An Informative Exploration of the Lexicon

Tiago Pimentel





Four parts:

- Phonotactic Complexity and Its Trade-offs
- Disambiguatory Signals are Stronger in Word-initial Positions
- (Non-)Arbitrariness of the Sign
 - Meaning to Form: Measuring Systematicity as Information
 - Finding Concept-specific Biases in Form--Meaning Associations
- Speakers Fill Lexical Semantic Gaps with Context

1. We use **bits per phoneme** as a measure of phonotactic complexity, and find it has **robust trade-offs** with **word length**.



2. We analyse **bits per phoneme** in first and second halves of words, and find a cross-linguistic **tendency** to **frontload information**.

Language



3. We operationalise **systematicity of the sign** as a **mutual information between word forms and meanings**, and find small (but significant) values both **within and across languages**.

Language	H(W)	MI(W;V)	U(W V)
English	3.401	0.11	3.24%
German	3.195	0.168	5.26%
Dutch	3.245	0.156	4.82%
Within Langu	lages —		



4. We operationalise **lexical ambiguity** as the **conditional entropy of a meaning given a word**, and find **consistent trade-offs** with the word's contextual entropy.



Four parts:

- Phonotactic Complexity and Its Trade-offs
- Disambiguatory Signals are Stronger in Word-initial Positions
- (Non-)Arbitrariness of the Sign
 - Meaning to Form: Measuring Systematicity as Information
 - Finding Concept-specific Biases in Form--Meaning Associations
- Speakers Fill Lexical Semantic Gaps with Context

Phonotactic Complexity and Its Trade-offs

Tiago Pimentel, Brian Roark, and Ryan Cotterell







Sequences of speech sounds allowed in a language.

For instance, in English:

- brick
- blick
- bnick

Several linguists believe all languages are equally complex.

Consequence: compensatory relationships between complexity measures should exist.

• e.g. vowel inventory size in a language correlates to a language's average word length.

Several measures of a language's phonological complexity exist:

- Size of Phoneme Inventory:
 - the number of vowel categories in a language;
- Markedness in Phoneme Inventory:
 - marked phonemes, such as clicks, make a language more complex;
- Number of Licit Syllables:
 - phonological constraints extend beyond individual units, so counting syllables seems like a logical next step in measuring its complexity;
- Word Length:
 - implicitly taken as a complexity measure when researchers examine its correlation with e.g. inventory size.

Phonotactic complexity as **bits per phoneme**

Bits per phoneme is the **entropy** of a language's word types.

We estimate it using a character-level language model's cross-entropy:

$$H(p_{lex}) \leq H(p_{lex}, q_{lex}) \approx -\frac{1}{N} \sum_{i=1}^{N} \log q_{lex}(\tilde{\mathbf{x}}^{(i)})$$

Bits per phoneme is the **entropy** of a language's word types.

We estimate it using a character-level language model's cross-entropy:

$$H(p_{lex}) \leq H(p_{lex}, q_{lex}) \approx -\frac{1}{N} \sum_{i=1}^{N} \log q_{lex}(\tilde{\mathbf{x}}^{(i)})$$

Tighter, the better the language model.

- Trigram language model;
- Character-level LSTM language model.

Linguistic rationale for bits-per-phoneme:

- Modest linguistic annotations
- Incorporates frequency of phenomena
- Captures interaction between phonemes
- Long-distance dependencies

Data: NorthEuraLex (Dellert and Jäger, 2017)

Concept-aligned lexicons

Composed of "basic" concepts

Concept	Language	Word	IPA
eye	portuguese	olho	/oʎu/
ear	finnish	korva	/kɔrυα/
give	north karelian	antua	/antʊa/
tooth	veps	hambaz	/hambaz/
black	northern sami	čáhppes	/t͡ʃaahppes/
immediately	hill mari	töpök	/tørøk/

Languages from 21 families in Europe/Asia

We omit Mandarin (no tone annotation)

Trade-Off



	Correlation		
Measure	Pearson r	Spearman ρ	
Number of:			
phonemes	-0.047	-0.054	
vowels	-0.164	-0.162	
consonants	0.030	0.045	
Bits/phoneme:			
unigram	-0.217	-0.222	
trigram	-0.682	-0.672	
LSTM	-0.762	-0.744	

Trade-Off



Bits-per-phoneme versus average word length in IPA.

Possible confound: positional effects.

• Phonemes later in a word in general have higher probability given the previous phonemes than those earlier in the string (van Son and Pols, 2003).

Truncated Words: Only consider first 3 characters in wordforms.

- Original (Full words): ρ=-0.744;
- Control (Truncated): ρ=-0.469

Correlation

Inter- and Intra- Family Trade-Offs

Classic measures of phonological complexity:

- correlate with word length across a varied set of languages,
- but do not within language families.

