# Beyond Words

Lower and upper bounds on the entropy of subword units in diverse languages
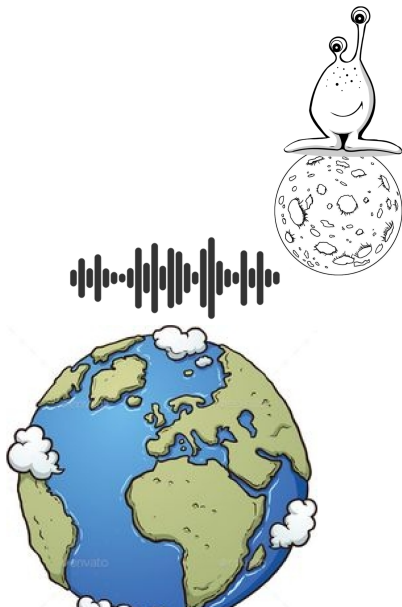
Christian Bentz

August 8, 2023

Department of General Linguistics, University of Tübingen

# Introduction

# The Marsian Scientist



*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.

| System | Example* |
|---|---|
| English | And they remembered his words , |
| Greek (Modern) | Και ενεθυμηθησαν τους λογους αυτου . |
| Hindi | तब उस की बातें उन को स्मरण आईं । |
| Georgian | და მოეჳსკენნეს სიტყუანი მისნი . |
| Korean | 저희가 예수의 말씀을 기억하고 |
| Burmese | မိန့်တော်မူခဲ့သောစကားများကိုပြန်၍သတိရ ၍ ၊ – |
| Russian | И они вспомнили эти слова Его . |
| Inuktitut | ᐊᕐᖁᓱᓚ ᐃᖅᑭᕐᖁᔭᑦ ᐅᖖᐅᕐᓯᕐᖁᓂᖅ , |
| Kalaallisut | Taava oqaasii eqqaalerpaat . |
| Bird Song | uj kd ro su sv sw sx gf jr dw kd tc jt ag ta |
| Morse | phh pppp p hp s pp hp s h pppp p s hphp hhh |
| DNA | GGTAGTTAGGGTCTGAAAAAGATTTTGCG |
| Weather | SWCCSSSSSSSSSSSCSOFSPPPFPPFPP |
| Random | hihhe bh fif cd gbgdiiigc ghigbbg af icegeebiifg |
| Shuffled | swr a j e eitimii hfeooa ti i d qs sfi roeviebg ep |

*For natural languages: Verse number 42024008 of the New Testament.

# Methods: Unigram Entropy

### Definition

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

- V: number of word types,
- $p(x_i)$: probability of word type.

### Example

$\text{in}_1 \text{ the}_2 \text{ beginning}_3 \text{ god}_4 \text{ created}_5 \text{ the}_6 \text{ heavens}_7$
$\text{and}_8 \text{ the}_9 \text{ earth}_{10} \text{ and}_{11} \text{ the}_{12} \text{ earth}_{13} \text{ was}_{14}$
$\text{waste}_{15} \text{ and}_{16} \text{ empty}_{17}$ [...]

| unit | freq |
|------|------|
| the | 4 |
| and | 3 |
| earth | 2 |
| in | 1 |
| beginning | 1 |
| god | 1 |
| ... | ... |

$\widehat{H}^{ML}(X) = -(\frac{4}{17} \log_2(\frac{4}{17}) + \frac{3}{17} \log_2(\frac{3}{17}) + \cdots + \frac{1}{17} \log_2(\frac{1}{17})) \sim 3.2$

