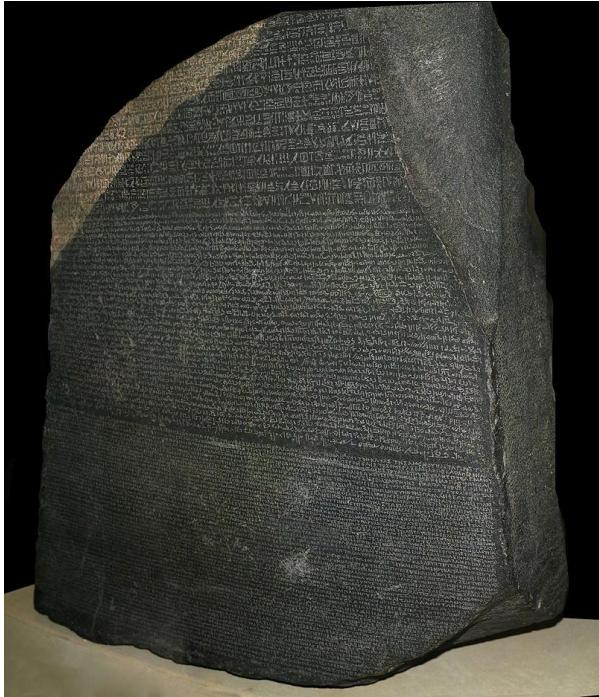


# Introduction to NLP

141

A Brief Overview of Languages and Linguistics



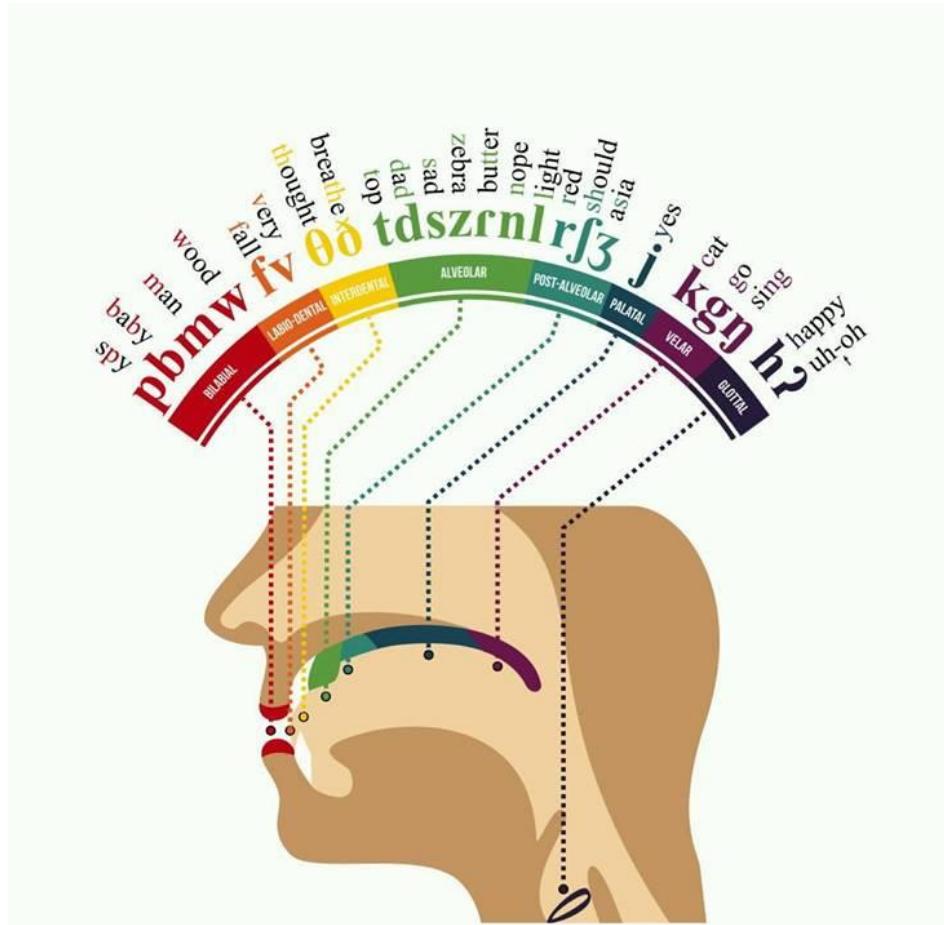
خ	ح	ج	ث	ت	ب	ا
kha	haa	jiim	thaa	taa	baa	alif
ص	ش	س	ز	ر	ذ	د
saad	shiin	siin	zaay	raa	thaal	daal
ق	ف	غ	ع	ظ	ط	ض
qaaf	faa	ghayn	ayn	thaa	taa	daad
ي	و	ه	ن	م	ل	ك
yaa	waaw	ha	nuun	miim	laam	kaaf



