H> <fir.P1.2;H> <fir.P1.2;H> <fir.P1.3;H> sory ckhar or y kair chtaiin shar are cthar cthar dan <fir.P1.3;H> syaiir sheky or ykaiin shod cthoary cthes daraiin sa

<f1r.P1.4;H> ooiin oteey oteos roloty cth?ar daiin otaiin or okan

Animal

Birdsong

- 5 species of birds
- 33 sequences

Arriaga et al. (2015). Bird-DB: A database for annotated bird song sequences.





Fig. 1. Audio recording of black-headed grosbeak showing annotation with Prac

#type: aninal
#specification: Black-Headed Grosbeak (bhg)
#spctiptode: Latn
msource: http://taylor8.biology.ucla.edu/birdD8Query/Files/Files_TextGrids/2013/Jun/GTNP13-69.TextGrid (last access: 13.07.2019),
phrases appended in bhg_Phrases.pdf
#encoding: urf-8
#ccopyright: NA
#ccomments: Each 'unit' corresponds to a sound produced by the bird. White spaces indicate intervals of different durations between the
units. These annotations resulted from a recording made in one piece from a specific bird. Comments that would sometimes follow the bird
recording user encoved.

cline_j> uj kdi ro su sv sw sqf jr dw kd tc jt ag ta tb uj kdi no su sv sw sqf jr dw kd tc jr ag ta td te uj kdi no su sv sw sx sqf jr dw kd tc jr ag ta td te uj kdi no su sv sw sqf jr dw kd tc jt ag ta td te uj kdi no su sv sx sw sh jr dw kd tc jt ag ta td te uj kdi no su sv sw sqf jr dw kd tc jt ag ta td te uj kdi no su sv sw sk sqf jr dw kd tc jt ag ta td te uj kdi no su sv sw sqf jr dw kd tc jt ag ta td te uj kdi no su sv sw sk sqf jr dw kd tc jt ag ta td te uj kdi no su sv sw sqf jr dw kd tc jt ag ta td te uj kdi no su sv sx sw sh jr dw kd tc jt ag ta td te uj kdi no su sv sw sk sqf jr dw kd tc jt ag ta td te uj kdi no su sv sx sw sk sqf jr dw kd tc jt ag ta td te uj kdi no su sv sx sw sh jr dw kd tc jt ag ta jt ag tu jt kdi no su sv sx sw sh jr dw kd tc jt ag ta jt dw sy st sw sw sjr dw kd tc jt ag ta jt dw kd tc jt ag ta td te uj kdi no su sv sx sqf sw jr dw kd tc jt ag ta jt dt tc jt ag ge ta uj kdi no su sv sx sw sk gf jr dw kd tc jt ag ta td te uj kdi no su sv sx sqf sw jr dw kd tc jt ag ta td tc

Natural

DNA

• 30 sequences



#type:	natural
#specification:	DNA (dna)
#scriptcode:	Latn
#source:	https://www.ncbi.nlm.nih.gov/nuccore/BA0001000001.1?report=fasta (last access: 04.07.2019), Input formats found on https://
blast.ncbi.nlm.u	nih.gov/Blast.cgi?CMD=Web&PAGE TYPE=BlastDocs&DOC TYPE=BlastHelp (last access: 04.07.2019)
#encoding:	utf-8
#copyright:	NA
#comments:	Raphanus Sativus (radish) DNA, "N" stands for any nucleic acid (see input formats)
line_1>	GGTAGTTAGGGTCTGAAAAAGATTTTGCGTTTTGATAGTTAAATCGATGTAGAACTAAAGTCCCGGTGCA
<line_2></line_2>	AATGAGATTCATGCAGTTTGTAGAATAAGTGTGGGATTTGGTAAAAATAATTTGCAGTTTTGATAAAAAAA
line_3>	ATGTAAGACGATTTTGATTACAAAAATGATGAAAATATATAT
<line_4></line_4>	TTCTTAGAGATAATGTCTTAATATTTTATTTTATATATAATATTAGAATATTAGAACTAATATGAAATAACAA
<line_5></line_5>	ττςτααττατατατταττατταττατταττατταςταααατταττ
<line_6></line_6>	TGTAAATAAAGTGATAAATAAAAATAAAAATTAATGATTTCAATTTCCAATAACATATTTATAACACAACT
line_7>	GTTAAAATGCATCGTAACAATCATTTATATATATATATAT
<line 8=""></line>	TAAAATGTCAAGATTACCATTTTCTGGTTTTTTCAATGTTGATACAACTTATTAGAAACAAAC

Methods

- Sequence: All the characters in the entire file.
- **String**: A string of characters delimited by *white spaces*.
- Character: Individual UTF-8 character.
- Unit: Information encoding unit, i.e. string or character.

$$\widehat{H}^{ML}(X) = -\sum_{i=1}^{W} \widehat{p}^{ML}(x_i) \log_2 \widehat{p}^{ML}(x_i)$$
(1)