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

## Methods: Entropy Rate

**Definition**

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

- $n$: number of word tokens,
- $L_i$: length (+1) of the longest contiguous substring starting at position $i$ which is also present in $i = 2$ to $i - 1$.
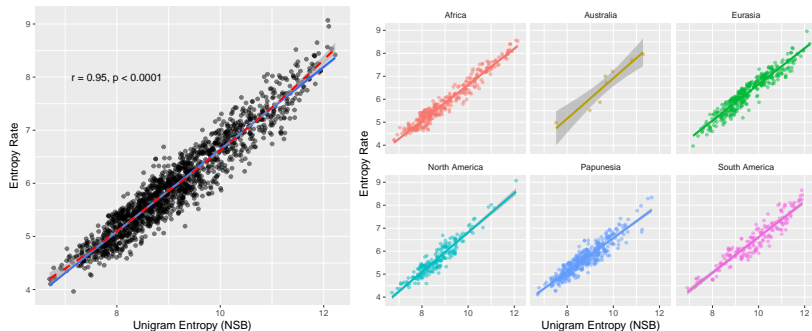
**Example**

$\text{in}_1 \ \text{the}_2 \ \text{beginning}_3 \ \text{god}_4 \ \text{created}_5 \ \text{the}_6 \ \text{heavens}_7$ $\text{and}_8 \ \text{the}_9 \ \text{earth}_{10} \ \boxed{\text{and}}_{11} \ \underline{\text{the}}_{12} \ \underline{\text{earth}}_{13} \ \text{was}_{14}$ $\text{waste}_{15} \ \text{and}_{16} \ \text{empty}_{17} \ [...]$

$$L_{11} = 3(+1) = 4$$

$$\frac{\log_2(11)}{4} \sim \frac{3.46}{4} \sim 0.87$$

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

4

# Entropy Across Languages of the World



Bentz, Alikaniotis, Cysouw, & Ferrer-i-Cancho (2017). The entropy of words - learnability and expressivity across more than 1000 languages.

Lavi-Rotbain & Arnon (2019). Children learn words better in low entropy.

Tal, Grossman, & Arnon (2022). Infant-directed speech becomes less redundant as infants grow: implications for language learning.

Lavi-Rotbain & Arnon (2023). Zipfian distributions in child-directed speech.

5

# Linguistic Interpretation

(1) Hawaiian (haw, PBC 41006018)

   *A   ua   olelo aku o   Ioane ia ia*   [...]
   Then PERF say   to   SUBJ Johan he.DAT [...]

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

(2) Turkish (tur, PBC 41006004)

   *Ýsa   da   on-lar-a*   [...] *de-di*
   Jesus also 3P-PL-DAT [...] say-3SG.PERF

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

(3) Iñupiatun (esk, PBC 41006004)

   *Aglaan Jesus-ŋum   itna-ġ-ni-ġai*   [...]
   But   Jesus-ERG this-say-report-3S.to.3PL

   "But Jesus said to them (it is reported) [...]"

Bentz (2018). Adaptive languages: An information-theoretic account of language diversity.

**Word Elephant**



What if we use other units of
measurement?

## Three Elephants in the Room

**Word Elephant**



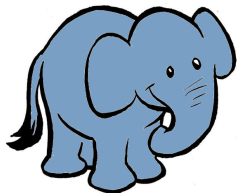What if we use other units of measurement?

**Spoken Elephant**
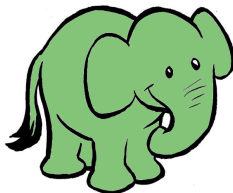


What about spoken language?

# Three Elephants in the Room
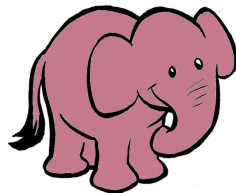
**Word Elephant**



What if we use other units of measurement?

**Spoken Elephant**



What about spoken language?

**Meaning Elephant**


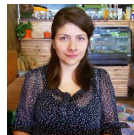
How does this relate to meaning anyways?