<b>Аа</b>	a (as in cat)	<b>Кк</b>	k (as in kick)	<b>Фф</b>	f (as in foot)
<b>Бб</b>	b (as in bus)	<b>Лл</b>	l (as in love)	<b>Хх</b>	h (like 'ch' in Bach)
<b>Вв</b>	v (as in very)	<b>Мм</b>	m (as in marry)	<b>Чч</b>	ts (as in puts)
<b>Гг</b>	g (as in good)	<b>Нн</b>	n (as in no)	<b>Цц</b>	ch (as in check)
<b>Дд</b>	d (as in dog)	<b>Оо</b>	o (as in hot)	<b>Шш</b>	sh (as in shut)
<b>Ее</b>	e (as in egg)	<b>Пп</b>	p (as in pot)	<b>Щщ</b>	sht (like 'shed' in pushed)
<b>Жж</b>	zh (like 's' in leisure)	<b>Рр</b>	r (as in red)	<b>ъъ</b>	a (like 'u' in but)
<b>Зз</b>	z (as in zoo)	<b>Сс</b>	s (as in sit)	<b>ьъ</b>	(consonant softening sound)
<b>Ии</b>	i (as in instant)	<b>Тт</b>	t (as in tree)	<b>Юю</b>	yu (like you)
<b>Йй</b>	y (as in young)	<b>Үү</b>	u (as in yule)	<b>Яя</b>	ya (as in yank)



# Consonants in English



**ENGLISH IPA**

WWW.LANGUAGEBASECAMP.COM

# IPA Chart (consonants)

辅音

CONSONANTS (PULMONIC)