Bits per phoneme correlates in both cases.

	Correlation			
Measure	Pearson r	· Spea	arman ρ	
Number of:				
phonemes	-0.214	(-	0.095	
vowels	-0.383	-	0.367	
consonants	-0.147		0.092	
Bits/phoneme:				
		ſ		
	Snoon			
	Spear	man ρ		
Family	Spear LSTM	man ρ Vowels	# Langs	
Family Dravidian	Spear LSTM	man ρ Vowels -0.894	# Langs	
Family Dravidian Indo-European	Spear LSTM	man ρ Vowels -0.894 -0.218	# Langs	
Family Dravidian Indo-European Nakh-Daghestanian	Spear LSTM	man ρ Vowels -0.894 -0.218 -0.530	# Langs 4 37 6	
Family Dravidian Indo-European Nakh-Daghestanian Turkic	Spear LSTM	man ρ Vowels -0.894 -0.218 -0.530 -0.773	# Langs 4 37 6 8	

* Statistically significant with p < 0.01† Statistically significant with p < 0.1

An Informative Exploration of the Lexicon

Tiago Pimentel





Four parts:

- Phonotactic Complexity and Its Trade-offs
- Disambiguatory Signals are Stronger in Word-initial Positions
- (Non-)Arbitrariness of the Sign
 - Meaning to Form: Measuring Systematicity as Information
 - Finding Concept-specific Biases in Form--Meaning Associations
- Speakers Fill Lexical Semantic Gaps with Context

Disambiguatory Signals are Stronger in Word-initial Positions

Tiago Pimentel, Ryan Cotterell, Brian Roark







Research Question

Are word-initial segments more informative for disambiguation than word-final ones?

Introduction

Is it easier to guess the ending of "dino****"?



Introduction

Is it easier to guess the ending of "dino****"?

Introduction

Is it easier to guess the ending of "dino****"?

Or the prefix of "****saur"?



Psycholinguistic experiments with human subjects:

Psycholinguistic experiments with human subjects:

- listeners find word-initial consonant deletions more disruptive than word-final (Bagley, 1900)
- mispronunciations are more likely in word endings (Fay and Cutler, 1977)
- recognizing written words with flipped initial characters is harder than with final ones (Bruner and O'Dowd, 1958)

Psycholinguistic experiments with human

- This does not measure how informative segments are. Only how useful they are for humans. listeners find word-initial consonant del than word-final (Bagley, 1900)
- mispronunciations are p Cutler, 1977)
- recognizip

thar

Intial characters is harder

cts:

Information-theoretic measurements on natural corpora:

• Estimate a segment's information as its contextual entropy:

$$\mathrm{H}(W_i \mid W_{< i}) = -\sum\limits_{\mathbf{w} \in \Sigma^*} p(w_i \mid \mathbf{w}_{< i}) \log p(w_i \mid \mathbf{w}_{< i})$$

Information-theoretic measurements on natural corpora:

• Estimate a segme formation as its contextual entropy:

$$\mathrm{H}(W_i \mid W_{< i}) = -\sum\limits_{\mathbf{w} \in \Sigma^*} p(w_i \mid \mathbf{w}_{< i}) \log p(w_i \mid \mathbf{w}_{< i})$$

Information-theoretic measurements on natural corpora: Character

• Estimate a segme entropy:

H
$$(W_i \mid W_{< i}) = -\sum\limits_{\mathbf{w} \in \Sigma^*} p(w_i \mid \mathbf{w}_{< i}) \log p(w_i \mid \mathbf{w}_{< i})$$

mation

Ctera

htextual

Information-theoretic measurements on natural corpora:

• Estimate a segment's information as its contextual entropy:

$$\mathrm{H}(W_i \mid W_{< i}) = -\sum\limits_{\mathbf{w} \in \Sigma^*} p(w_i \mid \mathbf{w}_{< i}) \log p(w_i \mid \mathbf{w}_{< i})$$

Information-theoretic measurements on natural corpora:

• Estimate a segment's information as its contextual entropy:

$$\mathrm{H}(V_i | W_{< i}) = -\sum\limits_{\mathbf{w} \in \Sigma^*} p(w_i \mid \mathbf{w}_{< i}) \log p(w_i \mid \mathbf{w}_{< i})$$

This inherently confounds **context size** and **word position**.

Information-theoretic measurements on natural corpora:

• Estimate a segment's information as its contextual entropy:

$$\mathrm{H}(V_i | W_{< i}) = -\sum\limits_{\mathbf{w} \in \Sigma^*} p(w_i \mid \mathbf{w}_{< i}) \log p(w_i \mid \mathbf{w}_{< i})$$

This inherently confounds **context size** and **word position**.
Left-to-right Conditional Entropy

Conditioning information can only reduce entropy! $\mathrm{H}(W_t \mid W_{< t}) \leq \mathrm{H}(W_t \mid W_{t-1}) \leq \mathrm{H}(W_t)$



Left-to-right Conditional Entropy

Conditioning information can only reduce entropy! $\mathrm{H}(W_t \mid W_{< t}) \leq \mathrm{H}(W_t \mid W_{t-1}) \leq \mathrm{H}(W_t)$

Left-to-right conditional entropies, thus:

 confound the amount of conditional information with word position.