# Subword Units

# Subword Entropy

| tinechcaquiznequi (Nahuatl) |
| me quieres oir (Spanish) |
| *you want to hear me* |

| In cihuamizton ipan ahcopechtli ca.(Nahuatl) |
| La gata estaba encima de la mesa. (Spanish) |
| *The (female) cat is on the table* |

| pejke san motlajtlachiliyaj (Nahuatl) |
| empezaron a mirarse nada mas (Spanish) |
| *they started to just look at each others* |

Table 1: Examples of Nahuatl-Spanish parallel sentences



Ximena
Gutierrez-Vasques

| ti-nech-caqui-z-nequi |
| 2.SG.S-1.S.O-'hear'-FUT-'want' |
| "Tú me quieres oír" (Spanish) |
| *you want to hear me* |
| Lexical correspondence: oir-caqui |

Table 2: Example of Nahuatl-Spanish
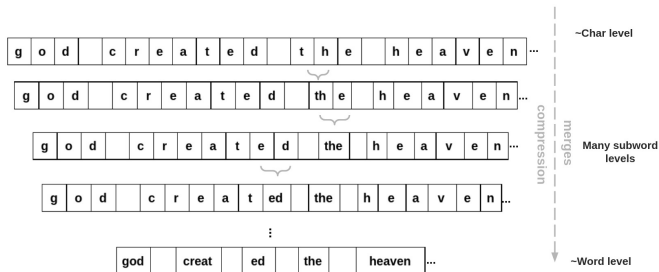
Gutierrez-Vasques, Sierra, & Pompa (2016). Axolotl: A web accessible parallel corpus for Spanish-Nahuatl.
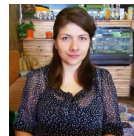
Ximena Gutierrez-Vasques & Victor Mijangos. (2018). Comparing morphological complexity of Spanish, Otomi and Nahuatl.

8

# Subwords of Byte Pair Encoding (BPE)



~Char level

compression

merges

Many subword levels

~Word level

Ximena
Gutierrez-Vasques

Olga Pelloni
(Sozinova)

Gutierrez-Vasques, Bentz, Sozinova, & Samardžić (2021). From characters to words: the turning point of BPE merges. *EACL*.

Gutierrez-Vasques, Bentz, & Samardžić (2023). Languages through the Looking Glass of BPE Compression. *Computational Linguistics*.
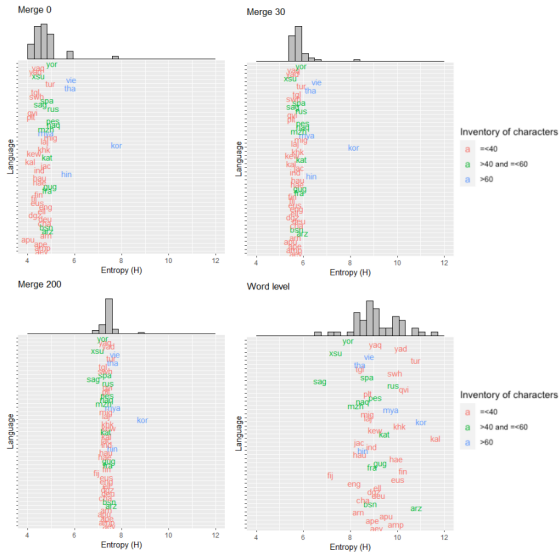
Tanja Samardžić

# From Characters to Words

# Character Entropy



Merge 0

**Inventory of characters**
- a  =<40
- a  >40 and =<60
- a  >60

Axes: Language (y-axis), Entropy (H) (x-axis)

| System | Example* |
|---|---|
| English | And they remembered his words , |
| Greek (Modern) | Και ενεθυμηθησαν τους λογους αυτου . |
| Hindi | तब उस की बातें उन को स्मरण आईं । |
| Georgian | და მოეჰსენებეს სიტყუანი მისნი . |
| Korean | 저희가 예수의 말씀을 기억하고 |
| Burmese | မိန့်တော်မူခဲ့သောစကားများကိုပြန်သတိရ ၍ �၊ - |
| Russian | И они вспомнили эти слова Его . |
| Inuktitut | ᐊᖕᒑᓛ ᐃᖁᒥᕐᓴᔅ ᐅᖃᐅᓯᖏᓐ , |
| Kalaallisut | Taava oqaasii eqqaalerpaat . |