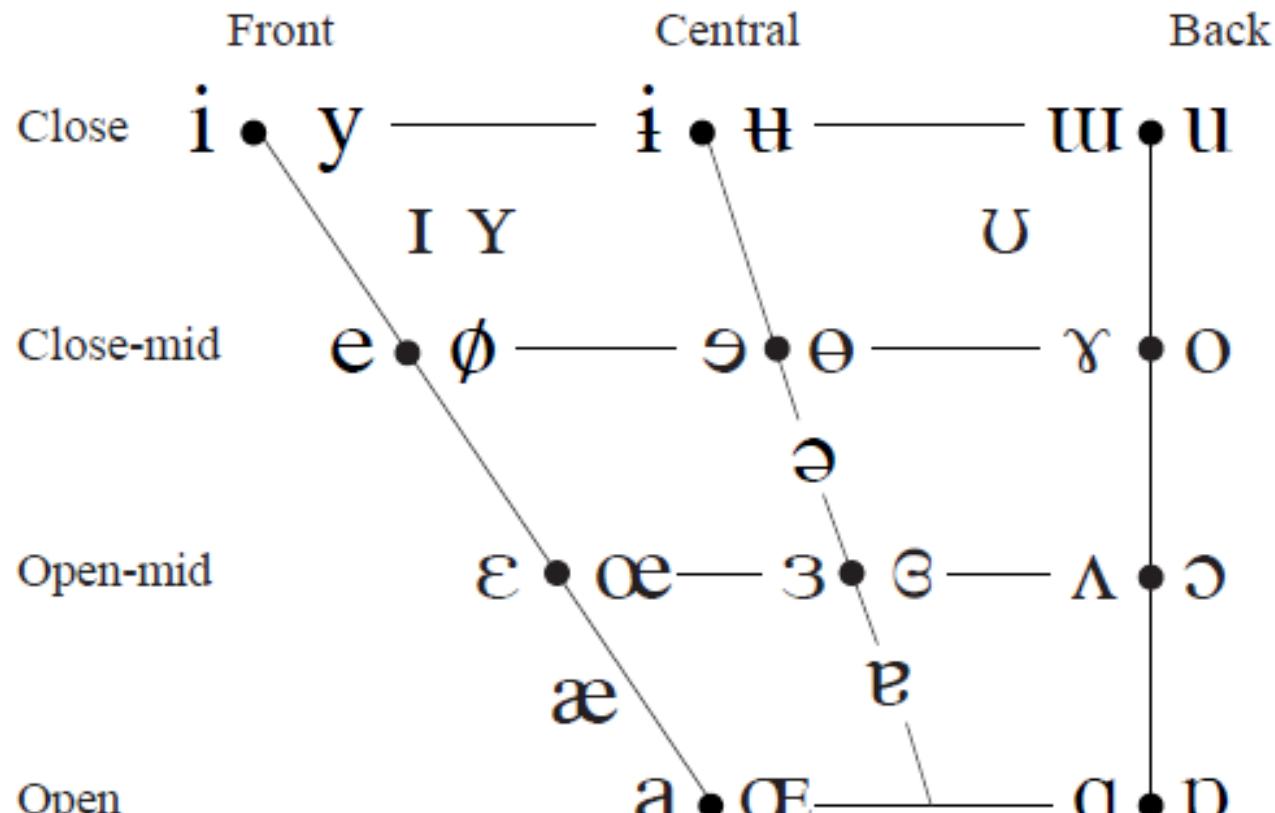
© 2005 IPA

	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Glottal
Plosive	p b			t d		t̪ d̪	c ɟ	k g	q ɢ		ʔ
Nasal	m	n̪		n		ɳ	j̪n̪	ŋ	ɳ		
Trill	B			r					R		
Tap or Flap		v̪		f		t̪					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ɟ	xɣ	χʁ	ħʕ	hɦ
Lateral fricative			ɬ ɬ̪								
Approximant		v̪		ɺ		ɻ	j̪	w̪			
Lateral approximant			ɭ		ɭ̪	ɻ̪		L			

Where symbols appear in pairs, the one to the right represents a voiced consonant. Shaded areas denote articulations judged impossible.

