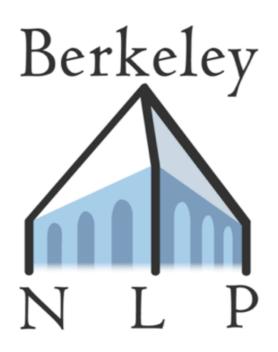
Linguistics: Speech and Lexical Semantics

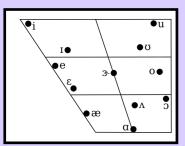


EECS 183/283a: Natural Language Processing

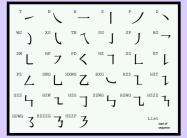




Stokoe notation of ASL



GA English vowels



Chinese stroke primitives, from Jiang et al. 2024

Phonemes/ Graphemes/





ASL mouth morphemes Large (CHA) and Small (OO)



Wikipedia (Annie Yang)

fán		
反	_	"anti-"

- 反 恐 [反恐] "anti-terror"
- făn jiàoquánde 反 教 权 的 [反教權的] "anti-clerical"
- fǎn fàxīsī
 反 法西斯 [反法西斯] "anti-fascist"

Morphemes /Lexemes



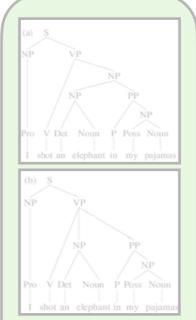
"language" in ASL

		ative		
	singular	plural		
		wir fegen		
present				
	er fegt	sie fegen		
		wir fegten		
preterite				
	er fegte			
imperative				

conjugation of German "fegen" (to sweep) (Wiktionary)

with affixes (Wikipedia)

Words





Constituents /Phrases





Utterances/ Sentences

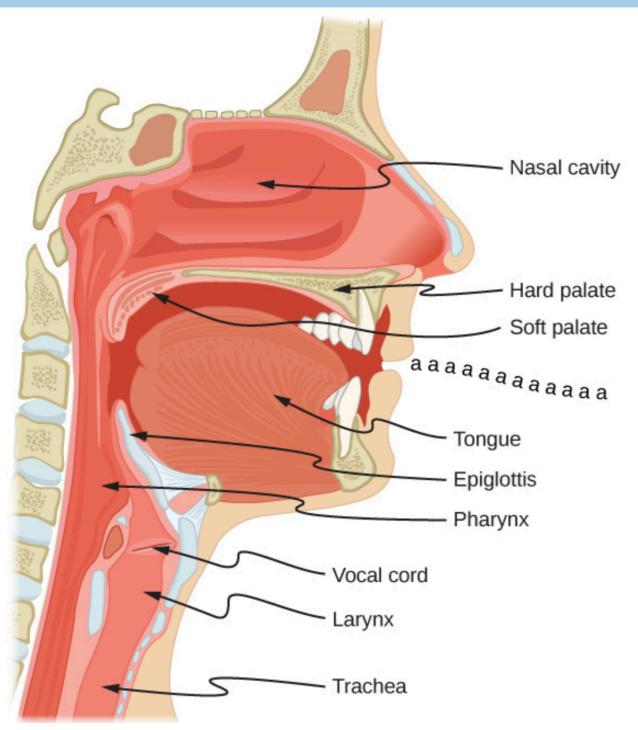
ਕੇਰਲ ਵਿੱਚ ਰਸਾਇਣਾਂ ਨਾਲ ਭਰਿਆ ਜਹਾਜ਼ ਡੱਬਿਆ, ਅਰਬ ਸਾਗਰ 'ਚ ਜੇ ਤੇਲ ਰਿਸਿਆਂ



Discourse/ Dialogue

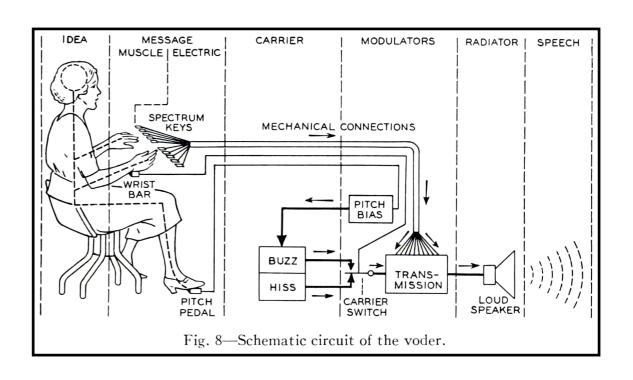
Speech Production





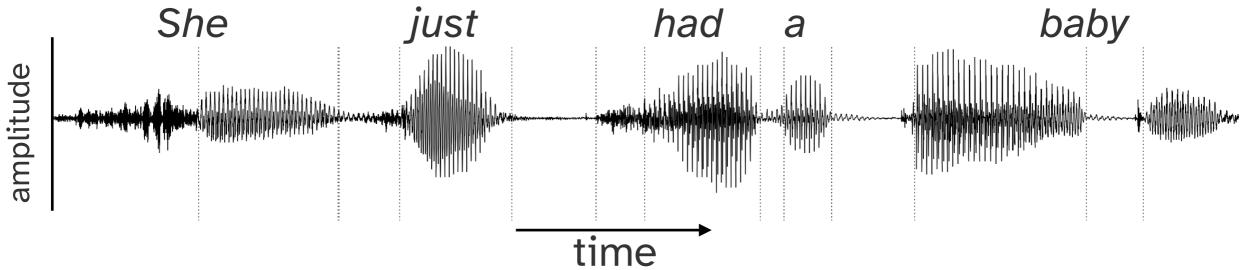
Vocal articulators that produce speech

- Air passing through vocal articulators produces speech
- Vocal folds, tongue, jaw, lips, velum are both independently and jointly controlled to produce different sounds
 - E.g., vocal fold vibration causes voicing
- The output of vocal articulation is an acoustic pressure wave



Speech Representations

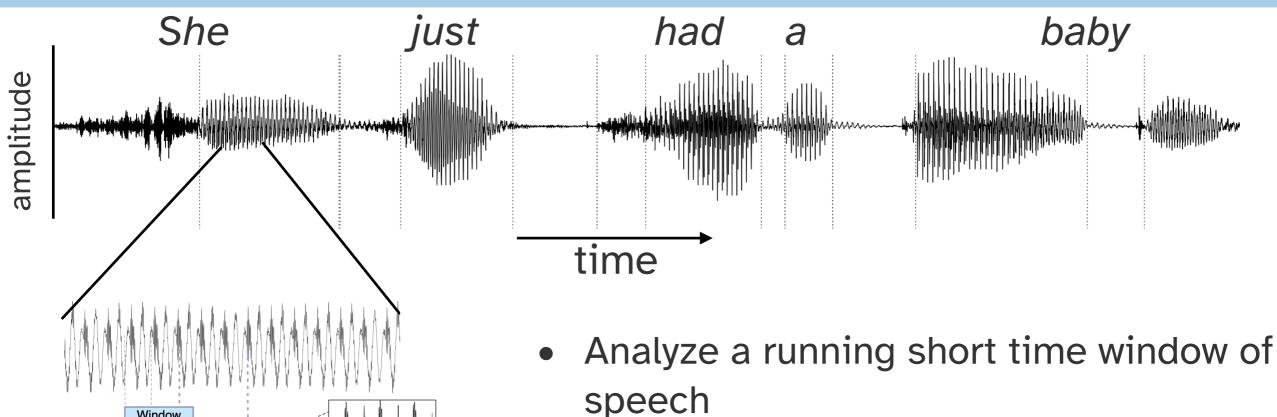




- The result of the vocal articulation is an acoustic pressure wave
- Speech can thus be represented as an acoustic waveform
- Waveforms are continuous time series cannot be easily analyzed or interpreted, or computed with
- Signal processing can give more interpretable information

Speech Waveform





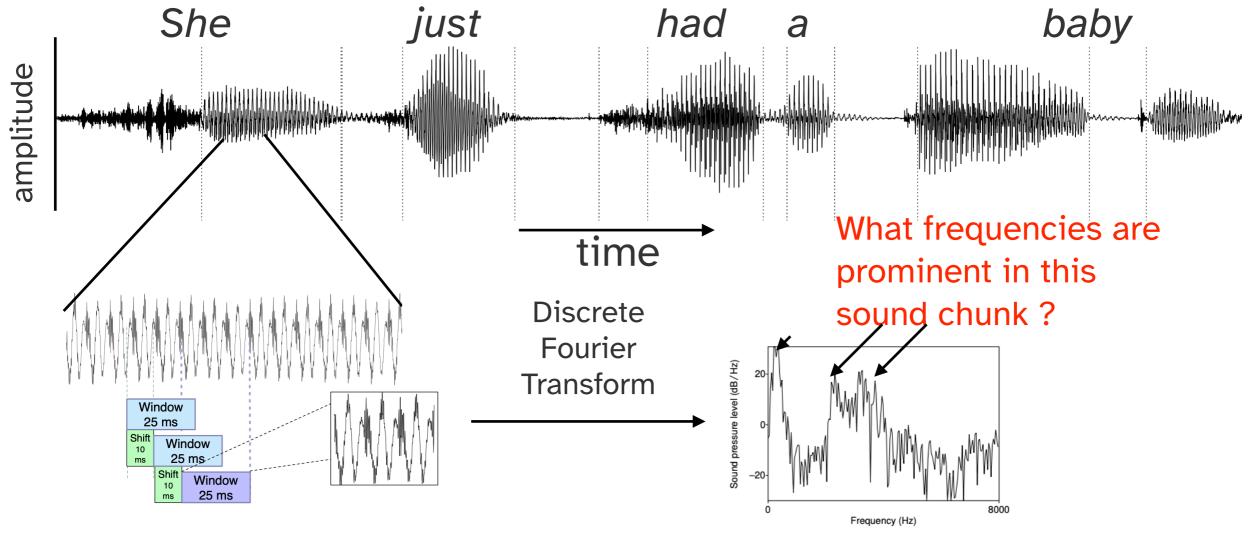
 Run Fourier Transform to convert time to time-frequency representation