# Subword Entropy (200 Merges)



Merge 200 — scatter plot of Language vs Entropy (H). Inventory of characters: a =<40, a >40 and =<60, a >60.

| English (eng) | | | | Turkish (tur) | | | |
|---|---|---|---|---|---|---|---|
| prod. | subword | merge | examples | prod. | subword | merge | examples |
| 110 | ing<\w> | 29 | begin**ning** | 355 | lar | 8 | on**lara** |
| 67 | eth<\w> | 103 | nazar**eth**, eat**eth** | 309 | ler | 13 | gün**ler**de |
| 38 | est<\w> | 166 | l**est**, car**est** | 150 | ler<\w> | 32 | gitti**ler** |
| 26 | led<\w> | 122 | fil**led**, cal**led** | 145 | lar<\w> | 38 | adam**lar** |
| 26 | com | 137 | **com**ing, **com**e | 131 | den<\w> | 40 | sen**den** |
| 23 | oun | 129 | r**oun**d, f**oun**d | 129 | yor | 69 | öğreti**yor**d |
| 21 | for | 86 | **for**sook | 116 | dan<\w> | 60 | tarafın**dan** |
| 21 | ent<\w> | 91 | w**ent**, garm**ent** | 110 | ini<\w> | 64 | indiğ**ini** |
| 21 | ght<\w> | 102 | tau**ght**, mi**ght** | 107 | ların | 96 | ağ**ların**ı |
| 19 | ing | 90 | th**ing**s, br**ing**ing | 80 | ine<\w> | 50 | üzer**ine** |

# Word Entropy (2000 Merges)



Merge 2000

Inventory of characters
a  =<40
a  >40 and =<60
a  >60

Sanumá (xsu)

(4) Sama        töpö se  kite
    1PL.EXCL 3PL hit FUT

    "We will hit them." Borgman (1990)

Georgian (kat)

(5) და  მოეკსენნეს
    da  mo-e-qsenn-es
    and PREV-3P.PL-mention-3P.PL.AOR
    სიტყუანი        მისნი
    sit'q'va-n-i    misni
    word-PL-NOM 3P.POSS.NOM.SG

    Literal translation: "And they mentioned his words."
    English verse: "And they remembered his words."