# IPA Chart (vowels)

## VOWELS

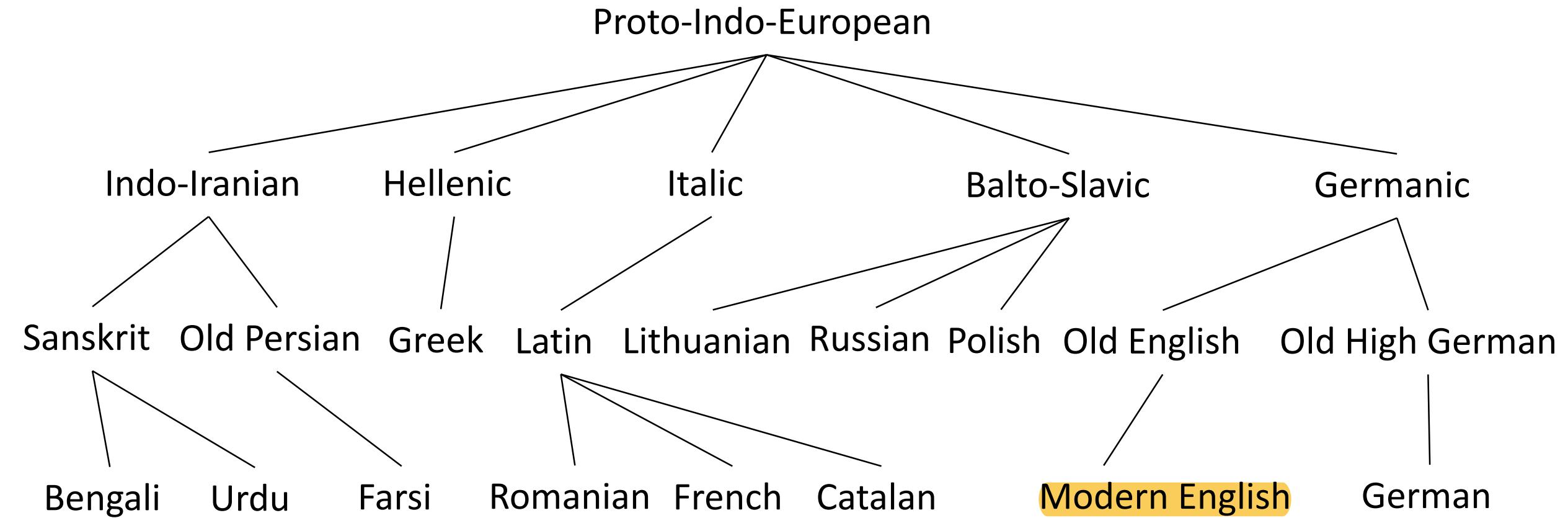


Where symbols appear in pairs, the one to the right represents a rounded vowel.

# Indo-European Words for Two



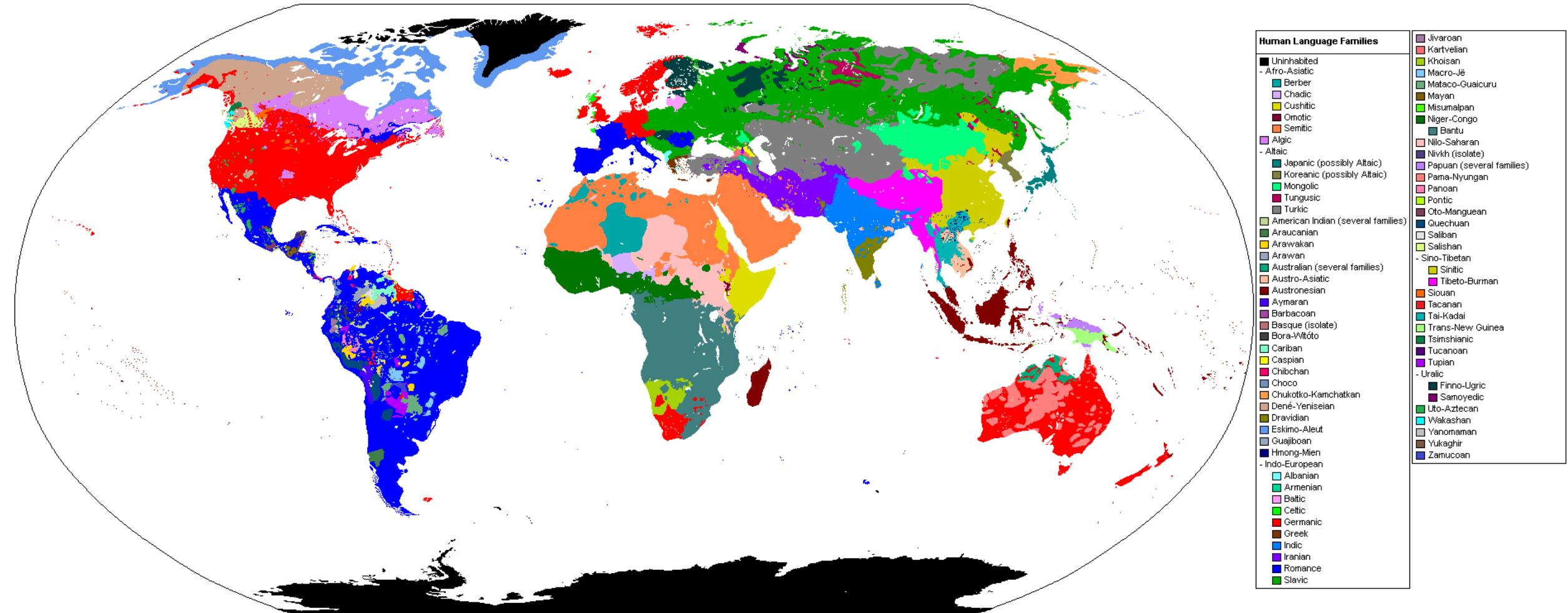
# Some Indo-European languages



# Some non-Indo-European Languages

- Altaic
  - Turkish
- Uralic (Finno-Ugric)
  - Finnish
  - Hungarian
- Semitic
  - Arabic
  - Hebrew
- Uto-Aztecanc