• ML: Maximum likelihood or "plug-in" estimator

Shannon, Claude E. (1948). A mathematical theory of communication. Cover & Thomas (2006). Elements of information theory, p. 14.

unit	freq	
а	5	•
А	1	
b	2	$\widehat{H}^{ML}(X) = -(\frac{5}{51}\log_2(\frac{5}{51}) + \frac{1}{51}\log_2(\frac{1}{51}) + \dots) \sim$
d	3	3.97
е	5	
f	1	

 $\rm All_1\ human_2\ beings_3\ are_4\ born_5\ free_6\ and_7\ equal_8\ in_9\ dignity_{10}\ and_{11}\ rights_{12}\ [...]$

unit	freq	
All	1	
and	2	
are	1	$\widehat{H}^{ML}(X) = -(\frac{1}{12}\log_2(\frac{1}{12}) + \frac{2}{12}\log_2(\frac{2}{12}) + \dots) \sim$
beings	1	3.41
born	1	
dignity	1	

The entropy rate as "per symbol entropy of *n* random variables".

$$h(\mathcal{X}) = \lim_{n \to \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n)$$
(2)

The entropy rate as "per symbol entropy of *n* random variables".

$$h(\mathcal{X}) = \lim_{n \to \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n)$$
(2)

Alternative definition as "the conditional entropy of the last random variable given the past":

$$h(\mathcal{X}) = \lim_{n \to \infty} H(X_n | X_1, X_2, ..., X_{n-1}).$$
(3)

Assumptions:

- $\cdot \ \mathcal{X}$ is a stationary stochastic process
- *n*: number of word tokens

Cover & Thomas (2006) Elements of information theory, p. 74-75. Dębowski (2020) Information theory meets power laws.

$$\hat{h}(X) = \frac{1}{n} \sum_{i=2}^{n} \frac{\log_2 i}{L_i},$$
(4)

- n: number of unit tokens
- L_i : length (+1) of the longest contiguous subsequence starting at position *i* which is also present in i = 2 to i 1

Gao, Kontoyiannis & Bienenstock (2008). Estimating the entropy of binary time series, equation (6).

 $L_{43} = 3(+1) = 4$ $\frac{\log_2(43)}{4} \sim \frac{5.43}{4} \sim 1.36$

R package *Hrate* (https://github.com/dimalik/Hrate) Bentz, Alikaniotis, Cysouw, & Ferrer-i-Cancho (2017). The entropy of words - learnability and expressivity across more than 1000 languages.

Example Plots (UDHR First Sentence in English)



Results

Stabilization Analyses (100 Units)

Character Entropies (100 Units)



String Entropies (100 Units)



- Extending the database
- Use further measures (repetition rate, mutual information, Zipf parameters etc.)
- Create and test classification system

Extending the Database: The 100 Language Corpus





https://www.spur.uzh.ch/en/departments/research/textgroup/MorphDiv.html



Tanja Samardžić



Olga Sozinova



Ximena Gutierrez-Vasques



Extending the Database: Palaeolithic Signs

Current



Ewa Dutkiewicz



Chris Bentz



Saetbyul Lee



Gabriele Russo





Stephanie Samson



David Matzig

SignBase



Dutkiewicz, et al. (2020) SignBase, a collection of geometric signs on mobile objects in the Paleolithic.

www.signbase.org

Home Objects References Description Contributors Download Imprint Data Protection Conta	Home	Objects References	Description	Contributors	Download	Imprint	Data Protection	Conta
---	------	--------------------	-------------	--------------	----------	---------	-----------------	-------

Welcome to SignBase!

SignBase is an open access database for geometric motifs on mobile objects in Prehistory. Its focus lies on finds of the Eurasian Paleolithic and African Middle Stone Age. In these time periods, geometric motifs – also referred to as signs, patterns, or marks – are abundant in parietal art as well as on mobile objects. The term 'geometric' denotes simple nonfigurative forms such as dots, lines, and crosses, as well as more complex patterns. This includes frequent semi-abstract depictions such as vulvae, but excludes figurative depictions of animals, humans, etc. Decorted mobile objects are mostly made of osseous material, like ivory, bone or antiev, while also featuring other organic and inorganic materials.



The relevant artifacts come from stratified archaeological contexts and are assigned to the particular techno-complexes. The objects are the core elements of the database, carrying a unique identifier. With this identifier comes information about geographic and archaeological provenience, the type of object and material, size and preservation, literature references, as well as a picture if available. The geometric motifs on each object are described in detail using a specifically developed encoding. The database aims to enable quantitative comparative studies on the development of graphical expressions before the emergence of writing systems.

SignBase gives latitudinal and longitudinal information for each object, and hence allows analyses of geographical distribution and geographic clustering of signs. Densities of sites with geometric signs, or densities of particular sign types give an overview of the abundance and distribution of signs. In the long run, this might be linked to population turn-overs, and respective cultural entities of the Paleolithic.



Thank You