# Corpus-Based Typology

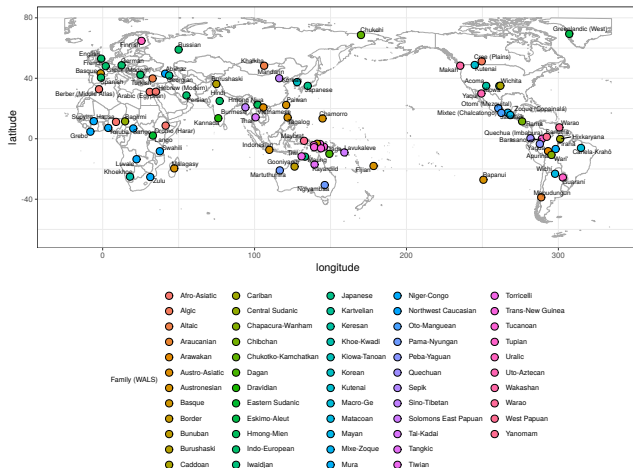| ISO | Name | Prod. | Synthesis | Reference | Fusion |
|-----|------|-------|-----------|-----------|--------|
| vie | Vietnamese | -1.33 | analytic | Haspelmath 2010 | isolating |
| tha | Thai | -1.33 | analytic | Moravcsik 2012 | isolating/concat. |
| sag | Sango | -1.29 | analytic | Karan 2006 | concatenative |
| yor | Yoruba | -1.21 | analytic | Haspelmath 2010 | tonal/isolating |
| eng | English | -0.94 | analytic | Haspelmath 2010 | concatenative |
| fij | Fijian | -0.89 | analytic | Dixon 1988 | isolating |
| pes | Persian/Farsi | -0.78 | synthetic | Greenberg 1960 | concatenative |
| fra | French | -0.19 | synthetic | Dixon 2003 | concatenative |
| ell | Greek (Modern) | 0.02 | synthetic | Dixon 2003 | concatenative |
| rus | Russian | 0.07 | synthetic | Aikhenvald 2007 | concatenative |
| swh | Swahili | 0.2 | synthetic | Haspelmath 2010 | concatenative |
| yaq | Yaqui | 0.46 | synthetic | Guerrero 2019 | concatenative |
| tur | Turkish | 0.54 | synthetic | Bickel 2007 | concatenative |
| gug | Paraguayan Guaraní | -0.19 | polysynthetic | Aikhenvald 2017 | concatenative |
| arn | Mapudungun | 0.73 | polysynthetic | Bickel 2017 | concatenative |
| amp | Alamblak | 1.02 | polysynthetic | Bruce 1984 | concatenative |
| apu | Apurinã | 1.3 | polysynthetic | Facundes 2014 | concatenative |
| bsn | Barasano | 1.43 | polysynthetic | Gomez 2004 | concatenative |
| kal | Kalaallisut | 3.25 | polysynthetic | Haspelmath 2010 | concatenative |

**Summary**

Languages have relatively high entropy divergence at the **character level** (c. 4-8 bits per character), and at the **word level** (c. 7 to 11 bits per word), but they have similar entropies at specific **subword levels** (c. 7-8 bits per subword).
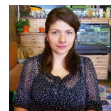
# Spoken vs. Written

# Text Data Diversity Sample (TeDDi)



Family (WALS)

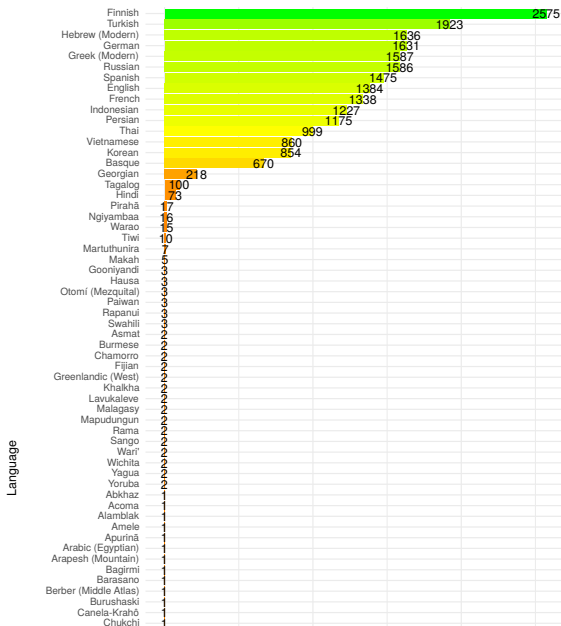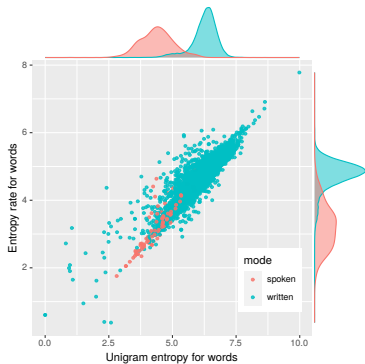| | | | | |
|---|---|---|---|---|
| Afro-Asiatic | Cariban | Japanese | Niger-Congo | Torricelli |
| Algic | Central Sudanic | Kartvelian | Northwest Caucasian | Trans-New Guinea |
| Altaic | Chapacura-Wanham | Keresan | Oto-Manguean | Tucanoan |
| Araucanian | Chibchan | Khoe-Kwadi | Pama-Nyungan | Tupian |
| Arawakan | Chukotko-Kamchatkan | Kiowa-Tanoan | Peba-Yaguan | Uralic |
| Austro-Asiatic | Dagan | Korean | Quechuan | Uto-Aztecan |
| Austronesian | Dravidian | Kutenai | Sepik | Wakashan |
| Basque | Eastern Sudanic | Macro-Ge | Sino-Tibetan | Warao |
| Border | Eskimo-Aleut | Mataocan | Solomons East Papuan | West Papuan |
| Bununan | Hmong-Mien | Mayan | Tai-Kadai | Yanomam |
| Burushaski | Indo-European | Mixe-Zoque | Tangkic | |
| Caddoan | Iwaidjan | Mura | Tiwian | |