Window

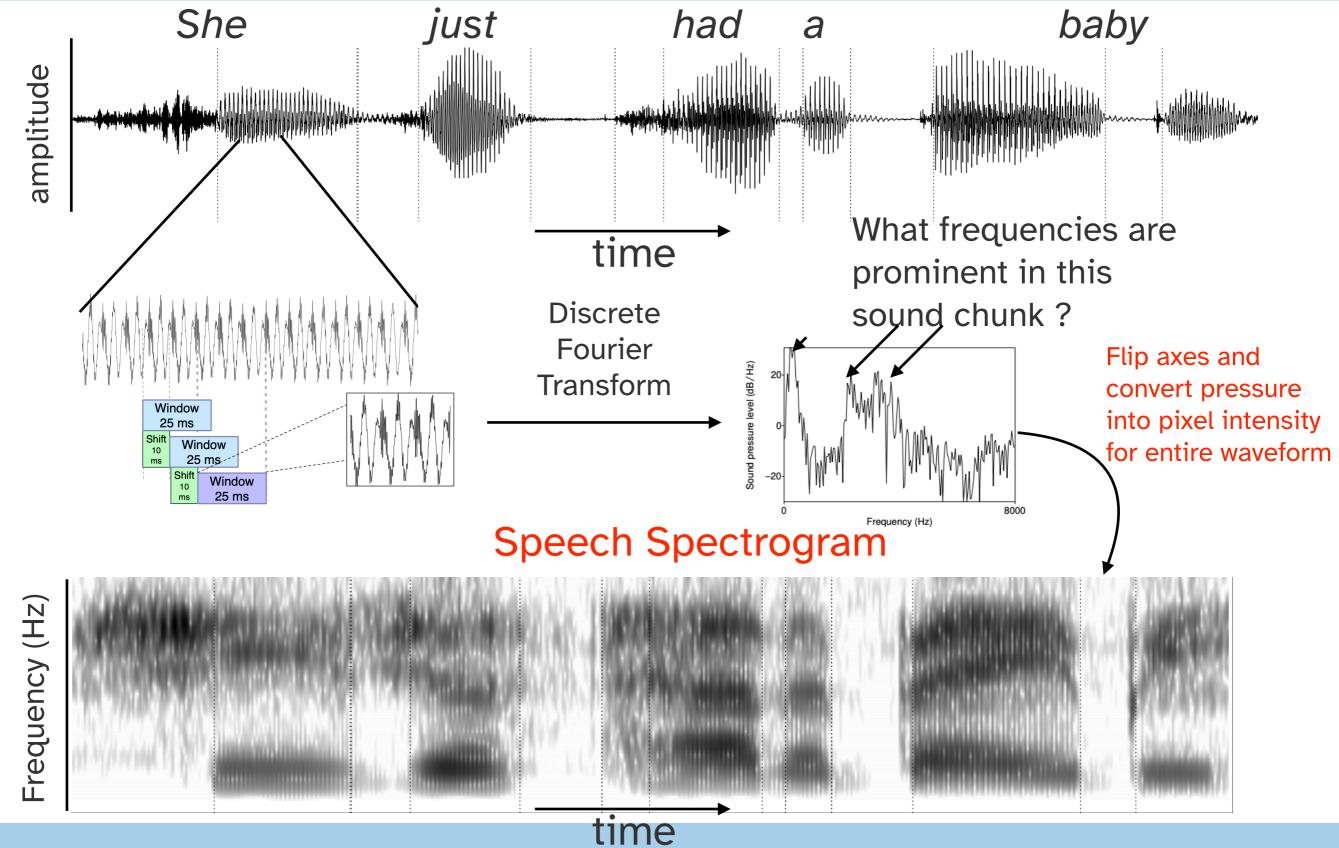
Speech Waveform





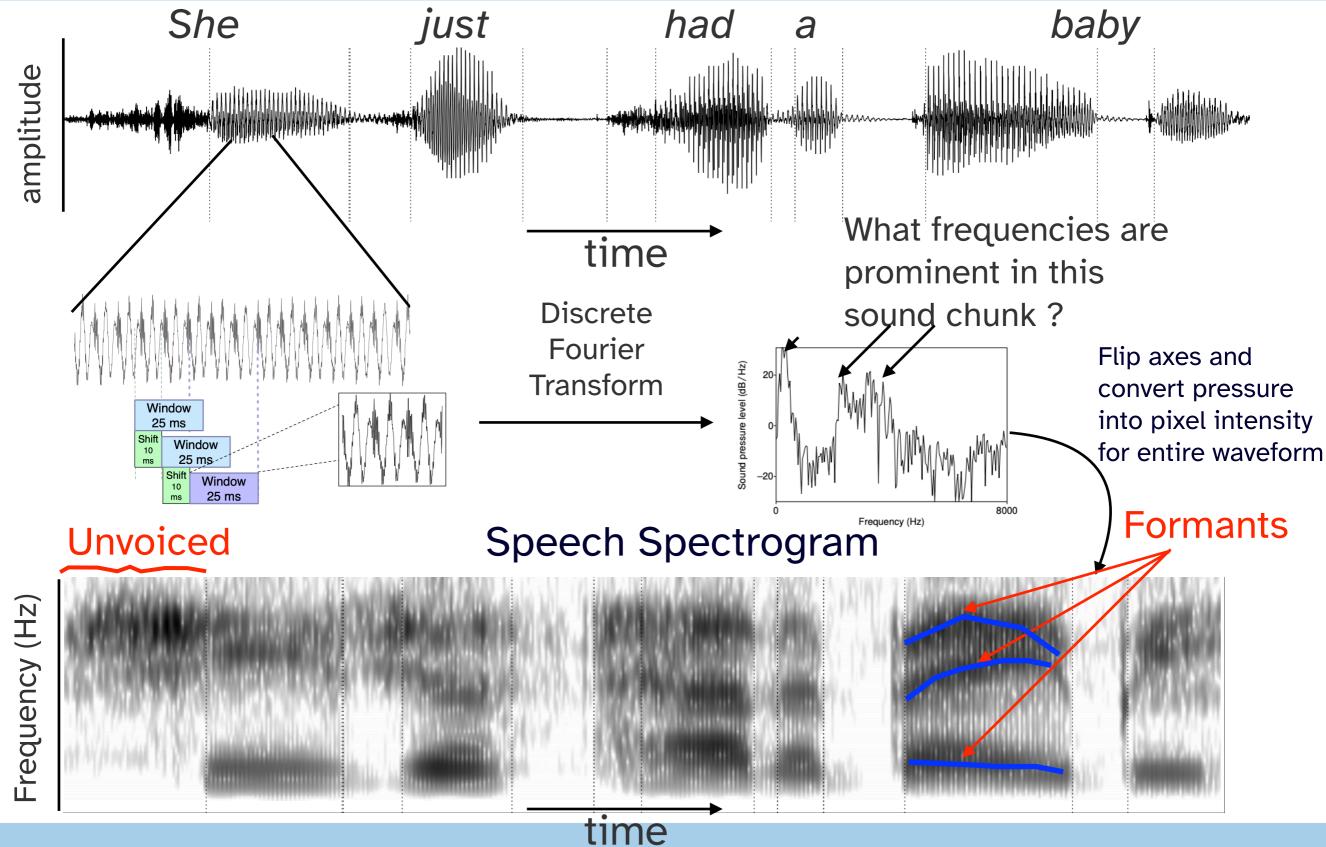
Speech Spectrogram





Speech Spectrogram

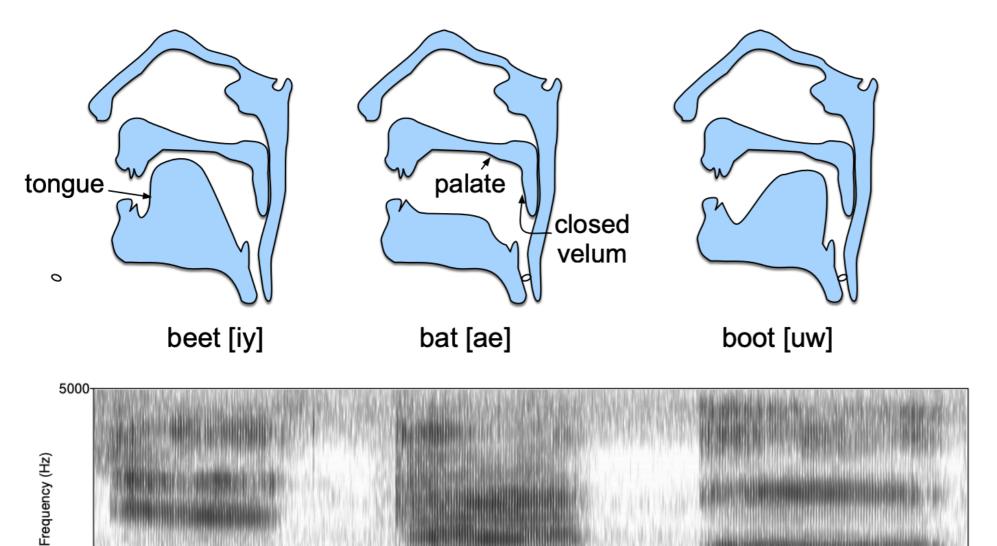




Phonetics



Study of speech sounds — their physical production, spectral and perceptual properties



ae (s)

iy

Articulatory Phonetics

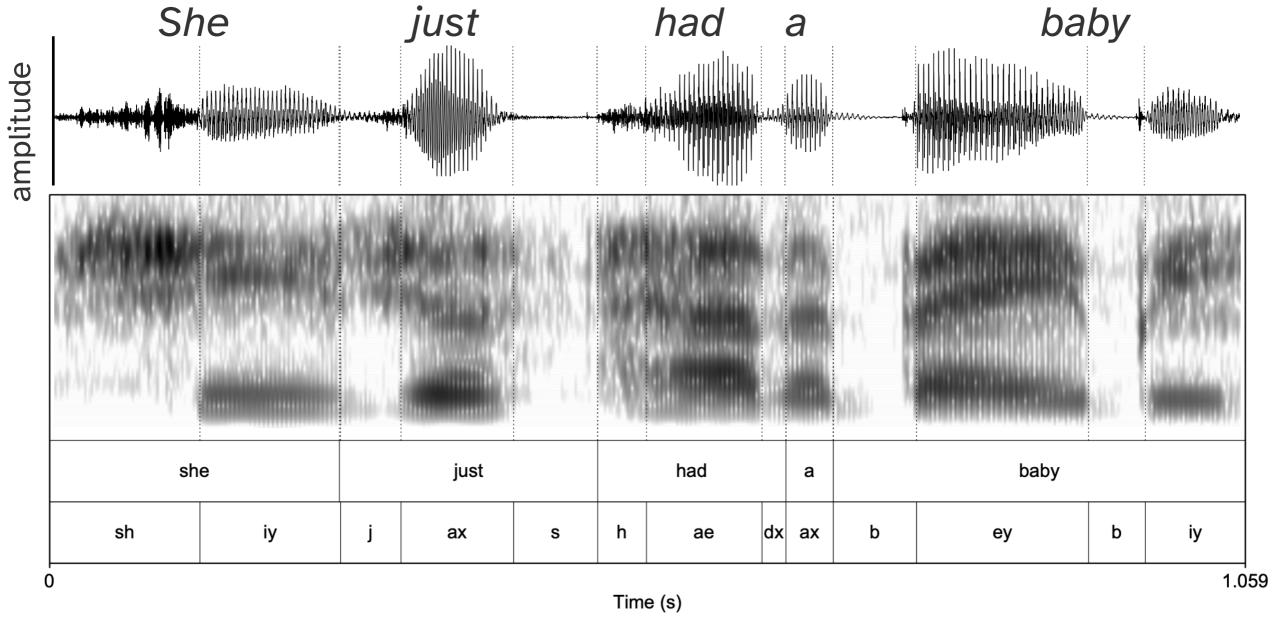
Acoustic Phonetics

2.81397

UW

Acoustic Phonetics

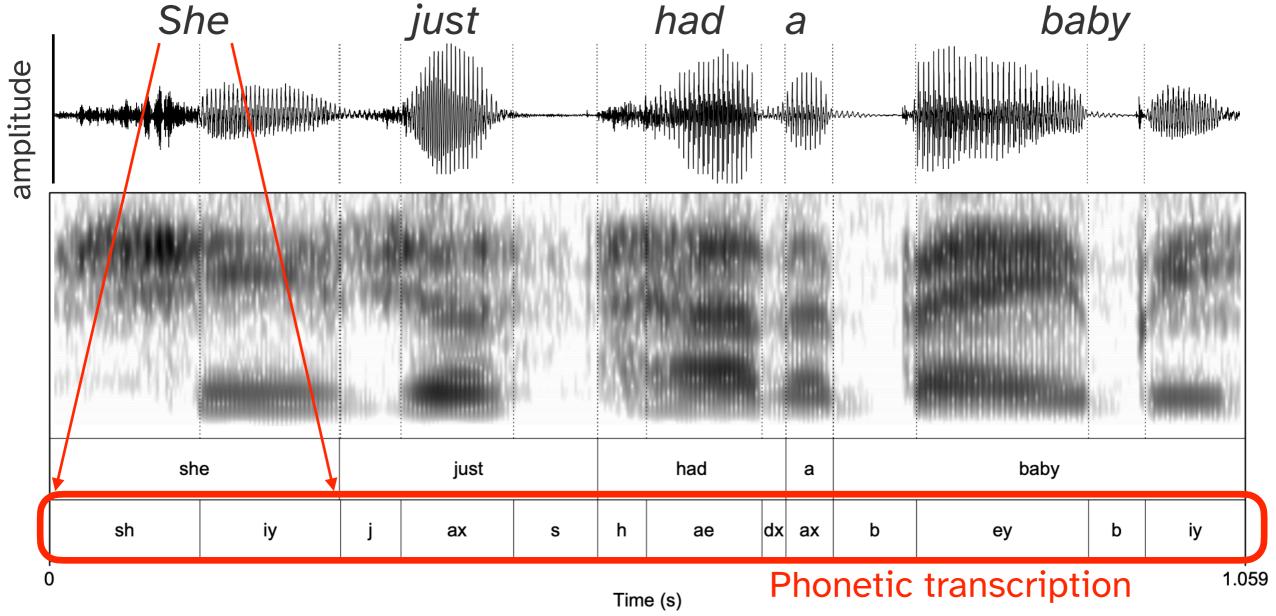




- Spectrogram reveals some segmental structure with distinct properties
- These are phonemes perceptually distinct speech sounds

Acoustic Phonetics





- Spectrogram reveals some segmental structure with distinct properties
- These are phonemes perceptually distinct speech sounds

International Phonetic Alphabet (IPA

- Phoneticians compiled a common set of sounds used to codify different speech sounds (across languages)
- English spelling (aka orthography) is not phonetic: about 40 distinct phonemes represented by 26 graphemes
 - About 16 vowel sounds
 - About 24 consonant sounds



THE INTERNATIONAL PHONETIC ALPHABET (revised to 2020)

CONSONAN	rs (PULM	ONIC)								⊚⊕⊚	2020 IPA	
	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retroflex Palatal		Velar	Uvular	Pharyngeal	Glottal	
Plosive	рb			t d		t d	С Ј	k g	q G		3	
Nasal	m	m		n		η	n	ŋ	N			
Trill	В			\mathbf{r}					R			
Tap or Flap		V		ſ		τ						
Fricative	φβ	f v	θð	s z	J 3	ş z	çj	ху	Χк	ħΥ	h fi	
Lateral fricative				4 В								
Approximant		υ		J		J	j	щ				
Lateral				1		1	λ	L				

Symbols to the right in a cell are voiced, to the left are voiceless. Shaded areas denote articulations judged impossible

CONSONANTS (NON-PULMONIC)

Clicks	Voiced implosives	Ejectives
O Bilabial	6 Bilabial	texamples:
Dental	d Dental/alveolar	p' Bilabial
! (Post)alveolar	f Palatal	t' Dental/alveolar
+ Palatoalveolar	g Velar	k' Velar
Alveolar lateral	G Uvular	S' Alveolar fricative

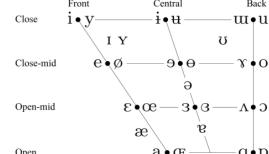
OTHER SYMBOLS

M Voiceless labial-velar fricative C Z Alveolo-palatal fricatives J Voiced alveolar lateral flap Simultaneous and X U Voiced labial-palatal approximant

H Voiceless epiglottal fricative

Yoiced epiglottal fricative can be represented by two symbols joined by a tie bar if necessary P Epiglottal plosive

VOWELS



founa tifan

SUPRASEGMENTALS Primary stress

1	Secondary stress									
I	Long	(eĭ							
•	Half-long	(e^{\centerdot}							
0	Extra-short	Č	ĕ							
	Minor (foot) gro	up								
ĺ	Major (intonation) group									
	Syllable break	J	$_{ m Ji.ækt}$							
\smile	Linking (absence	e of a bre	ak)							
	TONES AND W		CCENTS NTOUR							
é	or Extra	ě or /	Rising							
é		ê \	Falling							

DIACRITICS

0	Voiceless	ů ď	Breathy voiced b. a. Dental t. d.
~	Voiced	ş ţ	\sim Creaky voiced $\stackrel{b}{b}$ $\stackrel{a}{a}$ $\stackrel{\Box}{\Box}$ Apical $\stackrel{\Box}{\Box}$ $\stackrel{\Box}{\Box}$
h	Aspirated	$t^{h} d^{h}$	Linguolabial t d Laminal t d
,	More rounded	ò	w Labialized t^{w} d^{w} $^{\sim}$ Nasalized $ ilde{e}$
c	Less rounded	ç	$^{\mathrm{j}}$ Palatalized t^{j} d^{j} $^{\mathrm{n}}$ Nasal release d^{n}
	Advanced	ų	$^{\gamma}$ Velarized t^{γ} d^{γ} l Lateral release d^{l}
_	Retracted	ė	$^{\Gamma}$ Pharyngealized \mathbf{t}^{Γ} \mathbf{d}^{Γ} $^{\gamma}$ No audible release \mathbf{d}^{γ}
	Centralized	ë	~ Velarized or pharyngealized }
×	Mid-centralized	ě	Raised Q (I = voiced alveolar fricative)
	Syllabic	ņ	Lowered $\underset{\tau}{\mathbf{e}}$ ($\underset{\tau}{\beta}$ = voiced bilabial approximant)
^	Non-syllabic	é	Advanced Tongue Root Q
~	Rhoticity	or ar	Retracted Tongue Root P

Some diacritics may be placed above a symbol with a descender, e.g. $\check{\Pi}$

International Phonetic Alphabet (IPA



 Phoneticians compiled a common set of sounds used to codify different speech sounds (across languages)

Central Rotokas (Papua New Guinea): 12 consonants, 10 vowels

	Bilabial	Alveolar	Velar
Voiceless	р	t	k
Voiced	b	d	g

	Front	Central	Back
Close	i (iː)		u (u:)
Close-mid	e (eː)		o (oː)
Open		a (aː)	

Archi (Dagestan, Russia)

			Labial	Der	ntal	(Pos	-	(Pre-)velar		ı	Jvulaı	Epiglottal	Glottal	
				plain	lab.	ıb. plain lab. plain lab. plain lab. phar. phar.+		phar.+lab.							
Nasal		m	n												
	voice	d	b	d	d ^{w2}			g	g ^w						
Plosive	voicele	ess	р	t	tw			k	kw	q	qw	d٤	d _{2m}	?	?1
Piosive	fortis		p:1	t:1				k:1	kwr2	q':¹		q°'ı			
	ejective		p'	ť				k'	k ^w '	q'	qw'	q°'	q ^{°w}		
	voiceless	lenis		îs	îsw2	€	€	k ™	k ™w						
Affricate		fortis		îs:3											
Affricate		lenis		îs'	îsw'	ÎĴ'	î∫w'	k "	k ⊞w'						
	ejective	fortis		îs':1		î∫'r²									
		lenis		s	s ^{w2}	ſ	ſw	91D 864	es W	χ	χw	Χ°	X,m		h
Fricative	voiceless	fortis		SI	sw _z 2	ſː	∫wː	12 ¥	NR W.	Χī	Χ _w ï	χ'n	X _{em} ī		
	voiced			z	Zw	3	3 ^w	Ļ1		R	R _m 5	R _c	R _{JM}		
	Trill			r										н	
App	Approximant			-1		j			w						

THE INTERNATIONAL PHONETIC ALPHABET (revised to 2020)

CONSONANT	rs (P	ULM	ONIC))																⊚⊕⊚	2020	IPA
	Bila	abial	Labiodental Dental Alveolar Postalveolar Retroflex Palatal Velar				Alveolar Postalve		al Alveolar Postalveolar I		Postalveolar		lar	Uv	ular	Phary	ngeal	Glottal				
Plosive	р	b					t	d			t	d	С	J	k	g	q	G			3	
Nasal		m		ŋ				n				η		n		ŋ		N				
Trill		В						r										\mathbf{R}				
Tap or Flap				\mathbf{V}				\mathbf{l}				\mathfrak{r}										
Fricative	ф	β	f	V	θ	ð	s	\mathbf{Z}	ſ	3	ş	Z,	ç	j	x	γ	χ	\mathbf{R}	ħ	?	h	ĥ
Lateral fricative							ł	В														
Approximant				υ				J				Ţ		j		щ						
Lateral approximant								1				l		Λ		L						

Symbols to the right in a cell are voiced, to the left are voiceless. Shaded areas denote articulations judged impossible.