# Language Families



By Industrius at English Wikipedia. Later version(s) were uploaded by Mttl at English Wikipedia. (Image:BlankMap-World.png by User:Vardion) [GFDL ([www.gnu.org/copyleft/fdl.html](http://www.gnu.org/copyleft/fdl.html))], via Wikimedia Commons

# Language Diversity

[Afro-Asiatic](#) (374)  
[Alacalufan](#) (2)  
[Algic](#) (44)  
[Altaic](#) (66)  
[Amto-Musan](#) (2)  
[Andamanese](#) (13)  
[Arafundi](#) (3)  
[Arai-Kwomtari](#) (10)  
[Arauan](#) (5)  
[Araucanian](#) (2)  
[Arawakan](#) (59)  
[Arutani-Sape](#) (2)  
[Australian](#) (264)  
[Austro-Asiatic](#) (169)  
**Austronesian** (1257)  
[Aymaran](#) (3)  
[Barbacoan](#) (7)  
[Basque](#) (1)  
[Bayono-Awbono](#) (2)  
[Border](#) (15)  
[Caddoan](#) (5)  
[Cahuapanan](#) (2)

[Carib](#) (31)  
[Central Solomons](#) (4)  
[Chapacura-Wanham](#) (5)  
[Chibchan](#) (21)  
[Chimakuan](#) (1)  
[Choco](#) (12)  
[Chon](#) (2)  
[Chukotko-Kamchatkan](#) (5)  
[Chumash](#) (7)  
[Coahuiltecan](#) (1)  
[Constructed language](#) (1)  
[Creole](#) (82)  
[Deaf sign language](#) (130)  
[Dravidian](#) (85)  
[East Bird's Head-Sentani](#) (8)  
[East Geelvink Bay](#) (11)  
[East New Britain](#) (7)  
[Eastern Trans-Fly](#) (4)  
[Eskimo-Aleut](#) (11)  
[Guahiban](#) (5)  
[Gulf](#) (4)

[Harakmbet](#) (2)  
[Hibito-Cholon](#) (2)  
[Hmong-Mien](#) (38)  
[Hokan](#) (23)  
[Huavean](#) (4)  
[Indo-European](#) (439)  
[Iroquoian](#) (9)  
[Japonic](#) (12)  
[Jivaroan](#) (4)  
[Kartvelian](#) (5)  
[Katukinan](#) (3)  
[Kaure](#) (4)  
[Keres](#) (2)  
[Khoisan](#) (27)  
[Kiowa Tanoan](#) (6)  
[Lakes Plain](#) (20)  
**Language isolate** (50)  
[Left May](#) (2)  
[Lower Mamberamo](#) (2)  
[Lule-Vilela](#) (1)  
[Macro-Ge](#) (32)  
[Mairasi](#) (3)

[Maku](#) (6)  
[Mascoian](#) (5)  
[Mataco-Guaicuru](#) (12)  
[Mayan](#) (69)  
[Maybrat](#) (2)  
[Misumalpan](#) (4)  
[Mixed language](#) (23)  
[Mixe-Zoque](#) (17)  
[Mongol-Langam](#) (3)  
[Mura](#) (1)  
[Muskogean](#) (6)  
[Na-Dene](#) (46)  
[Nambiquaran](#) (7)  
**Niger-Congo** (1532)  
[Nilo-Saharan](#) (205)  
[Nimboran](#) (5)  
[North Bougainville](#) (4)  
[North Brazil](#) (1)  
[North Caucasian](#) (34)  
[Oto-Manguean](#) (177)  
[Panoan](#) (28)  
[Pauwasi](#) (5)  
[Peba-Yaguan](#) (2)  
[Penutian](#) (33)  
[Piawi](#) (2)  
[Pidgin](#) (17)  
[Quechuan](#) (46)  
[Ramu-Lower Sepik](#) (32)  
[Salishan](#) (26)  
[Salivan](#) (3)  
[Senagi](#) (2)  
[Sepik](#) (56)  
**Sino-Tibetan** (449)  
[Siouan](#) (17)  
[Sko](#) (7)  
[Somahai](#) (2)  
[South Bougainville](#) (9)  
[South-Central Papuan](#) (22)  
[Tacanan](#) (6)  
[Tai-Kadai](#) (92)  
[Tarascan](#) (2)  
[Tequistlatecan](#) (2)  
[Tor-Kwerba](#) (24)