<u>Why not Left-to-right?</u>

Consider an artificial language where every word contains a copy of its first half:

- e.g., foofoo, barbar, foobarfoobar, etc.
- initial and final halves have identical disambiguatory strength; they are the same!
- conditional surprisal would be nearly zero for final halves.

Does **left-to-right conditional entropy** measure a **property of the lexicon** or simply the fact that **conditioning reduces entropy**?

Results - Forward Surprisal

All but one language in the three analysed datasets had larger word-initial surprisal Significant word-initial Significant word-final Forward Languages Dataset CELEX 3 3 0 NorthEuraLex 107 106 0 Wikipedia 41 41 0

Results - Forward Surprisal

All but one language in the three analysed datasets had larger word-initial surprisal

l (bits)	• w • no • ce	ikipedia ortheurale elex	x	a a a a a a a a a a a a a a a a a a a
al Surprisal				•••
Fin 5	1 martin	20 ¹⁰ 2010 2010	- 19 X	
	2	3	4	5
	Ini	tial Surp	risal (b	oits)

Dataset	Languages	Forward			
CELEX	3	3	0		
NorthEuraLex	107	106	0		
Wikipedia	41	41	0		

Results - Backward Surprisal

Many languages have higher word-final surprisals.

risal (bits) 5 5	• • • • • • • • • • • • • • • • • • •	•	
Final Surp		 wikip north celex 	edia euralex
	3	.4	5
	Initial Su	rprisal (bi	its)

Dataset	Languages	Forw	/ard	Backward		
CELEX	3	3	0	0	3	
NorthEuraLex	107	106	0	11	31	
Wikipedia	41	41	0	0	39	

There seems to be both:

- a large effect of the amount of conditional information
- a lexical effect of front-loading disambiguatory signals





We propose the use of three measures to control for context size:

We propose the use of three measures to control for context size:

• Unigram surprisal: $H(W_t)$



We propose the use of three measures to control for context size:

• Unigram surprisal: $H(W_t)$



Controlling for Full information a segment conveys We propose the use of the control for context size: Unigram surprisal: $H(W_t)$ 0 \bigcirc Without knowing S anything else, how relevant is s?

We propose the use of three measures to control for context size:

- Unigram surprisal: $H(W_t)$
- Cloze Surprisal: $H(W_t | W_{\neq t})$



We propose the use of three measures to control for context size:

- Unigram surprisal: $H(W_t)$
- Cloze Surprisal: $H(W_t | W_{\neq t})$



Controlling for C We propose thon-redundant Non-redundant xt Size nomation a segment

convevs for conte

measures to control

- Unigram arprisal: (W_t)
- Cloze Surprisal: $H(W_t | W_{\neq t})$



We propose the use of three measures to control for context size:

- Unigram surprisal: $H(W_t)$
- Cloze Surprisal: $H(W_t | W_{\neq t})$
- Position-specific Surprisal: $H(W_t | T = t, |W|)$



We propose the use of three measures to control for context size:

- Unigram surprisal: $H(W_t)$
- Cloze Surprisal: $H(W_t | W_{\neq t})$
- Position-specific Surprisal: $H(W_t | T = t, |W|)$



<u>Controlling</u> for Context Size

Information when primed We propose the use of three me by position and length for context size:

- Unigram surprisal: $H(W_t)$
- Cloze Surprisal: $H(W_t | W_{\neq t})$
- Position-specific Surprisal: $H(W_t \mid T = t, |W|)$



We propose the use of three me Informati' control for context size:

Sze Surprise
 Position-spe Inspired by Nooteboom and van der Vlugt's (1988) experiments.

Knowing the position, how relevant is s?

primed'

9th



CELEX (Baayen et al., 2015):

NorthEuraLex(Dellert et al., 2019)

Wikipedia

<u>Data</u>

CELEX (Baayen et al., 2015):

• English, German and Dutch;

NorthEuraLex(Dellert et al., 2019)

• 107 languages from 21 language families;

Wikipedia

• 41 typologically diverse languages;

<u>Data</u>

CELEX (Baayen et al., 2015):

- English, German and Dutch;
- monomorphemic words.

NorthEuraLex(Dellert et al., 2019)

- 107 languages from 21 language families;
- concept aligned word lists for these languages.