Moran, Bentz, Gutierrez-Vasques, Sozinova, & Samardžić (2022). TeDDi
Sample: Text Data Diversity Sample for Language Comparison and Multilingual
NLP.


Tanja Samardžić
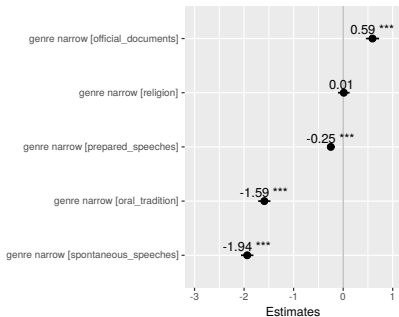

Olga Pelloni (Sozinova)


Ximena Gutierrez-Vasques


Steven Moran

16

# TeDDi Overview



| Language | Count |
|---|---|
| Finnish | 2575 |
| Turkish | 1923 |
| Hebrew (Modern) | 1636 |
| German | 1631 |
| Greek (Modern) | 1587 |
| Russian | 1586 |
| Spanish | 1475 |
| English | 1384 |
| French | 1338 |
| Indonesian | 1227 |
| Persian | 1175 |
| Thai | 999 |
| Vietnamese | 860 |
| Korean | 854 |
| Basque | 670 |
| Georgian | 218 |
| Tagalog | 100 |
| Hindi | 73 |
| Pirahã | 17 |
| Ngiyambaa | 16 |
| Warao | 15 |
| Tiwi | 10 |
| Martuthunira | 7 |
| Makah | 5 |
| Gooniyandi | 3 |
| Hausa | 3 |
| Otomí (Mezquital) | 3 |
| Paiwan | 3 |
| Rapanui | 3 |
| Swahili | 3 |
| Asmat | 2 |
| Burmese | 2 |
| Chamorro | 2 |
| Fijian | 2 |
| Greenlandic (West) | 2 |
| Khalkha | 2 |
| Lavukaleve | 2 |
| Malagasy | 2 |
| Mapudungun | 2 |
| Rama | 2 |
| Sango | 2 |
| Wari' | 2 |
| Wichita | 2 |
| Yagua | 2 |
| Yoruba | 2 |
| Abkhaz | 1 |
| Acoma | 1 |
| Alamblak | 1 |
| Amele | 1 |
| Apurinã | 1 |
| Arabic (Egyptian) | 1 |
| Arapesh (Mountain) | 1 |
| Bagirmi | 1 |
| Barasano | 1 |
| Berber (Middle Atlas) | 1 |
| Burushaski | 1 |
| Canela-Krahô | 1 |
| Chukchi | 1 |