[Torricelli](#) (56)  
[Totonacon](#) (12)  
**Trans-New Guinea** (477)  
[Tucanoan](#) (25)  
[Tupi](#) (76)  
[Unclassified](#) (73)  
[Uralic](#) (37)  
[Uru-Chipayo](#) (2)  
[Uto-Aztecian](#) (61)  
[Wakashan](#) (5)  
[West Papuan](#) (23)  
[Witotoan](#) (6)  
[Yanomam](#) (4)  
[Yele-West New Britain](#) (3)  
[Yeniseian](#) (2)  
[Yuat](#) (6)  
[Yukaghir](#) (2)  
[Yuki](#) (2)  
[Zamucoan](#) (2)  
[Zaparoan](#) (7)

# NACLO Problem

- <http://nacloweb.org/resources/problems/2012/N2012-D.pdf>
- <http://nacloweb.org/resources/problems/2012/N2012-DS.pdf>
- Problem by Dragomir Radev

Many languages are related to each other for historical reasons. They may have a common ancestor or they may have borrowed words from each other. Linguists group languages into families and branches, based on their common ancestry.

Here is a list of translations of the first article of the Universal Declaration of Human Rights in 17 languages:

Your task is to identify similarities among these languages and group them into seven clusters (groups) of related languages as sketched in the diagram below:

[http://unicode.org/udhr/assemblies/first\\_article\\_all.html](http://unicode.org/udhr/assemblies/first_article_all.html)

Here is a list of translations of the first article of the Universal Declaration of Human Rights in 17 languages:

- A. (English) All human beings are born free and equal in dignity and rights. They are endowed with reason and conscience and should act towards one another in a spirit of brotherhood.
- B. (Latin) Omnes homines dignitate et iure liberi et pares nascuntur, rationis et conscientiae participes sunt, quibus inter se concordiae studio est agendum.
- C. Vsi ljudje se rodijo svobodni in imajo enako dostojanstvo in enake pravice. Obdarjeni so z razumom in vestjo in bi morali ravnati drug z drugim kakor bratje.
- D. Dieub ha par en o dellezegezh hag o gwirioù eo ganet an holl dud. Poell ha skiant zo dezho ha dleout a reont bevañ an eil gant egile en ur spered a genvreudeuriezh.
- E. Tuots umans naschan libers ed equals in dignità e drets. Els sun dotats cun intellet e conscienza e desan agir tanter per in uin spiert da fraternità.
- F. Toate ființele umane se nasc libere și egale în demnitate și în drepturi. Ele sunt înzestrăte cu rațiune și conștiință și trebuie să se comporte unii față de altele în spiritul fraternității.
- G. Genir pawb yn rhydd ac yn gydradd â'i gilydd mewn urddas a hawliau. Fe'u cynysgaeddir â rheswm a chydwybod, a dylai pawb ymddwyn y naill at y llall mewn ysbryd cymodlon.
- H. Visi žmonės gimsta laisvi ir lygūs savo orumu ir teisėmis. Jiems suteiktas protas ir sąžinė ir jie turi elgtis vienas kito atžvilgiu kaip broliai.
- I. Totu sos èsseres umanos naschint lliberos e iguals in dignitat e in drets. Issos tenent sa resone e sa cussèntzia e depent operare s'unu cun s'àteru cun ispiritu de fraternitat.
- J. Gizon-emakume guztiak aske jaiotzen dira, duintasun eta eskubide berberak dituztela; eta ezaguera eta kontzientzia dutenez gero, elkarren artean senide legez jokatu beharra dute.
- K. Kai rahvas roittahes vällinny da taza-arvozinu omas arvos da oigevuksis. Jogahizele heis on annettu mieli da omatundo da heil vältämättäh pidäy olla keskenäh, kui vellil.