Wikipedia

- 41 typologically diverse languages;
- no phonetic information (only graphemes)

- a cross-linguistic tendency to front-load disambiguatory information
- not a universal phenomena—some languages have more informative word-final segments

Dataset	Languages	Forward		Backward		Unigram		Position-Specific		Cloze	
CELEX	3	3	0	0	3	2	0	2	1	2	1
NorthEuraLex	107	106	0	11	31	71	1	24	4	45	1
Wikipedia	41	41	0	0	39	39	1	31	1	35	2



- a cross-linguistic tendency to front-load disambiguatory information
- not a universal phenomena—some languages have more informative word-final segments

Dataset	Languages	Forward		Backward		Unigram		Position-Specific		Cloze	
CELEX	3	3	0	0	3	2	0	2	1	2	1
NorthEuraLex	107	106	0	11	31	71	1	24	4	45	1
Wikipedia	41	41	0	0	39	39	1	31	1	35	2



- a cross-linguistic tendency to front-load disambiguatory information
- not a universal phenomena—some languages have more informative word-final segments

Dataset	Languages	Forward		Backward		Unigram		Position-Specific		Cloze	
CELEX	3	3	0	0	3	2	0	2	1	2	1
NorthEuraLex	107	106	0	11	31	71	1	24	4	45	1
Wikipedia	41	41	0	0	39	39	1	31	1	35	2



- a cross-linguistic tendency to front-load disambiguatory information
- not a universal phenomena—some languages have more informative word-final segments

Dataset	Languages	Forward		Backward		Unigram		Position-Specific		Cloze	
CELEX	3	3	0	0	3	2	0	2	1	2	1
NorthEuraLex	107	106	0	11	31	71	1	24	4	45	1
Wikipedia	41	41	0	0	39	39	1	31	1	35	2







An Informative Exploration of the Lexicon

Tiago Pimentel





Four parts:

- Phonotactic Complexity and Its Trade-offs
- Disambiguatory Signals are Stronger in Word-initial Positions
- (Non-)Arbitrariness of the Sign
 - Meaning to Form: Measuring Systematicity as Information
 - Finding Concept-specific Biases in Form--Meaning Associations
- Speakers Fill Lexical Semantic Gaps with Context

(Non-)Arbitrariness of the Sign

+ Arya D. McCarthy; Brian Roark; Søren Wichmann; Damián Blasi; Ryan Cotterell

Introduction



Introduction

There are small but systematic patterns in these connections:

- Iconicity: Word forms that "resemble" their meanings, e.g. *meow*
- Systematicity of the sign: Similar meanings are more likely to have similar forms.
- Phonesthemes: Sub-morphemic units which are associated with some small semantic domain.



Meaning to Form: Measuring Systematicity as Information

Tiago Pimentel, Arya D. McCarthy, Damián Blasi, Brian Roark, Ryan Cotterell



Research Question

Can we quantify a language's systematicity of the sign?


Pearson correlation between word-pair distances:

- Phonological distance: raw word form edit distance.
- Semantic distance: word2vec cosine distance.

Problems:

- Hand defined distance metrics;
- Only linear relations between distances;
- No control for other factors (e.g. part-of-speech)





Advantages:

- No need to define distance metrics;
- Capture non-linear interactions;
- Straightforward to control for other factors;

MI(meanings; forms | POS) = H(forms | POS) - H(forms | meanings, POS)

But, how can we measure H(forms) and H(forms | meanings)?



We use two LSTMs to get the language's entropy

1. H(forms):

- Predict phone given previous ones;
- $p_{\theta}(form) = \prod p_{\theta}(w_t | w_{t-1})$
- $H(forms) \le H_{\theta}(forms) \approx -\sum \log p_{\theta}(form) / N$



We use two LSTMs to get the language's entropy

- 2. H(forms | meanings)
- Condition LSTM on meaning (word2vec embedding);
- $p_{\theta}(\text{form} | \text{meaning}) = \prod p_{\theta}(\mathbf{w}_{t} | \mathbf{w}_{t-1}, \mathbf{m})$
- $H_{\theta}(\text{forms} | \text{meanings}) \approx -\sum \log p_{\theta}(\text{form} | \text{meaning}) / N$



We now estimate the MI with the cross-entropies:

 $MI(meanings; forms) \approx H_{\theta}(forms) - H_{\theta}(forms | meanings)$

We also compute the uncertainty coefficient:

Unc(forms | meanings) = MI(meanings; forms) / H(forms)



Used only monomorphemic words.

Results:

- Statistically significant systematicity in all three languages.
- Systematicity effect is reduced when we condition on POS.

		System	aticity	Systematicity given POS tags			
Language	H(forms)	MI	Unc	MI	Unc		
English	3.401	0.11	3.24%	0.084	2.50%		
German	3.195	0.168	5.26%	0.154	4.84%		
Dutch	3.245	0.156	4.82%	0.089	2.84%		

<u>Results - NorthEuraLex</u>

Lexicon consists of "basic" concepts;

• We assume words are not multi-morphemic.