# Entropy for Spoken vs. Written

# Multiple Regression Model



Word Entropy Rate

Unigram Word Entropy

# Examples

Martuthunira (vma, spoken, "Mourning chant")

```
<line_1>        Ngunhu waruul wilangayi Purripurring|ura waruul wilangayi , ngunhaa
<segmentation>  Ngunhu waruul wilangayi Purripurri-ngura waruul wilangayi , ngunhaa
<glossing>      that.NOM still HES Purripurri-BELONG still HES that.NOM
<translation>   That fellow , who is one of Purripurri's mob , he came to me

<line_2>        waruul junarrilha nganaju wilangayi . Ngayu nhawulhanguru
<segmentation>  waruul juna-rri-lha nganaju wilangayi . Ngayu nhawu-lha-nguru
<glossing>      still spirit-INV-PAST 1SG.ACC HES 1SG.NOM see-PAST-PRES
<translation>   as a spirit . I saw

<line_3>        ngurnaa mangkarnkuwilangayi . Malyarranpalharru wilangayi . Ngunhu
<segmentation>  ngurnaa mangkarn-kuwilangayi . Malyarra-npa-lha-rru wilangayi . Ngunhu
<glossing>      that.ACC spirit-ACC HES sick-INCH-PAST-NOW HES that.NOM
<translation>   his ghost . And now I've gotten sick . He
```

Dench, A. C. (1994). Martuthunira: A language of the Pilbara region of Western Australia, p. 282-287.

# Examples

### Rama (rma, spoken, "Manatee hunting")

```
<line_1>        ipang ika kiikna paalpa baanalpi traali lakun aik .
<glossing>      island of men manatee they-look-for go.out lagoon in
<translation>   ' Men of Rama Cay go manatee hunting in the lagoon '

<line_2>        paalpa ansungka , paalpa ankungi .
<glossing>      manatee they-see-when manatee they-strike
<translation>   ' When they see a manatee , they strike it '

<line_3>        paalpa anmalngu .
<glossing>      manatee they-kill
<translation>   ' They kill the manatee '
```

Craig, C. (1986). The Rama language; a text with grammatical notes. Estudios de Lingüística Chibcha 5. 21-44.

Texts of spoken registers appear to have lower entropies on average than texts of written registers. This is likely related to more repetitive usage of words and word chunks (as a result of online memory constraints?).

# Meaning

# Information and Meaning



*The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point. [...] **semantic aspects of communication are irrelevant to the engineering problem**. The significant aspect is that the actual message is one selected from a set of possible messages.*

Shannon, Claude E. (1948). A mathematical theory of communication, p. 1.

## Example

```
Article 1
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.
```

Universal Declaration of Human Rights (UDHR) in English

```
Raeiclt 1
Rll humrn btings rat boan fatt and tqurl in digniey
rnd aighes. Ehty rat tndowtd wieh atrson rnd
conscitnct rnd should rce eowrads ont rnoehta in r
spiaie of baoehtahood.
```

Universal Declaration of Human Rights (UDHR) in ???

# Information and Meaning





*[...] two messages, one of which is heavily loaded with meaning and the other which is pure nonsense, can be exactly equivalent, from the present viewpoint, as regards information. It is this, undoubtedly, that Shannon means when he says that "the semantic aspects of communication are irrelevant to the engineering aspects."* **But this does not mean that the engineering aspects are necessarily irrelevant to the semantic aspects**.
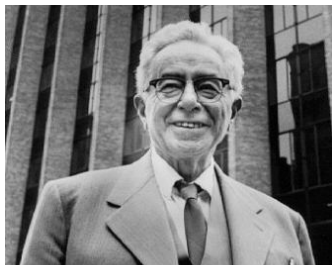
Shannon & Weaver (1949). The mathematical theory of communication, p. 8.

## Three Levels of Communication Problems



- **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). The mathematical theory of communication, p. 4.

# Entropy and Mutual Information

The entropy of signals is an upper bound on the mutual information between signals and meanings, i.e.