VOWELS

Close-mid

Open-mid

Open

Close

CONSONANTS (NON-PULMONIC)

Clicks	Voiced implosives	Ejectives
O Bilabial	6 Bilabial	texamples:
Dental	d Dental/alveolar	p' Bilabial
! (Post)alveolar	f Palatal	t' Dental/alveolar
+ Palatoalveolar	g Velar	k' Velar
Alveolar lateral	G Uvular	S' Alveolar fricative

OTHER SYMBOLS

M Voiceless labial-velar fricative W Voiced labial-velar approximant

U Voiced labial-palatal approximant H Voiceless epiglottal fricative

Yoiced epiglottal fricative P Epiglottal plosive

C Z Alveolo-palatal fricatives J Voiced alveolar lateral flap

Simultaneous and X

can be represented by two symbols joined by a tie bar if necessary

SUPRASEGMENTALS

Central

w•u

founa tifan

- 1	Secondary stress	
I	Long	er
•	Half-long	e^{\centerdot}
0	Extra-short	ĕ
	Minor (foot) group	
	Major (intonation) grou	p
	Syllable break	лі.ækt
\cup	Linking (absence of a b	reak)

TONES AND WORD ACCENTS

e or lagh e or lagh	Rising
é ↑ High ê ∨	Fallin
	High rising
è low è /	Low rising
	Rising
↓ Downstep	

↑ Upstep

Global fall

DIACRITICS

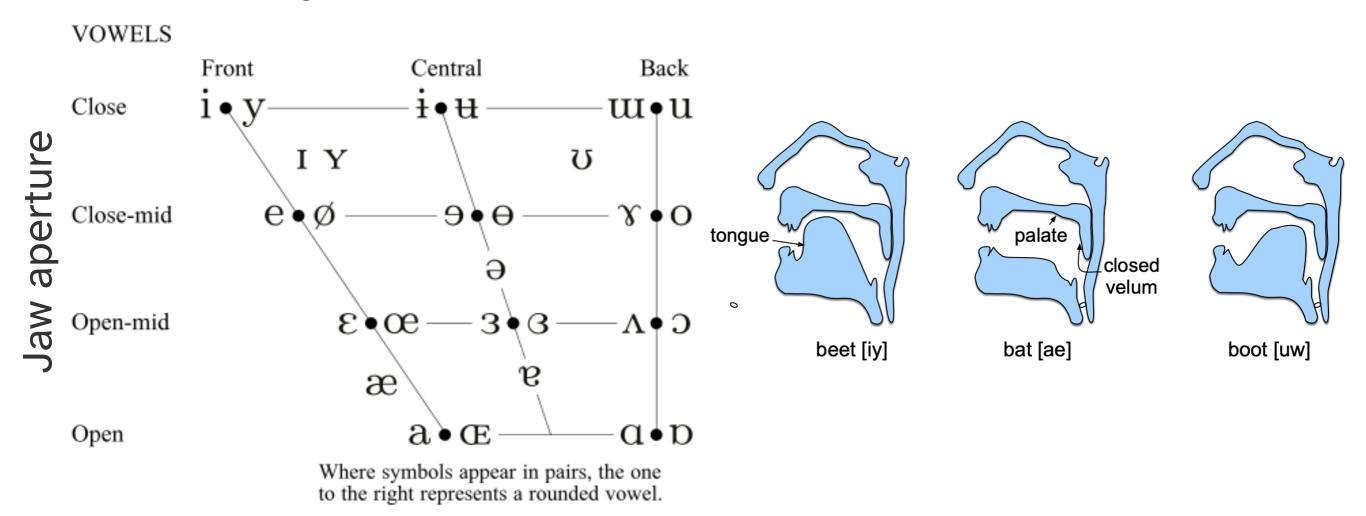
	Miles									
۰ ۱	Voiceless	ņ ģ		Breathy voiced	ÿ	ä	_	Dental	ţ	ď
~ \	Voiced	ş ţ	~	Creaky voiced	Ď	a	J	Apical	\mathbf{t}	$\dot{\mathbf{q}}$
h A	Aspirated	$t^h d^h $	_	Linguolabial	ţ	ğ	-	Laminal	ţ	d
, N	More rounded	Ş	w	Labialized	t^{w}	d^{w}	~	Nasalized		ẽ
, I	less rounded	ç	j	Palatalized	$\mathrm{t^{j}}$	d^{j}	n	Nasal release		d^{n}
	Advanced	ų	Y	Velarized	\mathbf{t}^{γ}	$\mathrm{d}^{\scriptscriptstyle{\gamma}}$	1	Lateral release		d^{l}
_ F	Retracted	ė	ſ	Pharyngealized	$\mathrm{t}^{\mathfrak{l}}$	\mathbf{q}_{ϵ}	٦	No audible releas	e	d^{\lnot}
(Centralized	ë	~	Velarized or phary	ngeali	zed	ł			
× 1	Mid-centralized	ě	_	Raised	ę	(<u>I</u> =	voic	ed alveolar fricativ	e)	
S	Syllabic	ņ		Lowered	ę	$(\bar{\beta} =$	voic	ed bilabial approxi	mant	:)
_ 1	Non-syllabic	é	4	Advanced Tongue	Root	ę				
√ F	Rhoticity	or ar	F	Retracted Tongue	Root	ę				

Some diacritics may be placed above a symbol with a descender, e.g. $\ddot{\mathbf{n}}$

International Phonetic Alphabet (IPA)

- Vowels are characterized by jaw position and tongue shape
- Some vowels also use lips (eg. sound uw in cool)

Tongue frontness



Manner of articulation

International Phonetic Alphabet (IPA)



alveolar

palatal

@ ⊕ @ 2020 IPA

velar

glottal

(nasal tract)

bilabial

- Consonants are characterized by place and manner of articulation
 - /p/ is caused by constriction at lips (labial)
 - /p/ is caused by sudden release of air (plosive) dental



CONSONANTS (PULMONIC)

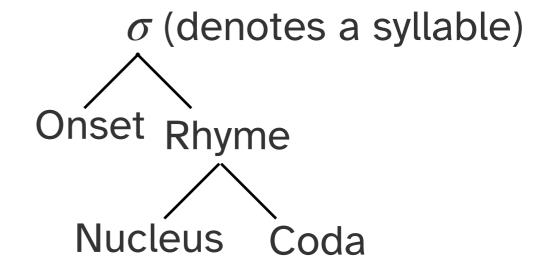
	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retro	oflex	Palatal	Velar	Uvular	Pharyngea	l Glottal
Plosive	рb			t d		t	d	С Ј	k g	q G		3
Nasal	m	m		n			η	n	ŋ	N		
Trill	В			r						R		
Tap or Flap		V		ſ			τ					
Fricative	φβ	f v	θ δ	s z	\int 3	ş	Z,	çj	ху	Χв	ħΥ	h fi
Lateral fricative				4 3								
Approximant		υ		J			J	j	щ			
Lateral approximant				1			l	Λ	L			

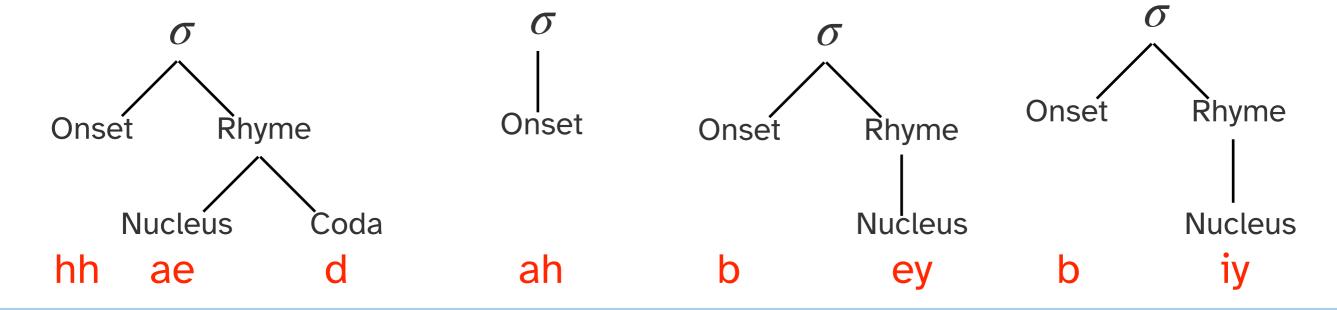
Symbols to the right in a cell are voiced, to the left are voiceless. Shaded areas denote articulations judged impossible.

Lexical Phonology



- Phonology is the study of rules that govern the organization of sounds in a language (Phonemes \rightarrow Syllables \rightarrow Words)
- English syllable structure: (C3)V(C4)
 /strεŋkθs/
- Hawaiian syllable structure: (C)V(C)
- Georgian syllable structure: (C8)V(C5)
 გვბრდღვნის /ˈgvbrdɣvnis/





Letters to Sounds



- Pronunciation dictionaries (often made by linguists) give the syllables and phonemes within each word in vocabulary
 - CMU Phonetic Dictionary gives the syllabic and phonetic spellings for >110K words in English
 - ML based phonetizers are built on such phonetic dictionaries

```
Graphemes She just had a baby

IPA ∫ix dʒ∧st hæd ə 'beɪbi

Arpabet sh iy jh ah s t h ae d ah b ey b iy

Arpabet is an ASCII friendly representation of IPA
```

Letters to Sounds





0:01

• (General American) IPA(key): /'b3kli/

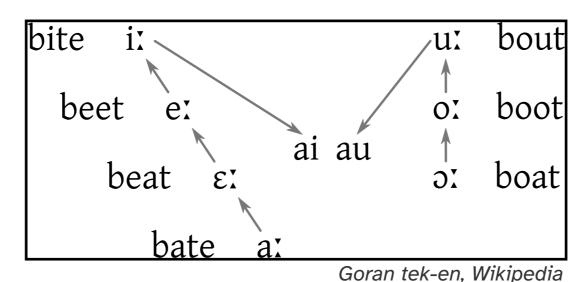
Audio (Southern England):



Sound Variation



- What might cause differences in pronunciation?
 - Dialectical differences
 - Native language when speaking a non-native language
 - Context formality, noise levels, etc.
 - Language change over time



We learn sounds before we are born!

	Vowel pronunciation						
Word	Late Middle English before the GVS	Modern English after the GVS					
b <i>i</i> te	[iː]	[aɪ]					
m <i>ee</i> t	[eː]	[iː]					
m <i>ea</i> t	[:3]						
ser <i>e</i> ne	[6,]						
m <i>a</i> te	[aː]	[eɪ]					
<i>ou</i> t	[uː]	[aʊ]					
b <i>oo</i> t	[oː]	[uː]					
b <i>oa</i> t	[ːc]	[00]					
st <i>o</i> ne	[0,]	[ooj					

Phonemes and Graphemes

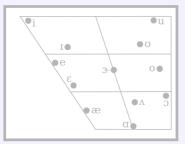


- Words are composed of atomic units based on sound (for spoken languages)
- Sounds are a function of how we move our vocal tracts and mouth anatomy
- Languages have distinct sets of possible sounds (phonemic inventory)
- And "rules" governing which sound sequences are likely (syllable structure)



$$\begin{split} & B_0 \ L & \text{i. } & X_1 X_1 \dot{a} \\ & B_0 \ L & \text{i. } & X_2 X_1 \dot{a} \\ & B_1 \ B_2 & \text{i. } & X_1 X_1 \dot{a} \\ & B_2 \ L & \text{i. } & X_2 X_1 \dot{a} \\ & B_3 \ L & \text{i. } & X_2 X_2 \dot{a} \\ & B_4 \ L & \text{i. } & X_3 X_1 \dot{a} \\ & B_4 \ L & \text{i. } & X_2 X_2 \dot{a} \\ & B_5 \ L & \text{i. } & X_3 X_1 \dot{a} \\ & B_6 \ L & \text{i. } & X_1 X_2 \dot{a} \\ & B_7 \ L & \text{i. } & X_2 X_2 \dot{a} \\ & B_8 \ L & \text{i. } & X_1$$

Stokoe notation of ASL



GA English vowels

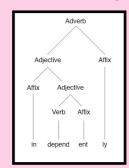
Chinese stroke primitives, from Jiang et al. 2024

Phonemes/ Graphemes/





ASL mouth morphemes Large (CHA) and Small (OO) learnhowtosign.com



Wikipedia (Annie Yang)

făn 反 — "anti-"

- făn kŏng ● 反 恐 [反恐] — "anti-terror"
- fǎn jiàoquánde • 反 教 权 的 [反教權的] — "anti-clerical"
- fǎn fàxīsī • 反 法西斯 [反法西斯] — "anti-fascist"

Wikipedia

Morphemes / /Lexemes /



"language" in ASL <u>lifeprint.com</u>

	indicative				
	singular	plural			
	ich fege	wir fegen			
present	du fegst	ihr fegt			
	er fegt	sie fegen			
	ich fegte	wir fegten			
preterite	du fegtest	ihr fegtet			
	er fegte	sie fegten			
	feg (du)	f+ (5h-)			
imperative	fege (du)	fegt (ihr)			

conjugation of German "fegen" (to sweep) (Wiktionary)

aiar = to teach

ajari = to teach (imperative, locative)

aja**riian** = to teach (jussive, locative)

ajarkan = to teach (imperative, causative/applicative ajarkanlah = to teach (jussive, causative/applicative)

ajar**lah** = to teach (jussive, active

ajaran = teachings

belajar = to learn (intransitive, active)

diajar = to be taught (intransitive)

diajari = to be taught (transitive, locative)

diajarkan = to be taught (transitive, causative/applicati

dinalajari - to be studied (locative)

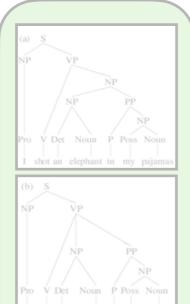
dipelajarkan = to be studied (causative/applicative)

mempelajari = to study (locative)

Indonesian "ajar" (to teach)

with affixes (Wikipedia)

Words



syntactic ambiguity, from UBC CPSC522



Universal Dependencies

Constituents
/Phrases



靜靜養貓 (YouTube)



Portal

Юная мышка Роза допела любимую песенку и теперь вслушивалась, как талло эко последних звуков. Роза была довольна собой наконец-то она нашла мего, откуда море отлично просматривалось. Вечерело, лучи солица, заходищего за ее спиной, окрасили воду « вазотителій и батровый цвета.