Use word2vec trained in English for all languages;

• Hard to train vectors for some languages.

<u>Results - NorthEuraLex</u>

Lexicon consists of "basic" concepts;

• We assume words are not multi-morphemic.



Use word2vec trained in English for all languages;

• Hard to train vectors for some languages.

Results:

- Significant systematicity in 87 of 106 languages;
- When we condition on POS tags, only 17 are statistically significant;
- Important to consider grammatical class on analysis.

Phonesthemes

Submorphemic affixal units

Usually flag a relatively small semantic domain

Classic example (Bergen, 2004):

- gl-
- related to light or vision;
- glimmer, glisten, glitter, gleam, glow and glint.

Should have higher mutual information values when compared to other k-grams.

Results - Phonesthemes

Results:

- We can find lists of known phonesthemes:
- all but two of our English phonesthemes are attested in prior work.
- Also find affixes which are pieces of fossilized morphology.

Language	Phonestheme	Examples
Dutch	/sx/- -/əl/ -/xt/ -/ɔp/	schelp, schild, schot, shacht, schaar kegel, nevel, beitel, vleugel, zetel beicht, nacht, vocht, plicht, licht stop, shop, drop, top, bob

Language	Phonestheme	Examples
English	/m/- /sl/-	infidel, intellect, institute, enigma, interim slop, slough, sluice, slim, slush
	-/kt/ -/mə/	aspect, object, fact, viaduct, tact panorama, asthma, trachoma, eczema, magma

Language	Phonestheme	Examples
German	/gə/- -/əln/ -/ln/ -/ən/	geschehen, Gebiet, gering, Geruecht, gesinnt rascheln, rumpeln, tummeln, torkeln, mogeln rascheln, rumpeln, tummeln, torkeln, mogeln goennen, saeen, besuchen, giessen, streiten



Finding Concept-specific Biases in Form–Meaning Associations

Tiago Pimentel, Brian Roark, Søren Wichmann, Ryan Cotterell, Damián Blasi



Research Question

Are there cross-linguistic associations between the forms and meanings of words? And how do we find them?





Data - ASJP

- Basic vocabulary wordlists
- Almost ³/₄ of world's languages (5189)!
- 100 basic concepts
 - body parts, colour terms, lower numerals, general properties (big, round), and some common flora and fauna (e.g. trees and dogs)







To maximize independence, we split our data per macro-area.

- 2 areas for training, 1 development, 1 test;
- Cross-validation.





Some language families cross macro-areas:

• Group them in the macro-area with more of the family's languages



Results



Results - Per concept

• Out of 100 concepts, 26 have significantly positive MI;



Results - Per language

- Out of 5189 languages, 85 have significantly positive MI;
 - At most 100 data points per language, hard statistical test after the corrections for multiple testing.
 - Cross-linguistic form-meaning biases are *potentially* not as rare or weak as believed. We can get significant language-level results with at most 100 concepts.
 - Further studies needed for stronger conclusions.



<u>Results - Per concept-token pair</u>

- Concept-token pairs with
 - particularly large MI across all four macro-areas.
 - *#* is the end-of-string token.
 - associations between [I] and "tongue" and between [p] and "full" (Blasi et al., 2016)
 - associations between [m] and [u] and "breast" (Jakobson, 1960; Traunmüller, 1994)
 - pronouns—e.g. I, we, you—and end-of-string [#]

Concept	Tokens	Concept	Tokens	Concept	Tokens	Concept	Tokens	Concept	Tokens
blood	S	eye	i	liver	klrt	path	d t	tree	#
bone	s u	fire	# t	louse	m n	say	#	two	r
breast	m u	fish	a s	mountain	bdglor	see	# e	water	#
come	# e	full	lopt	name	# i	skin	klprt	we	#ein
die	t	give	#	neck	0	star	klorstuw	who	#
dog	k	horn	k r	new	a	stone	k t	you	∦ain
drink	# u	Ι	# a n	night	Ndilmprtu	sun	е		
ear	elt	knee	Nbgkmortu	nose	Niu	tongue	delnr		
eat	#	leaf	alpt	one	k t	tooth	ei		

An Informative Exploration of the Lexicon

Tiago Pimentel







Four parts:

- Phonotactic Complexity and Its Trade-offs
- Disambiguatory Signals are Stronger in Word-initial Positions
- (Non-)Arbitrariness of the Sign
 - Meaning to Form: Measuring Systematicity as Information
 - Finding Concept-specific Biases in Form--Meaning Associations
- Speakers Fill Lexical Semantic Gaps with Context

Speakers Fill Lexical Semantic Gaps with Context

Tiago Pimentel, Rowan Hall Maudslay, Damián Blasi, Ryan Cotterell







Lexical Ambiguity

Words can mean more than one thing 😯

Consider the English word *buffalo*:

- You can pet a large buffalo (animal);
- You can visit Buffalo (US city);
- You can buffalo (intimidate) a person;







Lexical Ambiguity

Words can mean more than one thing 😯

Consider the English word **buffalo**:

- You can pet a large buffalo (animal);
- You can visit Buffalo (US city);
- You can buffalo (intimidate) a person;

Buffalo buffalo buffalo Buffalo buffalo!