L. Všetci ľudia sa rodia slobodní a sebe rovní , čo sa týka ich dôstojnosti a práv. Sú obdarení rozumom a majú navzájom jednat' v bratskom duchu.

M. Nascinu tutti l'omi libari è pari di dignità è di diritti. Pussediu a raghjoni è a cuscenza è li tocca ad agiscia trà elli di modu fraternu.

N. Saoláitear na daoine uile saor agus comhionann ina ndínit agus ina gcearta. Tá bauidh an réasúin agus an choinsiasa acu agus dlíd iad fén d'iompar de mheon bhrithreachais i leith a chéile.

O. Visi cilvēki piedzimst brīvi un vienlīdzīgi savā pašcienā un tiesībās. Viņi ir apveltīti ar saprātu un sirdsapziņu, un viņiem jāizturas citam pret citu brālības garā.

P. Kaikki ihmiset syntyvät vapaina ja tasavertaisina ar voltaan ja oikeuksiltaan. Heille on annettu järki ja omatunto, ja heidän on toimittava toisiaan kohtaan veljeyden hengessä.

Q. Wszyscy ludzie rodzą się wolni i równi pod względem swej godności i swych praw. Są oni obdarzeni rozumem i sumieniem i powinni postępować wobec innych w duchu braterstwa.

Here is a list of translations of the first article of the Universal Declaration of Human Rights in 17 languages:

English

A. (English) All human beings are born free and equal in dignity and rights. They are endowed with reason and conscience and should act towards one another in a spirit of brotherhood.

Latin

B. (Latin) Omnes homines dignitate et iure liberi et pares nascuntur, rationis et conscientiae participes sunt, quibus inter se concordiae studio est agendum.

Slovenian

C. Vsi ljudje se rodijo svobodni in imajo enako dostojanstvo in enake pravice. Obdarjeni so z razumom in vestjo in bi morali ravnati drug z drugim kakor bratje.

Breton

D. Dieub ha par en o dellezegezh hag o gwirioù eo ganet an holl dud. Poell ha skiant zo dezho ha dleout a reont bevañ an eil gant egile en ur spered a genvreudeuriezh.

Romansch

E. Tuots umans naschan libers ed equals in dignità e drets. Els sun dotats cun intellet e conscienza e desan agir tanter per in uin spiert da fraternità.

Romanian

F. Toate ființele umane se nasc libere și egale în demnitate și în drepturi. Ele sunt înzeștrăte cu rațiune și conștiință și trebuie să se comporte unii față de altele în spiritul fraternității.

Welsh

G. Genir pawb yn rhydd ac yn gydradd â'i gilydd mewn urddas a hawliau. Fe'u cynysgaeddir â rheswm a chydwybod, a dylai pawb ymddwyn y naill at y llall mewn ysbryd cymodlon.

Lithuanian

H. Visi žmonės gimsta laisvi ir lygūs savo orumu ir teisėmis. Jiems suteiktas protas ir sąžinė ir jie turi elgtis vienas kito atžvilgiu kaip broliai.

Sardinian

I. Totu sos èsseres umanos naschint lìberos e egaules in dinnidade e in deretos. Issos tenent sa resone e sa cussèntzia e depent operare s'unu cun s'àteru cun ispiritu de fraternidade.

Basque

J. Gizon-emakume guztiak aske jaiotzen dira, duintasun eta eskubide berberak dituztela; eta ezaguera eta kontzientzia dutenez gero, elkarren artean senide legez jokatu beharra dute.

Karelian

K. Kai rahvas roittahes vällinny da taza-arvozinnu omas arvos da oigevuksis. Jogahizele heis on annettu mieli da omatundo da heil vältämättäh pidäy olla keskenäh, kui vellil.