 Эгей, Роза, пожалуйте-ка ужинать. Не для того я старался, чтобы все это теперь простыло, значится... то есть отсырело. Нет уж, дудки!

Спутник Розы, крот Грумм, призывно помахал е дапой. Мышка подошла к маленькому костерку, на ко тором он стряпал, и принюхалась:

 Ого, лененики из дикого овса и суп из зелений Грумм улыбнулся, отчего его бархатная мордочки сморщилась, и помахал маленькой поварешкой, кото рую всегда носка за поясом наподобие меча;

Роза вядла глубокую раковину морского гресешка, наполненную ароматным суюм. Положный сюю леленну на плоский камень у костра, чтобы та не остыла, мышка отклебнула супа и покачала головой: — Ты хуже старой няньки, Грума Кланякинс Быось об акклад, если бы я позволила, ты бы меня и спать

Перед носом Розы замаячила поварешка.
— Так что же тебе, спать, это самос... и не надобно вовсе, так, что ли? Представь, что батюшка твой скажет, если я тебя домой доставлю всю усталую такую За потоятную э?

Utterances/ Sentences

B B C NEWS ਪੰਜਾਬੀ

ਾਂ ਚੀੜੀਓ ਪਾਠਕਾਂ ਦੀ ਪਸੰਦ ਚਾਰਜ ਕੌਮਾਂਸ਼ਰੀ

ਕੇਰਲ ਵਿੱਚ ਰਸਾਇਣਾਂ ਨਾਲ ਭਰਿਆ ਜਹਾਜ਼ ਡੁੱਬਿਆ, ਅਰਬ ਸਾਗਰ 'ਚ ਜੇ ਤੇਲ ਰਿਸਿਆਂ ਤਾਂ ਕੀ ਹੋਵੇਗਾ ਅਸਰ?

. . .

ਥੋਂ ਰਵਾਨਾ ਹੋਇਆ ਐੱਮਐੱਸਸੀ ਈਐੱਲਐੱਸਏ 3, 24 ਮਈ ਨੂੰ ਕੋਚੀ ਬੇਦਰਗਾਹ 'ਤੇ ਪਹੁੰਚਣ ਵਾਲਾ ਸੀ। (24 ਮਈ ਨੂੰ ਸਵੇਰੇ 12-15 ਵਜੇ ਭਾਰਤੀ ਤੱਟ ਰੱਖਿਅਕ ਨੂੰ ਜਹਾਜ਼ ਤੋਂ ਇੱਕ ਐਮਰਜੈਂਸੀ ਕਾਲ ਆਈ।

10 ਕੈਟੇਨਰਾਂ ਨੂੰ ਲੈ ਕੇ ਕੋਢੀ ਜਾਂਦੇ ਸਮੇਂ 184 ਮੀਟਰ ਲੰਬਾ ਐੱਮਐੱਸਸੀ ਈਐੱਲਐੱਸਏ 3 ਤੁੱਬਣ ਲੱਗਿਆ। ਦੋ ਕਾਰਗੇ ਜਹਾਜ਼ ਕੋਢੀ ਤੋਂ 38 ਸਮੁੰਦਰੀ ਮੀਲ ਦੱਖਣ-ਪੱਛਮ ਵਿੱਚ ਸੀ ਤਾਂ ਇਹ ਤਕਰਬੀਨ 26 ਡਿਕਰੀ ਤੱਕ ਤੋਂ ਵਿਸ਼ਾਹ ਸੀ।

ਾਰਤੀ ਤੋਂਟ ਕੈਖਿਅਕ ਨੇ ਤੁਰੰਤ ਨੇਡਲੇ ਜਹਾਜ਼ਾਂ ਨੂੰ ਬਚਾਅ ਕਾਰਜਾਂ ਲਈ ਭੇਜ ਦਿੱਤਾ ਸੀ। ਹਾਲਾਤ ਦੀ ਗਰਾਨੀ ਲਈ ਇੱਕ ਹਵਾਈ ਜਹਾਜ਼ ਵੀ ਮੌਜੂਦ ਸੀ। ਇਸ ਦੌਰਾਨ ਕਾਰਗੋ ਜਹਾਜ਼ ਲਗਾਰਾਰ ਝੁਕਦਾ ਗਿਆ ਸੇ ਕੁਝ ਤੱਕੇ ਸਮੇਦਰ ਵਿੱਚ ਇੱਗਣ ਲੱਗੇ।

ਭਾਰਤੀ ਜਲ ਲੈਨਾ ਨੇ 24 ਮਈ ਦੀ ਲਾਮ ਨੂੰ ਬਚਾਅ ਕਾਰਜ ਸ਼ੁਰੂ ਕੀਤੇ। ਦੇ ਜਹਾਜ਼, ਆਈਐੱਨਐੱਸ ਸ਼-ਅਤੇ ਆਈਐੱਨਐੱਸ ਸੁਜਾਰਾ, ਨੂੰ ਜਹਾਜ਼ ਵਿੱਚ ਸਵਾਰ 24 ਲੋਕਾਂ ਨੂੰ ਬਚਾਉਣ ਲਈ ਤੇਜਿਆ ਗਿਆ। ਆਈਐੱਨਐੱਸ ਸੁਜਾਰਾ ਸ਼ਾਮ 7 ਵਜੇ ਪਹੁੰਚਿਆ, ਜਦੋਂ ਕਿ ਆਈਐੱਨਐੱਸ ਸਤਪੁਰਾ ਰਾਤ 8 ਵਜੇ ਪਹੁੰਚ ਸੀਲਿਆ।

ਵੱਖਣ-ਪੱਛਮੀ ਮਾਨਸੂਨ, ਜੋ ਆਮ ਤੌਰ 'ਤੇ 1 ਜੂਨ ਨੂੰ ਸ਼ੁਰੂ ਹੁੰਦਾ ਹੈ, ਇਸ ਸਾਲ 24 ਮਈ ਨੂੰ ਸ਼ੁਰੂ ਹੋਇਆ ਸੰ ਇਸ ਲਈ ਸਮੁੰਦਰ ਦਾ ਮੌਸਮ ਖਰਾਬ ਸੀ।

ਆਈਐੱਟਿਐੱਸ ਸੁਜਾਰਾ 'ਦੇ ਕੈਪਟਨ ਅਰਜੂਨ ਸ਼ੇਖਰ ਨੇ ਖ਼ਬਰ ਦੇਸ਼ਸੀ ਦੋਐੱਨਆਈ ਨੂੰ ਦੱਸਿਆ,' ਸਾਨੂੰ ਪ੍ਰਤੀਯੂਲ ਹਾਲਾਤ ਦਾ ਸਾਹਮਣਾ ਕਰਨਾ ਪਿਆ। ਹਵਾ 74.08 ਕਿਲੋਮੀਟਰ ਪ੍ਰਤੀ ਘੰਟਾ (40 ਨਾਟ) ਦੀ ਰਫ਼ਟ ਨਾਲ ਵਾਗੇ ਹਹੀ ਸੀ। ਸਮੁੰਦਰ ਵਿੱਚ ਗੂੜਾ ਅਤੇ ਕੰਟੇਨਰ ਤੇਰ ਰਹੇ ਸਨ। ਇਸ ਕਾਰਨ ਰਾਤ ਨੂੰ ਜਹਾਜ਼ ਤੱਕ ਪਹੁੰਚਣਾ ਮੁਸ਼ਕਰ ਹੋ ਰਿਹਾ ਸੀ।"

ਾਜ਼ ਵਿੱਚ ਸਵਾਧਾ 2+ ਲੰਕਾਂ ਵਿੱਚੋਂ, 21 ਨੂੰ ਉਸ ਰਾਡ ਜ਼ਿਨ੍ਹਾਂ - ਜਿਸ਼ੀ ਜਾਨਲੰਦਾ ਸੰਦ ਦੇ ਬਾਚਾ ਲਿਆ ਲਿਆ। ਉਸ ਰਾਜਾ 'ਤੇ ਅਜੇ ਦੀ ਕੈਂਟੋਨਰ ਸਨ ਅਤੇ ਜਹਾਜ਼ ਪੂਰੀ ਭਰੁ ਾਂ ਨਹੀਂ ਕੁੱਬਿਆ ਸੀ, ਇਸ ਲਈ ਜਹਾਜ਼ ਮਾਸਟਰ, ਮੁੱਖ ਇੰਜੀਨੀਅਰ ਅਤੇ ਸਹਾਇਕ ਇੰਜੀਨੀਅਰ ਬਚਾਅ ਕਾਰਜਾਂ ਨੂੰ ਨੋਪਰੇ ਚਾਡਨ ਅਤੇ ਸਥਿਤੀ ਦੀ ਭਰਾਨੀ ਕਰਨ ਲਈ ਜਹਾਜ਼ 'ਤੇ ਹੀ ਰਹੇ।

ਤਿੰਨਾਂ ਨੇ ਭਾਰਤੀ ਤੱਟ ਰੱਖਿਅਕ ਅਤੇ ਭਾਰਤੀ ਜਲ ਸੈਨਾ ਦੀ ਨਿਗਰਾਨੀ ਹੇਠ ਜਹਾਜ਼ 'ਚ ਹੀ ਰਾਤ ਬਿਤਾਏ

BBC News in Punjabi



Portal 2 dialogues

CMAIR EID: Thank you. Does the Committee
have any questions? I do not see any. Thank you so
much for your presentation today.
Ms. RESNIK: We appreciate your time and
that you enabled us all to offer comments. Many
thanks.

CMAIR EID: Thank you.

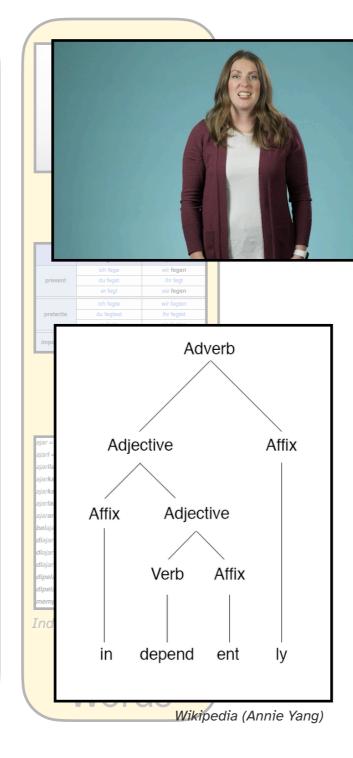
All right. We are now going to turn to
Carter Phillips, and we have now moved to Rule 29
comments.

MS. FMILLIPS: Judge Eid, can you see me and

US Supreme Court

Discourse/ Dialogue







ASL mouth morphemes Large (CHA) and Small (OO) <u>learnhowtosign.com</u>

Morphemes //Lexemes /

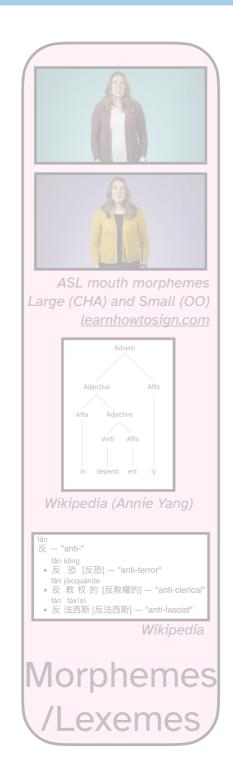
反 — "anti-"

fǎn kǒng

- 反 恐 [反恐] "anti-terror" făn jiàoquánde
- 反 教 权 的 [反教權的] "anti-clerical" făn fàxīsī
- 反 法西斯 [反法西斯] "anti-fascist"

Wikipedia







"language" in ASL lifeprint.com

	indicative				
	singular	plural			
	ich fege	wir fegen			
present	du fegst	ihr fegt			
	er fegt	sie fegen			
	ich fegte	wir fegten			
preterite	du fegtest	ihr fegtet			
	er fegte	sie fegten			
i	feg (du)	foot (ib.)			
imperative	fege (du)	fegt (ihr)			

conjugation of German "fegen" (to sweep) (Wiktionary)

ajar ilah = to teach (jussive, locative)
ajarkan = to teach (imperative, causative/applicative)
ajarkanlah = to teach (jussive, causative/applicative)
ajar lah = to teach (jussive, active)
<i>ajaran</i> = teachings
belajar = to learn (intransitive, active)
diajar = to be taught (intransitive)
diajari = to be taught (transitive, locative)
diajarkan = to be taught (transitive, causative/applicative
dipelajari = to be studied (locative)
dipelajarkan = to be studied (causative/applicative)
<i>mempelajari</i> = to study (locative)

ajar = to teach

ajari = to teach (imperative, locative)

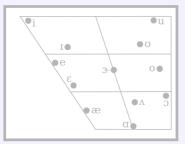
Indonesian "ajar" (to teach) with affixes (Wikipedia)

Words



$$\begin{split} & B_0 \ L & \text{i. } & X_1 X_1 \dot{a} \\ & B_0 \ L & \text{i. } & X_2 X_1 \dot{a} \\ & B_1 \ B_2 & \text{i. } & X_1 X_1 \dot{a} \\ & B_2 \ L & \text{i. } & X_2 X_1 \dot{a} \\ & B_3 \ L & \text{i. } & X_2 X_2 \dot{a} \\ & B_4 \ L & \text{i. } & X_3 X_1 \dot{a} \\ & B_4 \ L & \text{i. } & X_2 X_2 \dot{a} \\ & B_5 \ L & \text{i. } & X_3 X_1 \dot{a} \\ & B_6 \ L & \text{i. } & X_1 X_2 \dot{a} \\ & B_7 \ L & \text{i. } & X_2 X_2 \dot{a} \\ & B_8 \ L & \text{i. } & X_1$$

Stokoe notation of ASL



GA English vowels

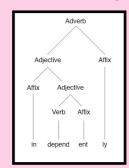
Chinese stroke primitives, from Jiang et al. 2024

Phonemes/ Graphemes/





ASL mouth morphemes Large (CHA) and Small (OO) learnhowtosign.com



Wikipedia (Annie Yang)

făn 反 — "anti-"

- făn kŏng ● 反 恐 [反恐] — "anti-terror"
- fǎn jiàoquánde • 反 教 权 的 [反教權的] — "anti-clerical"
- fǎn fàxīsī • 反 法西斯 [反法西斯] — "anti-fascist"

Wikipedia

Morphemes / /Lexemes /



"language" in ASL <u>lifeprint.com</u>

	indicative				
	singular	plural			
	ich fege	wir fegen			
present	du fegst	ihr fegt			
	er fegt	sie fegen			
	ich fegte	wir fegten			
preterite	du fegtest	ihr fegtet			
	er fegte	sie fegten			
	feg (du)	f+ (5h-)			
imperative	fege (du)	fegt (ihr)			

conjugation of German "fegen" (to sweep) (Wiktionary)

aiar = to teach

ajari = to teach (imperative, locative)

aja**riian** = to teach (jussive, locative)

ajarkan = to teach (imperative, causative/applicative ajarkanlah = to teach (jussive, causative/applicative)

ajar**lah** = to teach (jussive, active

ajaran = teachings

belajar = to learn (intransitive, active)

diajar = to be taught (intransitive)

diajari = to be taught (transitive, locative)

diajarkan = to be taught (transitive, causative/applicati

dinalajari - to be studied (locative)

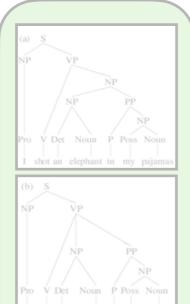
dipelajarkan = to be studied (causative/applicative)

mempelajari = to study (locative)

Indonesian "ajar" (to teach)

with affixes (Wikipedia)

Words



syntactic ambiguity, from UBC CPSC522



Universal Dependencies

Constituents
/Phrases



靜靜養貓 (YouTube)



Portal

Юная мышка Роза допела любимую песенку и теперь вслушивалась, как талло эко последних звуков. Роза была довольна собой наконец-то она нашла мего, откуда море отлично просматривалось. Вечерело, лучи солица, заходищего за ее спиной, окрасили воду « вазотителій и батровый цвета.