• Paraphrased as NY bisons intimidate other NY bisons







The Good Linguistic Question

Do speakers compensate for **lexical ambiguity** by making words more predictable (i.e. less uncertain) given their context in order to accomodate the listeners?



The Good Linguistic Question

Do speakers compensate for **lexical ambiguity** by making words more predictable (i.e. less uncertain) given their context in order to accomodate the listeners?

→ Put differently: Is there a negative correlation between contextual uncertainty and lexical ambiguity?

The Good Linguistic Question

Do speakers compensate for **lexical ambiguity** by making words more predictable (i.e. less uncertain) given their context in order to accomodate the listeners?

→ Put differently: Is there a negative not test for causality contextual uncertainty and the note: We do not test for causality side note: We do only correlation only correlation in any form.

Our Operationalisations

A Measure of Lexical Ambiguity

We operationalise **lexical ambiguity** as the half-pointwise entropy:





A Measure of Lexical Ambiguity

We operationalise le al ambiguity as the half-pointwise leaning by: H(M | W=w)



<u>A Measure of Lexical Ambiguity</u>



Equivalent, up to an additive constant, to mutual information (MI)

$$I(M; W = w) = H(M) - H(M | W = w)$$
constant



<u>A Measure of Lexical Ambiguity</u>

We operationalise le al ambiguity as the half-pointwise le al ambiguity as the

H(M | W=w)

Equation a word p to an additive constant, to Information Information (MI)

$$I(M; W = w) = H(M) - H(M | W = w)$$

constant



How to Measure Lexical Ambiguity?

WordNet


WordNet

• Discrete senses

BERT

• Continuous-meaning space

WordNet

- Discrete senses
- Hand-annotated



- Continuous-meaning space
- No hand annotation required!

WordNet

- Discrete senses
- Hand-annotated
- Only available in high-resource languages

BERT

- Continuous-meaning space
- No hand annotation required!
- Easily obtainable for new languages

WordNet

- Discrete senses
- Hand-annotated
- Only available in high-resource languages
- We assume an uniform distribution over senses

 $H(M | W=w) \approx \log_2(\#senses[w])$

BERT

- Continuous-meaning space
- No hand annotation required!
- Easily obtainable for new languages

WordNet

- Discrete senses
- Hand-annotated
- Only available in high-resource languages
- We assume an uniform distribution over senses

 $H(M | W=w) \approx \log_2(\#senses[w])$

BERT

- Continuous-meaning space
- No hand annotation required!
- Easily obtainable for new languages
- Assume embeddings are the word meaning: m ≈ BERT(p∘w∘s)

WordNet

- Discrete senses
- Hand-annotated
- Only available in high-resource languages
- We assume an uniform distribution over senses

 $H(M | W=w) \approx \log_2(\#senses[w])$

BERT

- Continuous-meaning space
- No hand annotation required!
- Easily obtainable for new languages
- Assume embeddings are the word meaning: m ≈ BERT(p∘w∘s)
- We use a Gaussian approximation (max-entropy upper bound)

H(M | W=w) ≈ H(N(µw, Σw)) .= ½ log₂ det (2πeΣw)

How well do the BERT and WordNet measures of lexical ambiguity correlate *with each other*?

How well do the BERT and WordNet measures of lexical ambiguity correlate *with each other*?

\rightarrow Relatively well!

Language	# Types	Pearson	Spearman
Arabic	836	0.25**	0.30**
English	6995	0.40**	0.40**
Finnish	1247	0.06*	0.07*
Indonesian	3308	0.12**	0.13**
Persian	2648	0.14**	0.13**
Portuguese	3285	0.13**	0.13**

How well do the BERT and WordNet measures of lexical ambiguity correlate *with each other*?

\rightarrow Relatively well!

Language	# Types	Pearson	Spearman
Arabic	836	0.25**	0.30**
English	6995	0.40**	0.40**
Finnish	1247	0.06*	0.07*
Indonesian	3308	0.12**	0.13**
Persian	2648	0.14**	0.13**
Portuguese	3285	0.13**	0.13**



How well do the BERT and WordNet measures of lexical ambiguity correlate *with each other*?

\rightarrow Relatively well!