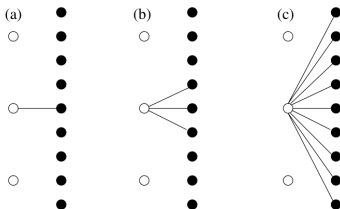$$H(S) \geq I(S, R) \tag{1}$$



**Figure 1.** Some mappings between signals (white circles) and stimuli (black circles) that are minima of $H(S)$ and $H(S|R)$ with $n = 3$ signals and $m = 9$ stimuli. (a)–(c) are minima of model A while (c) is the only valid minima of model B.

Ferrer-i-Cancho & Diaz-Guilera (2007). The global minima of the communicative energy of natural communication systems.
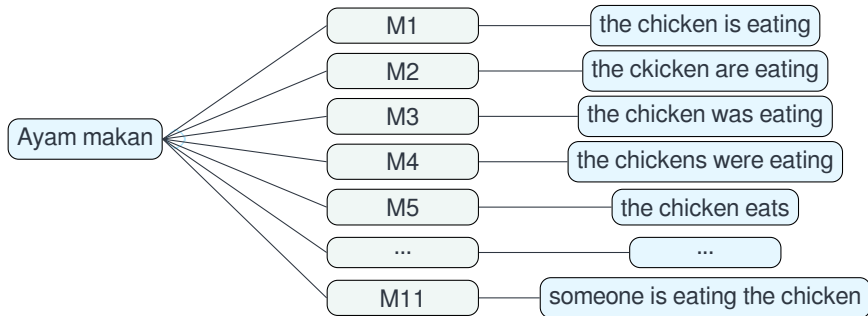
# Riau Indonesian and English

|  | English | Riau Indonesian |
|---|---|---|
|  | *The chicken is eating* | *Makan ayam/Ayam makan* |
| *Symmetry* | asymmetric: <br> agreement: *The chicken* → *is* <br> government: *is* → *-ing* | symmetric |
| *Number* (on CHICKEN) | marked: singular | unmarked: also means <br> 'The chickens are eating' |
| *Definiteness* (on CHICKEN) | marked: definite | unmarked: also means <br> 'A chicken is eating' |
| *Tense* (on EAT) | marked: present | unmarked: also means <br> 'The chicken was eating' <br> 'The chicken will be eating' |
| *Aspect* (on EAT) | marked: progressive | unmarked: also means <br> 'The chicken eats' <br> 'The chicken has eaten' |
| *Thematic role* (on CHICKEN) | marked: agent | unmarked: also means <br> 'Someone is eating the chicken' <br> 'Someone is eating for the chicken' <br> 'Someone is eating with the chicken' |
| *Ontological type* <br> (on CHICKEN EAT) | marked: activity | unmarked: also means <br> 'The chicken that is eating' <br> 'Where the chicken is eating' <br> 'When the chicken is eating' |

Gil (2005). Isolating-Monocategorial-Associational language.
Gil (2008). How much grammar does it take to sail a boat?

## Form-Meaning Mapping



Entropy of the **Signals**:

$H(S_{ind}) = -\log_2(1) = 0$ bits/chunk
$H(S_{eng}) = -\log_2(11) \sim 3.46$ bits/chunk

Conditional Entropy of **Meanings given the Signals**:

$H(M|S_{ind}) = -\sum_{s \in S_{ind}} p(s) \sum_{m \in \mathcal{M}} p(m|s) \log_2 p(m|s) = 3.46$ bits/chunk
$H(M|S_{eng}) = -\sum_{s \in S_{eng}} p(s) \sum_{m \in \mathcal{M}} p(m|s) \log_2 p(m|s) = 0$ bits/chunk

- The entropy is a **necessary but not sufficient condition** for meaningful communication.
- Languages certainly differ in entropy on the **signal side**.
- However, entropy in encoding and decoding messages is symmetrical by definition. If this holds true, then the **overall entropy is necessarily the same across languages**.

01010100 01101000 01100001
01101110 01101011 00100000
01111001 01101111 01110101

Thank You