 Эгей, Роза, пожалуйте-ка ужинать. Не для того я старался, чтобы все это теперь простыло, значится... то есть отсырело. Нет уж, дудки!

Спутник Розы, крот Грумм, призывно помахал е дапой. Мышка подошла к маленькому костерку, на ко тором он стряпал, и принюхалась:

 Ого, лененики из дикого овса и суп из зелений Грумм улыбнулся, отчего его бархатная мордочки сморщилась, и помахал маленькой поварешкой, кото рую всегда носка за поясом наподобие меча;

Роза вядла глубокую раковину морского гресешка, наполненную ароматным суюм. Положный сюю леленну на плоский камень у костра, чтобы та не остыла, мышка отклебнула супа и покачала головой: — Ты хуже старой няньки, Грума Кланякинс Быось об акклад, если бы я позволила, ты бы меня и спать

Перед носом Розы замаячила поварешка.
— Так что же тебе, спать, это самос... и не надобно вовсе, так, что ли? Представь, что батюшка твой скажет, если я тебя домой доставлю всю усталую такую За потоятную э?

Utterances/ Sentences

B B C NEWS ਪੰਜਾਬੀ

ਾਂ ਚੀੜੀਓ ਪਾਠਕਾਂ ਦੀ ਪਸੰਦ ਚਾਰਜ ਕੌਮਾਂਸ਼ਰੀ

ਕੇਰਲ ਵਿੱਚ ਰਸਾਇਣਾਂ ਨਾਲ ਭਰਿਆ ਜਹਾਜ਼ ਡੁੱਬਿਆ, ਅਰਬ ਸਾਗਰ 'ਚ ਜੇ ਤੇਲ ਰਿਸਿਆਂ ਤਾਂ ਕੀ ਹੋਵੇਗਾ ਅਸਰ?

. . .

ਥੋਂ ਰਵਾਨਾ ਹੋਇਆ ਐੱਮਐੱਸਸੀ ਈਐੱਲਐੱਸਏ 3, 24 ਮਈ ਨੂੰ ਕੋਚੀ ਬੇਦਰਗਾਹ 'ਤੇ ਪਹੁੰਚਣ ਵਾਲਾ ਸੀ। (24 ਮਈ ਨੂੰ ਸਵੇਰੇ 12-15 ਵਜੇ ਭਾਰਤੀ ਤੱਟ ਰੱਖਿਅਕ ਨੂੰ ਜਹਾਜ਼ ਤੋਂ ਇੱਕ ਐਮਰਜੈਂਸੀ ਕਾਲ ਆਈ।

10 ਕੈਟੇਨਰਾਂ ਨੂੰ ਲੈ ਕੇ ਕੋਢੀ ਜਾਂਦੇ ਸਮੇਂ 184 ਮੀਟਰ ਲੰਬਾ ਐੱਮਐੱਸਸੀ ਈਐੱਲਐੱਸਏ 3 ਤੁੱਬਣ ਲੱਗਿਆ। ਦੋ ਕਾਰਗੇ ਜਹਾਜ਼ ਕੋਢੀ ਤੋਂ 38 ਸਮੁੰਦਰੀ ਮੀਲ ਦੱਖਣ-ਪੱਛਮ ਵਿੱਚ ਸੀ ਤਾਂ ਇਹ ਤਕਰਬੀਨ 26 ਡਿਕਰੀ ਤੱਕ ਤੋਂ ਵਿਸ਼ਾਹ ਸੀ।

ਾਰਤੀ ਤੋਂਟ ਕੈਖਿਅਕ ਨੇ ਤੁਰੰਤ ਨੇਡਲੇ ਜਹਾਜ਼ਾਂ ਨੂੰ ਬਚਾਅ ਕਾਰਜਾਂ ਲਈ ਭੇਜ ਦਿੱਤਾ ਸੀ। ਹਾਲਾਤ ਦੀ ਗਰਾਨੀ ਲਈ ਇੱਕ ਹਵਾਈ ਜਹਾਜ਼ ਵੀ ਮੌਜੂਦ ਸੀ। ਇਸ ਦੌਰਾਨ ਕਾਰਗੋ ਜਹਾਜ਼ ਲਗਾਰਾਰ ਝੁਕਦਾ ਗਿਆ ਸੇ ਕੁਝ ਤੱਕੇ ਸਮੇਦਰ ਵਿੱਚ ਇੱਗਣ ਲੱਗੇ।

ਭਾਰਤੀ ਜਲ ਲੈਨਾ ਨੇ 24 ਮਈ ਦੀ ਲਾਮ ਨੂੰ ਬਚਾਅ ਕਾਰਜ ਸ਼ੁਰੂ ਕੀਤੇ। ਦੇ ਜਹਾਜ਼, ਆਈਐੱਨਐੱਸ ਸ਼-ਅਤੇ ਆਈਐੱਨਐੱਸ ਸੁਜਾਰਾ, ਨੂੰ ਜਹਾਜ਼ ਵਿੱਚ ਸਵਾਰ 24 ਲੋਕਾਂ ਨੂੰ ਬਚਾਉਣ ਲਈ ਤੇਜਿਆ ਗਿਆ। ਆਈਐੱਨਐੱਸ ਸੁਜਾਰਾ ਸ਼ਾਮ 7 ਵਜੇ ਪਹੁੰਚਿਆ, ਜਦੋਂ ਕਿ ਆਈਐੱਨਐੱਸ ਸਤਪੁਰਾ ਰਾਤ 8 ਵਜੇ ਪਹੁੰਚ ਸੀਲਿਆ।

ਵੱਖਣ-ਪੱਛਮੀ ਮਾਨਸੂਨ, ਜੋ ਆਮ ਤੌਰ 'ਤੇ 1 ਜੂਨ ਨੂੰ ਸ਼ੁਰੂ ਹੁੰਦਾ ਹੈ, ਇਸ ਸਾਲ 24 ਮਈ ਨੂੰ ਸ਼ੁਰੂ ਹੋਇਆ ਸੰ ਇਸ ਲਈ ਸਮੁੰਦਰ ਦਾ ਮੌਸਮ ਖਰਾਬ ਸੀ।

ਆਈਐੱਟਿਐੱਸ ਸੁਜਾਰਾ 'ਦੇ ਕੈਪਟਨ ਅਰਜੂਨ ਸ਼ੇਖਰ ਨੇ ਖ਼ਬਰ ਦੇਸ਼ਸੀ ਦੋਐੱਨਆਈ ਨੂੰ ਦੱਸਿਆ,' ਸਾਨੂੰ ਪ੍ਰਤੀਯੂਲ ਹਾਲਾਤ ਦਾ ਸਾਹਮਣਾ ਕਰਨਾ ਪਿਆ। ਹਵਾ 74.08 ਕਿਲੋਮੀਟਰ ਪ੍ਰਤੀ ਘੰਟਾ (40 ਨਾਟ) ਦੀ ਰਫ਼ਟ ਨਾਲ ਵਾਗੇ ਹਹੀ ਸੀ। ਸਮੁੰਦਰ ਵਿੱਚ ਗੂੜਾ ਅਤੇ ਕੰਟੇਨਰ ਤੇਰ ਰਹੇ ਸਨ। ਇਸ ਕਾਰਨ ਰਾਤ ਨੂੰ ਜਹਾਜ਼ ਤੱਕ ਪਹੁੰਚਣਾ ਮੁਸ਼ਕਰ ਹੋ ਰਿਹਾ ਸੀ।"

ਾਜ਼ ਵਿੱਚ ਸਵਾਧਾ 2+ ਲੰਕਾਂ ਵਿੱਚੋਂ, 21 ਨੂੰ ਉਸ ਰਾਡ ਜ਼ਿਨ੍ਹਾਂ - ਜਿਸ਼ੀ ਜਾਨਲੰਦਾ ਸੰਦ ਦੇ ਬਾਚਾ ਲਿਆ ਲਿਆ। ਉਸ ਰਾਜਾ 'ਤੇ ਅਜੇ ਦੀ ਕੈਂਟੋਨਰ ਸਨ ਅਤੇ ਜਹਾਜ਼ ਪੂਰੀ ਭਰੁ ਾਂ ਨਹੀਂ ਕੁੱਬਿਆ ਸੀ, ਇਸ ਲਈ ਜਹਾਜ਼ ਮਾਸਟਰ, ਮੁੱਖ ਇੰਜੀਨੀਅਰ ਅਤੇ ਸਹਾਇਕ ਇੰਜੀਨੀਅਰ ਬਚਾਅ ਕਾਰਜਾਂ ਨੂੰ ਨੋਪਰੇ ਚਾਡਨ ਅਤੇ ਸਥਿਤੀ ਦੀ ਭਰਾਨੀ ਕਰਨ ਲਈ ਜਹਾਜ਼ 'ਤੇ ਹੀ ਰਹੇ।

ਤਿੰਨਾਂ ਨੇ ਭਾਰਤੀ ਤੱਟ ਰੱਖਿਅਕ ਅਤੇ ਭਾਰਤੀ ਜਲ ਸੈਨਾ ਦੀ ਨਿਗਰਾਨੀ ਹੇਠ ਜਹਾਜ਼ 'ਚ ਹੀ ਰਾਤ ਬਿਤਾਏ

BBC News in Punjabi



Portal 2 dialogues

CMAIR EID: Thank you. Does the Committee
have any questions? I do not see any. Thank you so
much for your presentation today.
Ms. RESNIK: We appreciate your time and
that you enabled us all to offer comments. Many
thanks.

CMAIR EID: Thank you.

All right. We are now going to turn to
Carter Phillips, and we have now moved to Rule 29
comments.

MS. FMILLIPS: Judge Eid, can you see me and

US Supreme Court

Discourse/ Dialogue



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

- Text data can be viewed as a sequence of words
- First step in building a language technology: building a function that maps from arbitrary text data to that sequence

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

- Type-token distinction:
 - Type: a unique word in a text corpus
 - Token: an instance of a word type, appearing in a particular context

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

```
{'.', ',', ':', 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic',
'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form',
'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening',
'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches',
'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young'}
```



wordtypes (vocabulary)

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

```
{'.', ',', ':', 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic',
'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form',
'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening',
'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches',
'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young'}
```

wordtypes (vocabulary)

instances (tokens) of wordtype \'.'

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

```
{'.', ',', ':', 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic',
'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form',
'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening',
'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches',
'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young'}
```

wordtypes (vocabulary)

instances (tokens) of wordtype ','

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

```
{'.', ',', ':', 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic',
'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form',
'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening',
'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches',
'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young'}
```

wordtypes (vocabulary)

instances (tokens) of wordtype 'Every'

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

```
{'.', ',', ':', 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic',
'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form',
'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening',
'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches',
'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young'}
```

wordtypes (vocabulary)

instances (tokens) of wordtype 'and'

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

```
{'.', ',', ':', 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic',
'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form',
'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening',
'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches',
'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young'}
```

wordtypes (vocabulary)

instances (tokens) of wordtype 'language'

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

```
vocabulary =
['.', ',', ':', 'Every', 'The', 'a', 'adults', 'and', 'are', 'basic',
  'can', 'children', 'conversation', 'dialogue', 'even', 'far', 'form',
  'from', 'hold', 'illiterate', 'including', 'is', 'language', 'listening',
  'most', 'natural', 'of', 'preparing', 'reading', 'skills', 'speeches',
  'to', 'universal', 'use', 'user', 'whereas', 'writing', 'young']
```

vocabulary as a look-up table

```
tokenized_text = ['The', 'most', ..., 'skills, '.']
tokenized_indices = [vocabulary.index(token) for token in tokenized_text]

tokenized_indices: [4, 24, 25, 7, 9, 16, 26, 22, 33, 21, 13, 2, 3, 22,
34, 1, 20, 37, 11, 7, 19, 6, 1, 10, 18, 5, 12, 1, 35, 28, 1, 36, 1, 27,
30, 7, 14, 23, 31, 30, 8, 15, 17, 32, 29, 0]
```

text as a sequence of wordtype indices

Tokenize on Spaces



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

Simplest tokenizer (for English): splitting on spaces

```
tokenized = s.split(' ')
```

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue:', 'Every', 'language', 'user,', 'including', 'young', 'children', 'and', 'illiterate', 'adults,', 'can', 'hold', 'a', 'conversation,', 'whereas', 'reading,', 'writing,', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills.']
```

But this gets us some weird wordtypes:

'dialogue:'
'user,'
'skills.'



Not really words different from dialogue, user, skills

Rule-Based Tokenization (



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

nltk tokenizers, with special rules for punctuation

```
import nltk
tokenized = nltk.word_tokenize(s)
```

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```

But this still loses similarity between wordtypes

'skills'
'reading'
'speeches'



Lexically similar to, but morphologically distinct from skill, read, speech

Rule-Based Tokenization



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills. from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

nltk tokenizers, with special rules for punctuation

```
import nltk
tokenized = nltk.word tokenize(s)
```

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language',
'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',',
'including', 'young', 'children', 'and', 'illiterate', 'adults', ',',
'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',',
'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening',
'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```

And doesn't work for all languages

โดยที่การยอมรับนับถือเกียรติศักดิ์ประจำตัว และสิทธิเท่าเทียมกันและโอนมิได้ของ บรรดา สมาชิก ทั้ง หลายแห่งครอบครัว มนุษย์เป็นหลักมูลเหตุแห่งอิสรภาพ ความ ยุติธรรม และสันติภาพในโลก

Rule-Based Tokenization (



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

nltk tokenizers, with special rules for punctuation

```
import nltk
tokenized = nltk.word_tokenize(s)
```

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```

Once you've "trained" your tokenizer, you're stuck with it

```
vocab = ['.', ',', ':', 'Every', ..., 'writing', 'young']
vocab.index('ChatGPT') → not found!
```

Rule-Based Tokenization (



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

nltk tokenizers, with special rules for punctuation

```
import nltk
tokenized = nltk.word_tokenize(s)
```

```
['The', 'most', 'natural', 'and', 'basic', 'form', 'of', 'language', 'use', 'is', 'dialogue', ':', 'Every', 'language', 'user', ',', 'including', 'young', 'children', 'and', 'illiterate', 'adults', ',', 'can', 'hold', 'a', 'conversation', ',', 'whereas', 'reading', ',', 'writing', ',', 'preparing', 'speeches', 'and', 'even', 'listening', 'to', 'speeches', 'are', 'far', 'from', 'universal', 'skills, '.']
```

• Once you've "trained" your tokenizer, you're stuck with it



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

Strings are sequences of characters (bytes)!