Language	# Types	Pearson	Spearman
Arabic	836	0.25**	0.30**
English	6995	0.40**	0.40**
Finnish	1247	0.06*	0.07
Indonesian	3308	0.12**	Side
Persian	2648	0.14**	whe
Portuguese	3285	0.13**	0.13**

 $\overline{}$

<u>A Measure of Contextual Uncertainty</u>

We operationalise **contextual uncertainty** as the half-pointwise entropy:

H(W=w | C)

A Measure of Contextual Uncertainty

We operationalise **contermuter uncertainty** as the half-pointwise entremotion H(W=w | C)

Average uncertainty of a word in all its contexts

A Measure of Contextual Uncertainty



Average uncertainty of a word in all its contexts

May be approximated with a cloze language model

• This uses *bidirectional* context, which is different than most previous work

Our Empirical Findings

In WordNet

In WordNet

• Yes, in 5 of the 6 analysed languages!

Language		# Types	Pearson	Spearman
Arabic	(ar)	836	-0.14**	-0.15**
English	(en)	6995	-0.07**	-0.11**
Finnish	(fi)	1247	0.01	0
Indonesian	(id)	3308	-0.09**	-0.14**
Persian	(fa)	2648	-0.11**	-0.12**
Portuguese	(pt)	3285	-0.10**	-0.11**

In WordNet

• Yes, in 5 of the 6 analysed languages!

Language		# Types	Pearson	Spearman
Arabic	(ar)	836	-0.14**	-0.15**
English	(en)	6995	-0.07**	-0.11**
Finnish	(fi)	1247	0.01	0
Indonesian	(id)	3308	-0.09**	-0.14**
Persian	(fa)	2648	-0.11**	-0.12**
Portuguese	(pt)	3285	-0.10**	-0.11**



In **BERT**

In BERT

• Yes, in all the 18 analysed languages!

Language		# Types	Pearson	Spearman
Afrikaans	(af)	4505	-0.41**	-0.52**
Arabic	(ar)	10181	-0.33**	-0.41**
Bengali	(bn)	8128	-0.43**	-0.44**
English	(en)	7097	-0.33**	-0.35**
Estonian	(et)	4482	-0.40**	-0.44**
Finnish	(fi)	3928	-0.38**	-0.45**
Hebrew	(he)	13819	-0.34**	-0.37**
Indonesian	(id)	4524	-0.45**	-0.57**
Icelandic	(is)	3578	-0.44**	-0.46**
Kannada	(kn)	9695	-0.42**	-0.41**
Malayalam	(ml)	6203	-0.47**	-0.46**
Marathi	(mr)	5821	-0.39**	-0.40**
Persian	(fa)	6788	-0.39**	-0.49**
Portuguese	(pt)	5685	-0.31**	-0.45**
Tagalog	(tl)	3332	-0.45**	-0.50**
Turkish	(tr)	4386	-0.40**	-0.46**
Tatar	(tt)	2997	-0.34**	-0.39**
Yoruba	(yo)	417	-0.55**	-0.64**

In BERT

• Yes, in all the 18 analysed languages!



_anguage		# Types	Pearson	Spearman
Afrikaans	(af)	4505	-0.41**	-0.52**
Arabic	(ar)	10181	-0.33**	-0.41**
Bengali	(bn)	8128	-0.43**	-0.44**
English	(en)	7097	-0.33**	-0.35**
Estonian	(et)	4482	-0.40**	-0.44**
Finnish	(fi)	3928	-0.38**	-0.45**
Johrow	(ha)	10010	-0.34**	-0.37**
		-0.45**	-0.57**	
Tagalog			-0.44**	-0.46**
00 -		-0.42**	-0.41**	
00 -	•	-0.47**	-0.46**	
00 -	•	-0.39**	-0.40**	
00-	5.	-0.39**	-0.49**	
0-		-0.31**	-0.45**	
00-		-0.45**	-0.50**	
001	10	-0.40**	-0.46**	
			-0.34**	-0.39**

41/

(yo)

Yoruba

-0.55 **

-0 64**

In BERT

• Yes, in all the 18 analysed languages!