```
tokenized = s.encode()
```

Now our vocabulary is a fixed size (all possible Unicode)

characters)

81	82	83	84	85	86	87	88	89	00	01	02	949	950	951	952	953	954	955	956	95	9	758	959	960	_	0		м	Dia	D.	144		at
01	02	03	04								_	_	_	_	_	٠	_								F	£	Ŋń	₩	Pts	Rs	₩	回	ā
]	^	_	`	а	b	С	d	е	f	g	h	ρ	ς	σ	τ	U	Φ	X	Ψ	ú	ט				355	8356	8357	8358	8350	8360	8361	8362	8363
93	94	95	96	97	98	99	100	101	102	103	104	961	962	963	964	965	966	967	968	96	9				555	6330	6337	0330	0334	8300	6301	8302	6303
i	j	k	ι	m	n	0	р	q	r	s	t	&	₽	Q	R 2	R B	R R	Ř	SM	TEL	тм	ÿ	#	٦,	Ďр	Ŋ	₱	Ф	A	₹	¢	tt	\$
105	106	107	108	109	110	111	112	113	114	115	116	8472	8473	8474	8475 84	176 84	77 8478	8479	8480	8481	8482	8483	836	56 5	3367	8368	8369	8370	8371	8272	8373	8374	8375
u	v	w	х	у	z	{		}	~		€	7	Z	Ω	Ω	2	K	Å	Œ	Œ	е	e	030	0	3307	0300	0304	0370	6371	03/2	63/3	63/4	63/5



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

Strings are sequences of characters (bytes)!

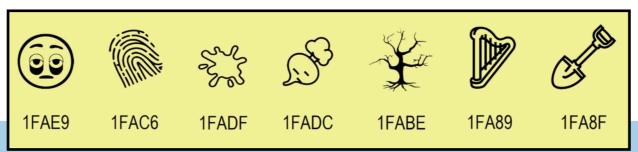
```
tokenized = s.encode()

54 68 65 20 6d 6f 73 74 20 6e 61 74 75 72 61 6c 20 61 6e 64 20 62 61 73 69 63 20 66 6f 72 6d 20 6f 66 20 6c 61 6e 67 75 61 67 65 20 75 73 65 20 69 73 20 64 69 61 6c 6f 67 75 65 3a 20 45 76 65 72 79 20 6c 61 6e 67 75 61 67 65 20 75 73 65 72 2c 20 69 6e 63 6c 75 64 69 6e 67 20 79 6f 75 6e 67 20 63 68 69 6c 64 72 65 6e 20 61 6e 64 20 69 6c 6c 69 74 65 72 61 74 65 20 61 64 75 6c 74 73 2c 20 63 61 6e 20 68 6f 6c 64 20 61 20 63 6f 6e 76 65 72 73 61 74 69 6f 6e 2c 20 77 68 65 72 65 61 73 20 72 65 61 64 69 6e 67 2c 20 77 72 69 74 69 6e 67 2c 20 70 72 65 70 61 72 69 6e 67 20 73 70 65 65 63 68 65 73 20 61 6e 64 20 65 76 65 6e 20 6c
```

69 73 74 65 6e 69 6e 67 20 74 6f 20 73 70 65 65 63 68 65 73 20 61 72 65 20 66 61 72 20 66 72 6f

 And we don't have any unknown wordtypes (until the Unicode Consortium adds new emoji)

6e 69 76 65 72 73 61 6c 20 73 6b 69 6c 6c 73 2e





The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

Strings are sequences of characters (bytes)!

```
tokenized = s.encode()
```

But: individual characters are not meaningful

```
['.', ',', ':', 'Every', ..., 'user', 'whereas', 'writing', 'young']
```

VS.



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

Strings are sequences of characters (bytes)!

```
tokenized = s.encode()

54 68 65 20 6d 6f 73 74 20 6e 61 74 75 72 61 6c 20 61 6e 64 20 62 61 73 69 63 20 66 6f 72 6d 20 6f 66 20 6c 61 6e 67 75 61 67 65 20 75 73 65 20 69 73 20 64 69 61 6c 6f 67 75 65 3a 20 45 76 65 72 79 20 6c 61 6e 67 75 61 67 65 20 75 73 65 72 2c 20 69 6e 63 6c 75 64 69 6e 67 20 79 6f 75 6e 67 20 63 68 69 6c 64 72 65 6e 20 61 6e 64 20 69 6c 6c 69 74 65 72 61 74 65 20 61 64 75 6c 74 73 2c 20 63 61 6e 20 68 6f 6c 64 20 61 20 63 6f 6e 76 65 72 73 61 74 69 6f 6e 2c 20 77 68 65 72 65 61 73 20 72 65 61 64 69 6e 67 2c 20 77 72 69 74 69 6e 67 2c 20 70 72 65 70 61 72 69 6e 67 20 73 70 65 65 63 68 65 73 20 61 6e 64 20 65 76 65 6e 20 6c 69 73 74 65 6e 69 6e 67 20 73 70 65 65 63 68 65 73 20 61 72 20 66 72 6f 6d 20 75 6e 69 76 65 72 73 61 6c 20 73 6b 69 6c 6c 73 2e
```

 It's now the ML model's job to learn to compose words from scratch



The most natural and basic form of language use is dialogue: Every language user, including young children and illiterate adults, can hold a conversation, whereas reading, writing, preparing speeches and even listening to speeches are far from universal skills.

from Pickering and Garrod 2004, "A mechanistic psychology of dialogue"

Strings are sequences of characters (bytes)!

```
tokenized = s.encode()
```

```
54 68 65 20 6d 6f 73 74 20 6e 61 74 75 72 61 6c 20 61 6e 64 20 62 61 73 69 63 20 66 6f 72 6d 20 6f 66 20 6c 61 6e 67 75 61 67 65 20 75 73 65 20 69 73 20 64 69 61 6c 6f 67 75 65 3a 20 45 76 65 72 79 20 6c 61 6e 67 75 61 67 65 20 75 73 65 72 2c 20 69 6e 63 6c 75 64 69 6e 67 20 79 6f 75 6e 67 20 63 68 69 6c 64 72 65 6e 20 61 6e 64 20 69 6c 6c 69 74 65 72 61 74 65 20 61 64 75 6c 74 73 2c 20 63 61 6e 20 68 6f 6c 64 20 61 62 64 20 63 66 67 2c 20 77 72 69 74 69 6e 67 2c 20 70 72 65 70 61 72 69 6e 67 2c 20 73 74 65 6e 69 74 6f 20 73 74 65 6c 74 73 2c 20 66 61 72 20 65 76 65 6e 20 66 69 73 74 65 6e 69 76 65 72 73 61 6c 20 73 70 65 65 63 68 65 73 20 61 6e 64 20 65 76 65 6e 20 6c 69 74 65 72 73 61 6c 20 73 70 65 65 73 20 61 72 65 20 66 61 72 20 66 72 6f 6d 20 75 6e 69 76 65 72 73 61 6c 20 73 6b 69 6c 6c 73 2e
```

But: input sequences are much longer

```
'language'
```

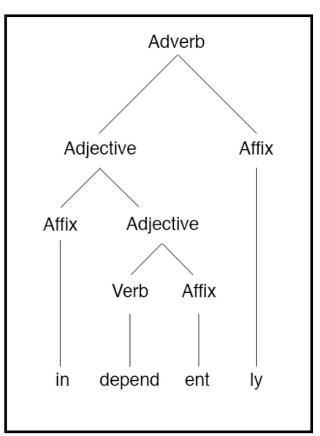
VS.

```
'l', 'a', 'n', 'g', 'u', 'a', 'g', 'e'
```

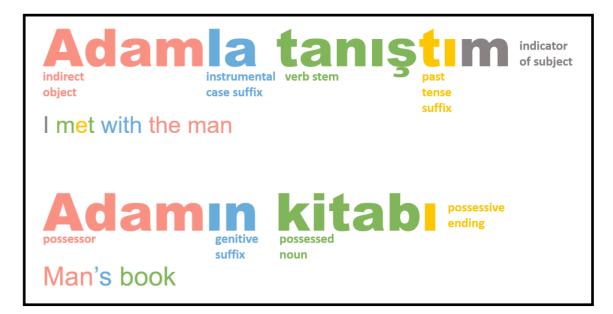
A Compromise: Subword Tokenization



- Main principle: words are (often) composed of subparts (morphemes)
- Our vocabulary should have entries for frequent words kept whole, because we have a lot of data about those words
- But it should also have entries for parts of less-frequent words, so our ML models can learn how to compose parts of words into whole words (especially unfamiliar words!)



Wikipedia (Annie Yang)





- Modern standard for building a tokenizer
- Inputs: collection of texts and target vocabulary size
- Initial vocabulary is the set of all bytes (characters) across the texts
- Until the target vocabulary size is reached, repeat the following:
 - Tokenize all of the texts using the current vocabulary
 - Find the most common bigram in the tokenized texts, then add it to the vocabulary as a new wordtype



```
Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary
```

```
('h', 'u', 'g', 'p', 'n', 'b', 's')
```



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)

bigrams + frequencies

'h' 'u' 15



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)

'h'	'u'	15
'u'	' g'	20



```
Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)
```

vocabulary tokenized texts

```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)
```

```
'h' 'u' 15
'u' 'q' 20
'p' 'u' 17
```



```
Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)
```

vocabulary tokenized texts

```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' <mark>'u' 'n'</mark>, 12), ('b' <mark>'u' 'n'</mark>, 4), ('h' 'u' 'g' 's', 5)
```

```
'h' '11'
         15
'u' 'q' 20
'p' 'u' 17
'u' 'n'
         16
```



```
Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)
```

vocabulary tokenized texts

```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)
```

```
'h' '11'
          15
'u' 'a'
         20
'p' 'u' 17
'u' 'n'
         16
'b' '11'
```



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)
```

```
'h' 'u'
         15
'u' 'a'
         20
'p' 'u' 17
'u' 'n' 16
'b' 'u'
'q' 's'
```



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)

bigrams + frequencies

'h'	'u'	15
'u'	'g'	20
'p'	'u'	17
'u'	'n'	16
'b'	'u'	4
'g'	's'	5

most frequent bigram; add to vocabulary



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)
```



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)
('h', 'u', 'g', 'p', 'n', 'b', 's', 'ug') - ('h' 'ug', 10), ('p' 'ug', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'ug' 's', 5)
```



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)
('h', 'u', 'g', 'p', 'n', 'b', 's', 'ug') ('h' 'ug', 10), ('p' 'ug', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'ug' 's', 5)
('h', 'u', 'g', 'p', 'n', 'b', 's', 'ug', 'un') ('h' 'ug', 10), ('p' 'ug', 5), ('p' 'un', 12), ('b' 'un', 4), ('h' 'ug' 's', 5)
```



Documents + frequencies: ('hug', 10), ('pug', 5), ('pun', 12), ('bun', 4), ('hugs', 5)

vocabulary tokenized texts

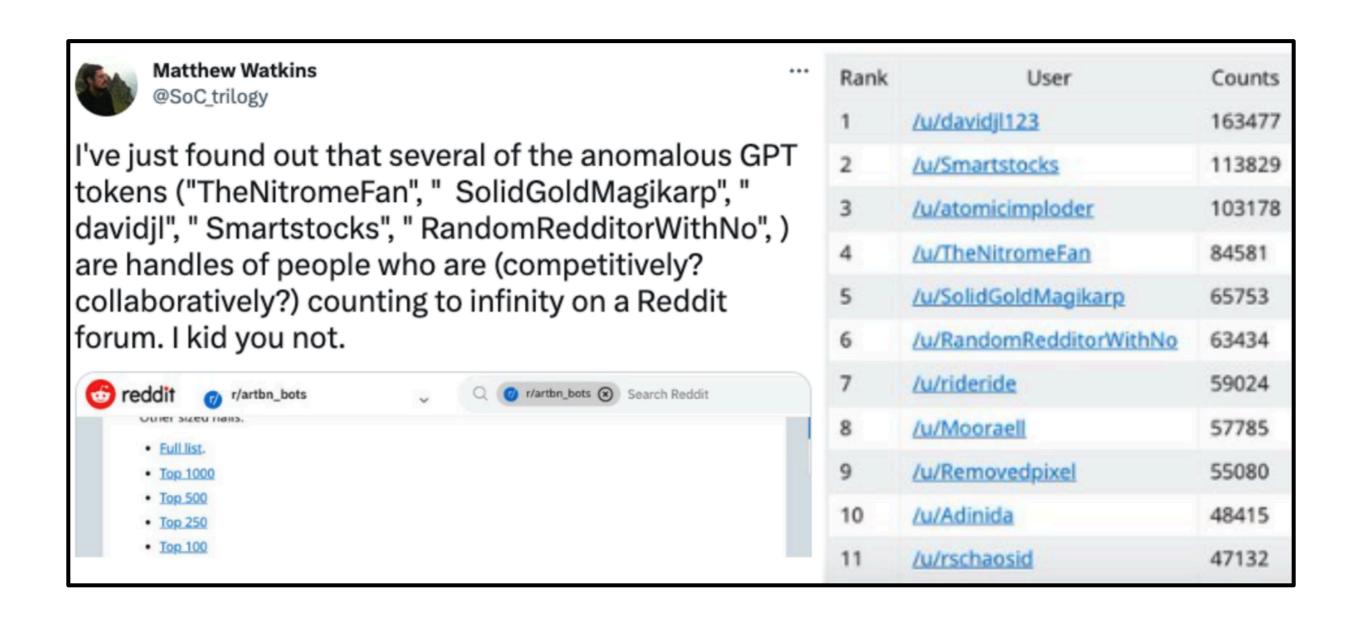
```
('h', 'u', 'g', 'p', 'n', 'b', 's') ('h' 'u' 'g', 10), ('p' 'u' 'g', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'u' 'g' 's', 5)
('h', 'u', 'g', 'p', 'n', 'b', 's', 'ug') ('h' 'ug', 10), ('p' 'ug', 5), ('p' 'u' 'n', 12), ('b' 'u' 'n', 4), ('h' 'ug' 's', 5)
('h', 'u', 'g', 'p', 'n', 'b', 's', 'ug', 'un') ('h' 'ug', 10), ('p' 'ug', 5), ('p' 'un', 12), ('b' 'un', 4), ('h' 'ug' 's', 5)
('h', 'u', 'g', 'p', 'n', 'b', 's', 'ug', 'un', 'hug') ('hug', 10), ('p' 'ug', 5), ('p' 'un', 12), ('b' 'un', 4), ('hug' 's', 5)
```

New word: "puns"



Vocabularies



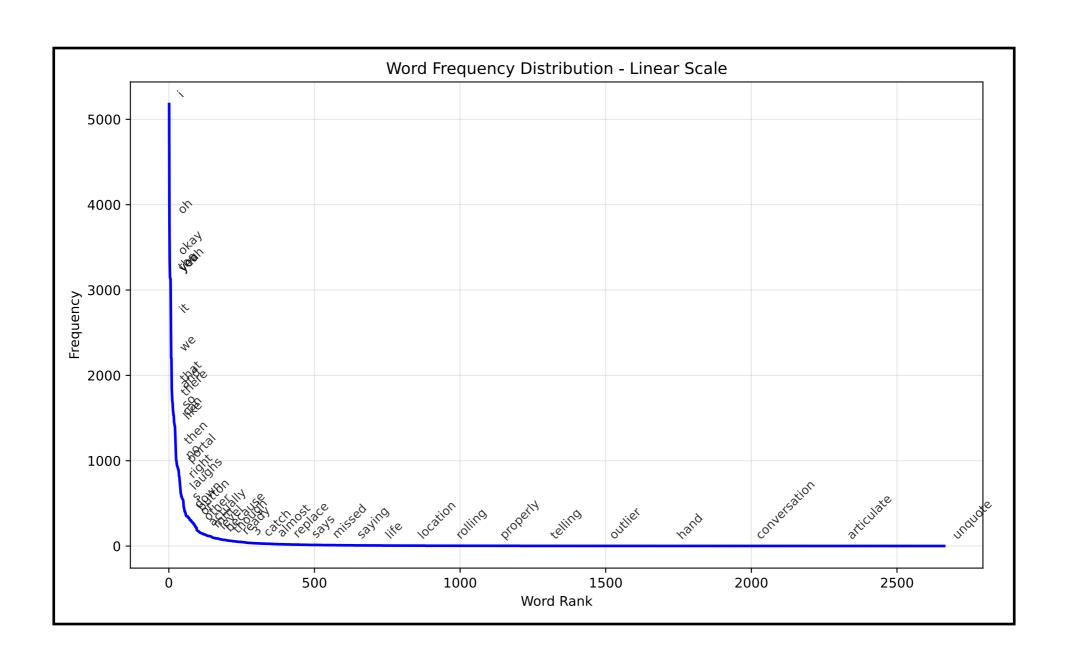


Vocabularies



Long-tail ("Zipfian") distribution of wordtypes

(from Portal 2 dialogues)

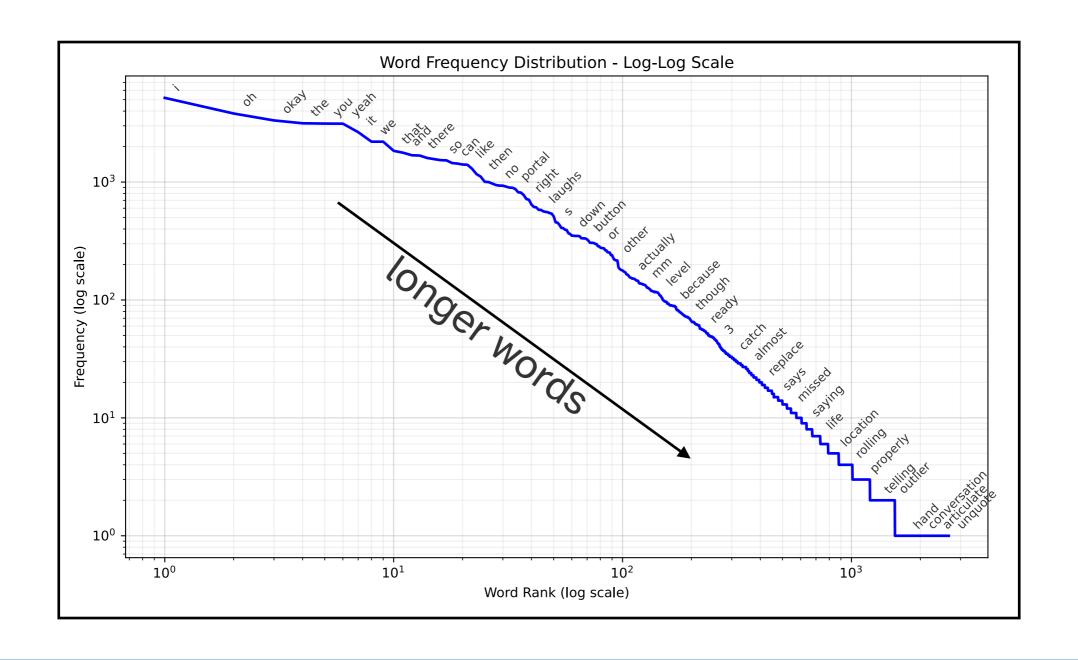


Vocabularies



Long-tail ("Zipfian") distribution of wordtypes

(from Portal 2 dialogues)



Meaning and Representation (



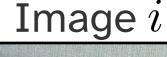
- Words are themselves not actual things, they only refer to things
 - Kind of like a pointer
 - Even if they aren't actual things, we can still do things to and know things about them
 - E.g., we know cats are a type of pet
- Core question: if a machine learning system is processing text, how should words be <u>represented</u> as input to and within the ML system?

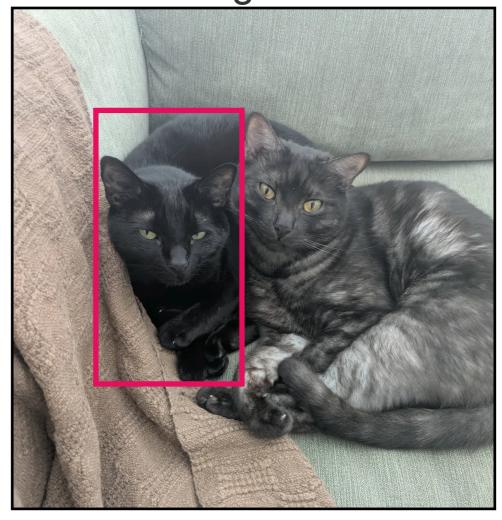
Word Meaning: Denotational Semantics



What color is Pepper?

- Let's start from the bottom up: what does each word mean?
 - What does "Pepper" mean?
- **Denotational semantics:** the symbol *refers to* something in the context in which the language is used





$$[\![\text{Pepper}]\!]^i = \{$$

Word Meaning: Denotational Semantics



What color is

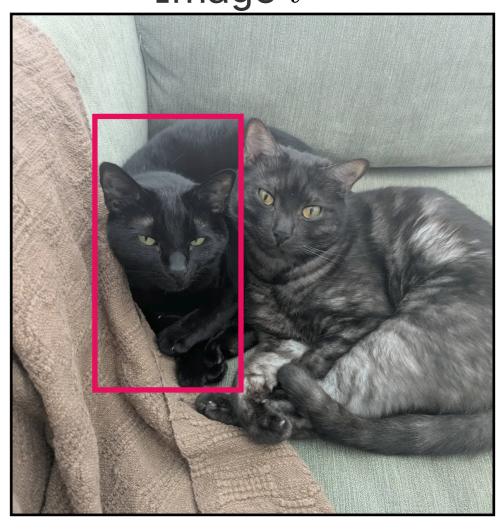
?

Black with green eyes

- Denotational semantics allows us to ground language to a world outside of language use
- But what happens when we don't have an available outside world?











Cats



Nick's Cats



Téa's Cats



Zineng's Cats



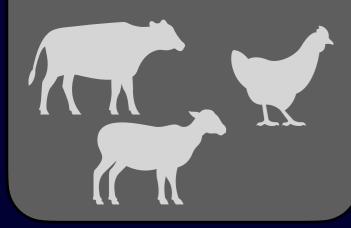
Alane's Cats Patrick



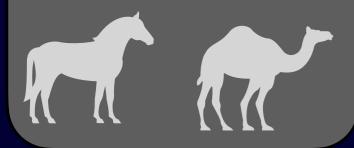


Domesticated Animals

Livestock



Draught Animals















Pets





Domesticated Animals

Hyponym / hypernym relation:

"pet" is a <u>hyper</u>nym of "dog";

"draught animal" is a <u>hypo</u>nym of

"domesticated animals"



Draught Animals



Pets

Dogs Mer's Dogs
Synsets:

{"draught animal", "beast of burden", "draft animal"}



What color is Pepper?

- From our ontology, we know
 - isa(Pepper, Alane's cat)
- ← (hypernym relation)

- isa(Alane's cat, cat)
- We might know something about cats, too:
 - Cats are typically not blue, pink, yellow, etc.
 - So the answer is probably not one of these



What color is pepper?



Missing context-dependence (polysemy)



- Missing context-dependence (polysemy)
- Missing meaning of new words

In cryptography, a **pepper** is a secret added to an input such as a password during hashing with a cryptographic hash function. This value differs from a salt in that it is not stored alongside a password hash, but rather the pepper is kept separate in some other medium, such as a Hardware Security Module.^[1] Note that the National Institute of Standards and Technology refers to this value as a **secret key** rather than a **pepper**. A pepper is similar in concept to a salt or an encryption key. It is like a salt in that it is a randomized value that is added to a seword hash, and it is similar to an encryption key in that it should be kept secret.

Webster, Craig (2009-08-03). "Securing Passwords with Salt, Pepper and Rainbows" ∠. Barking Iguana. Retrieved 2020-11-11.





Paullassiter, Wikipedia

- Missing context-dependence (polysemy)
- Missing meaning of new words
- Requires human labor

There are perhaps fifty thousand Capsicum cultivars grown worldwide.[1]

- Capsicum annuum, which includes bell peppers, cayennes, friggitello, jalapeños, paprika, and serrano.
 - Capsicum annuum 'New Mexico Group', common name Hatch or Anaheim, which includes Big Jim, Chimayó, and Sandia peppers.
- Capsicum baccatum, which includes the South American varieties, such as ají amarillo, ají limón, and criolla sella.
- *Capsicum chinense*, which includes all of the habaneros, [5] Scotch bonnets, Trinidad Scorpions, the Bhut Jolokia, and the Carolina Reaper.
- Capsicum frutescens, which includes the Tabasco pepper and many of the peppers grown in India;^[6] sometimes not distinguished as a species separate from C. annuum.^{[7][8]}
- Capsicum pubescens, which includes the rocoto and manzano pepper, are distinctive plants, having violet flowers, black seeds, and hairy dark green leaves, and grow as a large, multi-stemmed vine up to 5 meters long.



Falcodigiada, Wikipedia



Hankwang, Wikipedia



Daniel Risacher, Wikipedia



Endwints, Wikipedia



- Missing context-dependence (polysemy)
- Missing meaning of new words
- Requires human labor
- Subjective and culturally dependent











Word Meaning: Discrete Symbols



- Machine learning models expect numerical inputs
- Represent each word as a unique one-hot vector (vector with values 0, except for one value of 1)

$$\phi(cat) = [0 0 0 0 1 0 0 0 0]$$

 $\phi(kitten) = [0 0 1 0 0 0 0 0]$

- Vector dimension: the same size as the vocabulary
- Problem: no notion of similarity

How can I make my house safe for a new kitten?

How can I make my house safe for a new cat?

 (Is there some way we could use WordNet to get more informative numerical representations of words?)

Word Meaning: Continuous Vectors



Instead: represent words as continuous vectors

$$\phi$$
(cat) = [-0.6 1.3 -0.1 0.3]
 ϕ (kitten) = [-0.7 -1.5 0.0 0.3]

Implicitly provides notions of similarity:

$$||\phi(\text{cat}), \phi(\text{kitten})|| < ||\phi(\text{cat}), \phi(\text{dog})||$$

Similarity is learnable from text at scale:

```
... feeding time in multiple cat households can often be ...
... is to relieve the cat from the stress of ...
... effect on the feral cat population in rural areas ...
```