Language		# Types	Pearson	Spearman
Afrikaans	(af)	4505	-0.41**	-0.52**
Arabic	(ar)	10181	-0.33**	-0.41**
Bengali	(bn)	8128	-0.43**	-0.44**
English	(en)	7097	-0.33**	-0.35**
Estonian	(et)	4482	-0.40**	-0.44**
Finnish	(fi)	3928	-0.38**	-0.45**
Hobrow	(ha)	10010	-0.34**	-0.37**
		-0.45**	-0.57**	
Tagalog			-0.44**	-0.46**
800 -		-0.42**	-0.41**	
600	•	-0.47**	-0.46**	
400 -	•	-0.39**	-0.40**	
			-0.39**	-0.49**
			-0.31**	-0.45**
			-0.45**	-0.50**
400 1	10	15	-0.40**	-0.46**
			-0.34**	-0.39**

41/

(yo)

Yoruba

-0.55**

-0 64**

A Functionalist Derivation of the Trade-Off



• **Clarity** is the functionalist principle that a listener be able to reconstruct the speaker's intended meaning



• **Clarity** is the functionalist principle that a listener be able to reconstruct the speaker's intended meaning

• Information-theoretically, we operationalise **clarity** as:

H(M | W, C)



- **Clarity** is the functionalist principle that a listener be able to reconstruct the speaker's intended meaning
- Information-theoretic ning, we open ionalise clarity as:
 Meaning, we open ionalise clarity as:

which is the uncertainty of the meaning, given the context and the word



• **Robustness** is the functionalist principle that a speaker's utterance should be resilient to noise



• **Robustness** is the functionalist principle that a speaker's utterance should be resilient to noise

• We operationalise **robustness** as a tripartite MI:

I(M; C; W = w)



• **Robustness** is the functionalist principle that a speaker's utterance should be resilient to noise

• We operationalise **robustness** as a tripartite MI:

I(M; C; W = w)

which is the **redundant** information shared by meaning, context and word

• Assume language is clear, i.e. H(M | W, C) = 0

- Assume language is **clear**, i.e. **H(M | W, C) = 0**
- Assume language is **robust**, i.e. $I(M; C; W = w) \ge k$

- Assume language is **clear**, i.e. **H(M | W, C) = 0**
- Assume language is **robust**, i.e. $I(M; C; W = w) \ge k$

constant

- Assume language is **clear**, i.e. **H(M | W, C) = 0**
- Assume language is **robust**, i.e. $I(M; C; W = w) \ge k$

constant

• We show:

I(M; C; W = w) = H(M) - H(M | W = w) - H(W = w | C)

- Assume language is **clear**, i.e. **H(M | W, C) = 0**
- Assume language is **robust**, i.e. $I(M; C; W = w) \ge k$

constant

• We show: I(M; C; W = w) = H(M) - H(M | W = w) - H(W = w | C)constant

- Assume language is **clear**, i.e. **H(M | W, C) = 0**
- Assume language is **robust**, i.e. $I(M; C; W = w) \ge k$

We show: I(M; C; W = w) = H(M) - H(M | W = w) - H(W = w | C)constant

constant

• Thus, $H(M | W = w) + H(W = w | C) \le H(M) - k$

- Assume language is **clear**, i.e. **H(M | W, C) = 0**
- Assume language is **robust**, i.e. $I(M; C; W = w) \ge k$

constant






Language	H(W)	MI(W;V)	U(W V)
English	3.401	0.11	3.24%
German	3.195	0.168	5.26%
Dutch	3.245	0.156	4.82%

Test	H(W)	MI(W;V)	U(W V)
Africa	3.773	0.011	0.28%
Americas	3.901	0.007	0.17%
Eurasia	3.999	0.015 [‡]	0.38%
Pacific	3.755	0.016 [‡]	0.42%
Average	3.857	0.012 [‡]	0.31%









Language	H(W)	MI(W;V)	U(W V)
English	3.401	0.11	3.24%
German	3.195	0.168	5.26%
Dutch	3.245	0.156	4.82%

Test	H(W)	MI(W;V)	U(W V)
Africa	3.773	0.011	0.28%
Americas	3.901	0.007	0.17%
Eurasia	3.999	0.015 [‡]	0.38%
Pacific	3.755	0.016 [‡]	0.42%
Average	3.857	0.012 [‡]	0.31%









Language	H(W)	MI(W;V)	U(W V)
English	3.401	0.11	3.24%
German	3.195	0.168	5.26%
Dutch	3.245	0.156	4.82%

Test	H(W)	MI(W;V)	U(W V)
Africa	3.773	0.011	0.28%
Americas	3.901	0.007	0.17%
Eurasia	3.999	0.015 [‡]	0.38%
Pacific	3.755	0.016 [‡]	0.42%
Average	3.857	0.012 [‡]	0.31%









Language	H(W)	MI(W;V)	U(W V)
English	3.401	0.11	3.24%
German	3.195	0.168	5.26%
Dutch	3.245	0.156	4.82%

Test	H(W)	MI(W;V)	U(W V)
Africa	3.773	0.011	0.28%
Americas	3.901	0.007	0.17%
Eurasia	3.999	0.015 [‡]	0.38%
Pacific	3.755	0.016 [‡]	0.42%
Average	3.857	0.012 [‡]	0.31%









And thanks to all co-authors: Rowan Hall Maudslay, Brian Roark, Damián Blasi, Ryan Cotterell, Arya D. McCarthy, Søren Wichmann