```
... So I wrapped the kitten up in a towel ...
... they noticed that the kitten population was getting high ...
... were looking for a kitten , hoping that someone ...
```

Word Embeddings



- How do we get these vectors?
- Core principle: distributional hypothesis
 - Words that are used in similar contexts have similar meanings
 - Context: typically, other words in a text, but really anything can be context!

Distributional Hypothesis



A bottle of tesgüino is on the table.

Everybody likes tesgüino.

Tesgüino makes you drunk.

We make tesgüino out of corn.

- Found in bottles
- People like it
- It makes you drunk
- It's made of corn
- Maybe: similar to beer, wine, whiskey, etc.

Word Co-Occurrence



- Co-occurrence matrix, counting number of texts in which both words occur
- Word similarity ≈ cosine similarity of vector representations φ
- But embedding is huge: the size of the vocabulary!

	tesgüino	beer	wine	whiskey
drunk	3	400	450	600
bottle	10	600	1000	600
corn	15	0	0	400
hops	0	450	0	0
grapes	0	0	1000	0
dairy	0	0	0	0



- Input: corpus D including pairs (w, c) where w is a word and c is a context,
 e.g.:
 - "Everybody likes tesgüino" (w = "tesgüino", c = "likes") e.g., context is every (w = "tesgüino", c = "everybody") word in the sentence
- We want to model: given that we observe word w, what's likely in the context c?

$$p(c \mid w; \theta)$$

Our objective: find parameters that maximize the corpus probability

$$\arg\max_{\theta} \prod_{(w,c)\in\mathcal{D}} p(c \mid w; \theta)$$

• (There are other methods for word2vec (e.g., GloVe, CBOW), but we will only look at Skip-Gram in this class)



Objective: find

$$\arg\max_{\theta} \prod_{(w,c)\in\mathcal{D}} p(c \mid w; \theta)$$

where

$$p(c \mid w; \theta) = \frac{\exp(\text{score}(c, w))}{\sum_{c' \in \mathcal{C}} \exp(\text{score}(c', w))}$$

score: similarity of we use cosine similarity because, representations of in contrast to dot product, it is word and context unbiased towards vector

magnitude:
$$score(c, w) = cos(\phi(c), \phi(w)) = \frac{\phi(c) \cdot \phi(w)}{||\phi(c)|| ||\phi(w)||}$$



Objective: find

$$\arg\max_{\theta} \prod_{(w,c)\in\mathcal{D}} p(c \mid w; \theta)$$

$$p(c \mid w; \theta) = \frac{\exp(\operatorname{score}(c, w))}{\sum_{c' \in \mathcal{C}} \exp(\operatorname{score}(c', w))}$$

softmax function:
normalizes categorical
scores to a smooth
probability distribution



Objective: find

$$\arg\max_{\theta} \prod_{(w,c)\in\mathcal{D}} \frac{\exp(\operatorname{score}(c,w))}{\sum_{c'\in\mathcal{C}} \exp(\operatorname{score}(c',w))}$$

Problem: intractable to sum over all contexts



- * Solution: negative sampling
 - * Instead of summing over all possibly contexts in which a word could appear
 - * Approximate via learning from contexts in which a word doesn't appear
- * Model whether or not a word-context pair is likely to exist in the dataset

$$p((w,c) \in \mathcal{D} = \frac{1}{1 + \exp(-\operatorname{score}(c,w))}$$

similarity of representations of word and context



- * Solution: negative sampling
 - * Instead of summing over all possibly contexts in which a word could appear
 - * Approximate via learning from contexts in which a word doesn't appear
- * Model whether or not a word-context pair is likely to exist in the dataset

$$p((w,c) \in \mathcal{D} = \frac{1}{1 + \exp(-\operatorname{score}(c,w))}$$

sigmoid function: normalizes a scalar into a probability



NEW objective: find

$$\arg\max_{\theta} \prod_{(w,c)\in\mathcal{D}} p((w,c)\in\mathcal{D}) \prod_{(w,c)\in\mathcal{D'}} p((w,c)\not\in\mathcal{D})$$

where

$$p((w,c) \in \mathcal{D} = \frac{1}{1 + \exp(-\operatorname{score}(c,w))}$$
$$p((w,c) \notin \mathcal{D}) = 1 - p((w,c) \in \mathcal{D})$$

What is D?



NEW objective: find

$$\arg\max_{\theta} \prod_{(w,c)\in\mathcal{D}} p((w,c)\in\mathcal{D}) \prod_{(w,c)\in\mathcal{D}'} p((w,c)\not\in\mathcal{D})$$

* Negative samples: sample and train on l * |D| pairs (w', c') where

$$w' \sim p(W)$$
 $c' \sim p(C)$ unigram prior over words over contexts

- * How to train?
 - * Gradient descent
 - * In practice, we work with log probabilities, not direct probabilities, to avoid float underflow



What counts as context? Basically anything you want

A bottle of tesgüino is on the table.

Everybody likes tesgüino.

Tesgüino makes you drunk.

We make tesgüino out of corn.

	w = tesgüino
of w	1
bottle of w	1
bottle w	1
cup w	0
w is	1
w is on	1
w is on the table	1
w table	1
{corn, w}	1
make w out	1

Word Vectors



* In Skip-Gram, the parameters of our models are the vectors for words and the vectors for contexts

$$p((w,c) \in \mathcal{D} = \frac{1}{1 + \exp(-\cos(\phi(c), \phi(w)))}$$

context vectors word vectors

* Once we've trained a model, we can use these vectors for whatever we'd like!



a:b::c:?
$$\equiv d = \arg\max_{i \in \mathcal{W}} \frac{(\phi(b) - \phi(a) + \phi(c)) \cdot \phi(i)}{||(\phi(b) - \phi(a) + \phi(c)) \cdot \phi(i)||}$$

vector representing the relationship between a and b



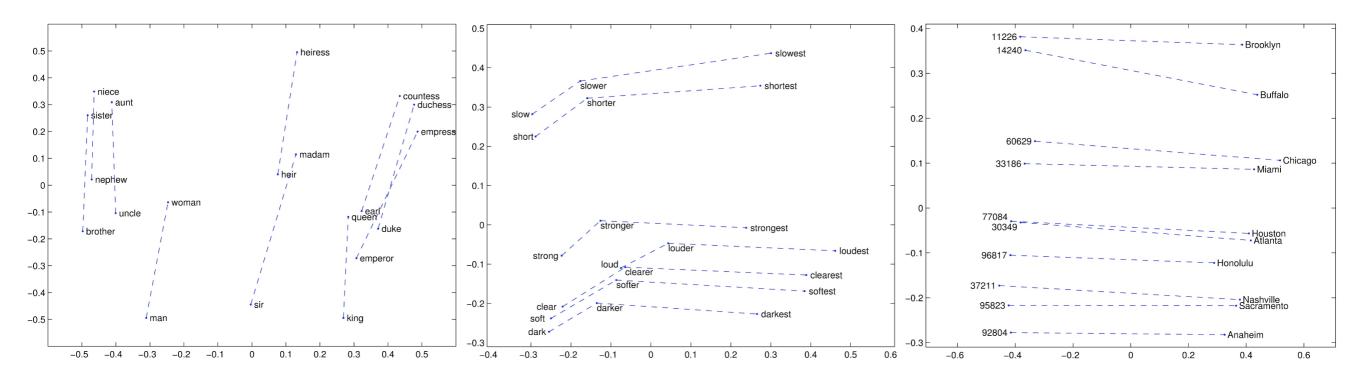
a:b::c:?
$$\equiv d = \arg\max_{i \in \mathcal{W}} \frac{(\phi(b) - \phi(a) + \phi(c)) \cdot \phi(i)}{||(\phi(b) - \phi(a) + \phi(c)) \cdot \phi(i)||}$$

apply this relationship to a new word c



a:b::c:?
$$\equiv d = \arg\max_{i \in \mathcal{W}} \frac{\left(\phi(b) - \phi(a) + \phi(c)\right) \cdot \phi(i)}{\left|\left|\phi(b) - \phi(a) + \phi(c)\right| \cdot \phi(i)\right|\right|}$$

cosine the expected and some similarity of representation of the other word analogy word i





Pairs (x, y) that maximize $\cos(\phi(\text{he}) - \phi(\text{she}), \phi(x) - \phi(y))$

Gender stereotype she-he analogies.

sewing-carpentry
nurse-surgeon
blond-burly
giggle-chuckle
sassy-snappy
volleyball-football

register-nurse-physician interior designer-architect feminism-conservatism vocalist-guitarist

diva-superstar cupcakes-pizzas housewife-shopkeeper softball-baseball

cosmetics-pharmaceuticals

petite-lanky

 $\begin{array}{c} charming\text{-}affable\\ hairdresser\text{-}barber \end{array}$

Gender appropriate she-he analogies.

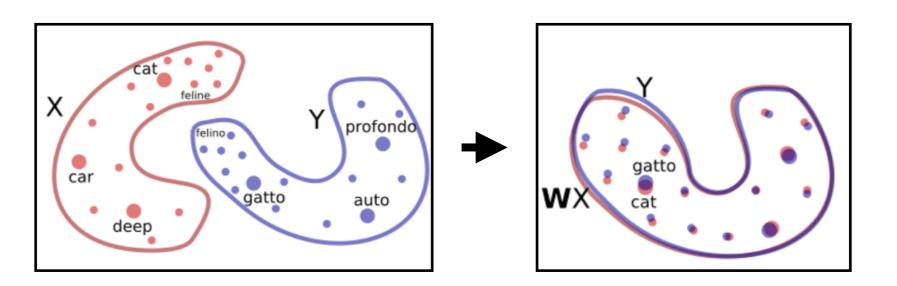
queen-king waitress-waiter sister-brother ovarian cancer-prostate cancer

mother-father convent-monastery

- Stereotypes are embedded in text data
- One way to "debias": identify category subspace (e.g., "gender" subspace), project words onto that subspace, then subtract those projections from the original word



- What if we have embeddings in different languages and want to align them? Given:
 - ullet ϕ^A : language A's embedding function
 - ϕ^B : language B's embedding function
 - $\mathcal{D} = \{w_i^A, w_i^B\}_{i=1}^M$: dataset pairing \mathbf{M} word translations
- Find a matrix W such that $\min_{W} \sum_{i=1}^{\infty} ||W\phi(w_i^A) \phi(w_i^B)||^2$



Challenge: Polysemy



Verb [edit]

run (third-person singular simple present runs, present participle running, simple past ran, past participle run or (nonstandard, colloquial) ran)

- 1. To move swiftly.
 - 1. (*intransitive*) To move forward quickly upon two feet by alternately making a short jump off either foot. [coordinate term ▲] [quotations ▼]

Coordinate term: walk

Run, and you might still catch the train!

2. (intransitive) To go at a fast pace; to move quickly. [quotations ▼]

I have been **running** all over the building looking for him.

Sorry, I've got to run; my house is on fire!

3. (transitive) To cover (a course or a distance) by running.

I can run a mile, but I can't run the cross-country course.

4. (transitive) To complete a running course or event in (a given time).

I was hoping to make the team, but I didn't run the qualifying time.

5. (*intransitive*) To move briskly or smoothly with a motion of sliding, rolling, sweeping etc.

The shuttle runs back and forth on these rollers.

As its name suggests, the monorail runs on a single rail

I felt her fingers running over my cheek

• • •

- 31. To tend, as to an effect or consequence; to incline. [quotations ▼]
- 32. To have a legal course; to be attached; to continue in force, effect, or operation; to follow; to go in company. [quality covenants run with the land.
- 33. To encounter or suffer (a particular, usually bad, fate or misfortune). [quotations ▼]
- 34. (golf) To strike (the ball) in such a way as to cause it to run along the ground, as when approaching a hole.
- 35. (video games, rare) To speedrun.
- 36. (sports, especially baseball) To eject from a game or match.

Jackson got himself **run** in the top of the sixth for arguing a borderline strike three call.

37. To press (a bank, etc.) with immediate demands for payment.

Noun [edit]

run (plural runs)

1. Act or instance of running, of moving rapidly using the feet. [quotations ▼]

I just got back from my morning run.

2. Act or instance of hurrying (to or from a place) (*not necessarily on foot*); dash or errand, trip. [quotations ▼]

I need to make a run to the store.

3. A pleasure trip. [quotations ▼]

Let's go for a run in the car.

- 4. Flight, instance or period of fleeing. [quotations ▼]
- 5. Migration of fish.
- 6. A group of fish that migrate, or ascend a river for the purpose of spawning.
- 7. A literal or figurative path or course for movement relating to:
 - 1. A (regular) trip or route. [quotations ▼]

The bus on the Cherry Street run is always crowded.

The route taken while running or skiing.

Which **run** did you do today?

• • •

21. A line of knit stitches that have unravelled, particularly in a nylon stocking. [quotations ▼]

I have a run in my stocking.

- 22. (nautical) The stern of the underwater body of a ship from where it begins to curve upward and inward.
- 23. (*mining*) The horizontal distance to which a drift may be carried, either by licence of the proprietor of a mine or by the nature of the formation; also, the direction which a vein of ore or other substance takes.
- 24. A pair or set of millstones.
- 25. One's gait while running; the way one runs.

I think they only have a weird run because their leg hurts.

Adjective [edit]

run (not comparable)

In a liquid state; melted or molten. [quotations ▼]

Put some **run** butter on the vegetables.

- 2. Cast in a mould. [quotations ▼]
- 3. Exhausted; depleted (especially with "down" or "out").
- 4. (of a zoology) Travelled, migrated; having made a migration or a spawning run. [quotations ▼]
- 5. Smuggled.

run brandy

Challenge: Polysemy



- Word embedding now has to somehow encode the meanings of all 67 senses of "run"...
 - Why is this hard?
- What could we do instead?
- Contextualized word embeddings: compute representation of individual words dependent on the context in which it appears
 - The resulting architectures were some of the main catalysts for the design of modern LLMs!

Opportunity: Sentence Representations



- If we have representations for parts of sentences (wordtypes), then can we get a representation of the whole sentence?
- One option: bag-of-words representation

$$\phi(\overline{x}) = \frac{1}{|\overline{x}|} \sum_{i=0}^{|\overline{x}|} \phi(x_i)$$

What's missing?

Word Embeddings in the Age of LLMs

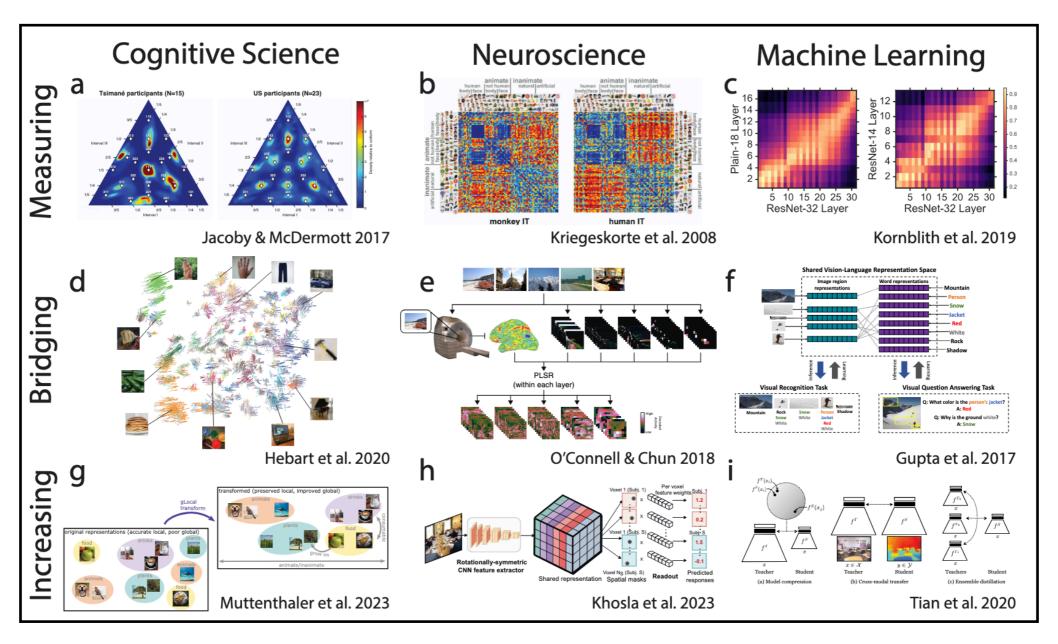


- Building a language technology pre-LLMs, for a target task:
 - Download some pre-trained word embeddings from the Internet
 - Initialize your language model with these word embeddings (all other parameters randomized)
 - Use your task-specific data to fine-tune the other parameters (possibly also fine-tuning the word embeddings)
 - Having a good starting point via word embeddings is really important, especially with little task-specific data
- But we don't do this anymore
 - Language models still learn embeddings specific to each wordtype
 - But we don't download them from the Internet we download the whole language model

Word Embeddings in the Age of LLMs



So why might we still care about word embeddings?

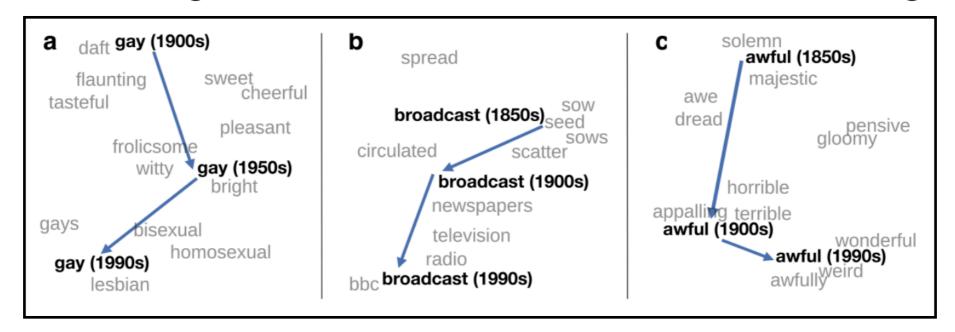


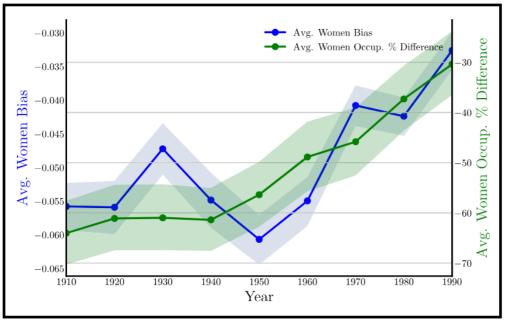
Evaluating how well an LLM's learned concept space correlates with human concept spaces

Word Embeddings in the Age of LLMs



So why might we still care about word embeddings?





1910	1950	1990
Charming	Delicate	Maternal
Placid	Sweet	Morbid
Delicate	Charming	Artificial
Passionate	Transparent	Physical
Sweet	Placid	Caring
Dreamy	Childish	Emotional
Indulgent	Soft	Protective
Playful	Colorless	Attractive
Mellow	Tasteless	Soft
Sentimental	Agreeable	Tidy

Analyze how language changes over time