Slovak	L. Všetci ľudia sa rodia slobodní a sebe rovní , čo sa týka ich dostôjnosti a práv. Sú obdarení rozumom a majú navzájom jednat' v bratskom duchu.
Corsican	M. Nascinu tutti l'omi libari è pari di dignità è di diritti. Pusseddu a raghjoni è a cuscenza è li tocca ad agis-cia trà elli di modu fraternu.
Irish	N. Saoláitear na daoine uile saor agus comhionann ina ndínit agus ina gcearta. Tá bauidh an réasúin agus an choinsiasa acu agus dlíd iad féin d'iompar de mheon bhrithreachais i leith a chéile.
Latvian	O. Visi cilvēki piedzimst brīvi un vienlīdzīgi savā pašcieņā un tiesībās. Viņi ir apveltīti ar saprātu un sirdsapziņu, un viņiem jāizturas citam pret citu brālības garā.
Finnish	P. Kaikki ihmiset syntyvät vapaina ja tasavertaisina arvoltaan ja oikeuksiltaan. Heille on annettu järki ja omatunto, ja heidän on toimittava toisiaan kohtaan veljeyden hengessä.
Polish	Q. Wszyscy ludzie rodzą się wolni i równi pod względem swej godności i swych praw. Są oni obdarzeni rozumem i sumieniem i powinni postępować wobec innych w duchu braterstwa.

# Language Families

1. CLQ Slavic
2. BEFIM Romance
3. J Basque
4. HO Baltic
5. DGN Celtic
6. KP Finno-Ugric
7. A English

## How English has changed over the last 1000 years: the 23rd Psalm

---

### *Modern (1989)*

The Lord is my shepherd, I lack nothing.  
He lets me lie down in green pastures.  
He leads me to still waters.

### *King James Bible (1611)*

The Lord is my shepherd, I shall not want.  
He maketh me to lie down in green pastures.  
He leadeth me beside the still waters.

### *Middle English (1100–1500)*

Our Lord gouerneth me, and nothyng shal defailen to me.  
In the sted of pastur he sett me ther.  
He norissed me upon water of fyllyng.

### *Old English (800–1066)*

Drihten me raet, ne byth me nanes godes wan.  
And he me geset on swythe good feohland.  
And fedde me be waetera stathum.

---



- How can I say "Togetherness"  
on german language?

- "Zusammengehörigkeitsgefühl"

# Language Diversity

- Articles
  - English vs. Russian
- Cases (e.g., in Latin)
  - Puer puellam vexat
- Sound systems
  - Glottal stop (the middle sound in “uh-oh”) - pro
  - Velar fricatives - articulated with the back of the tongue at the soft palate
    - Voiceless /χ/ - used e.g., in Russian
    - Voiced /γ/ - used e.g., in Modern Greek
- Social status (e.g., in Japanese)
  - otousan, お父さん = someone else's father
  - chichi, 父 = one's own father

# Language Universals

- Two types
  - unconditional
  - conditional
- Examples
  - All languages have verbs and nouns
  - All spoken languages have consonants and vowels
  - [Greenberg 1] “In declarative sentences with nominal subject and object, the dominant order is almost always one in which the subject precedes the object.”
  - [Greenberg 29] “If a language has inflection, it always has derivation.”

# WALS: the World Atlas of Language Structures

- <http://wals.info>
- Feature 83A: Order of Object and Verb
  - by Matthew S. Dryer
  - OV (713 languages), VO (705), no dominant order (101)
  - <http://wals.info/feature/83A#2/18.0/152.9>
- Other features:
  - 18A Absence of common consonants (by Ian Maddieson):  
no bilabials (5 languages), no fricatives (49), no nasals (12)
  - 67A Inflectional future tense (by Östen Dahl, Viveka Velupillai):  
yes (110), no (112)

# Links about World Languages

- Ethnologue
  - <http://www.ethnologue.com/>
- Number words in many languages
  - <http://www.zompist.com/numbers.shtml>
- Endangered languages
  - <http://www.endangeredlanguages.com/>
- Google fights to save 3,054 dying languages
  - <http://www.cnn.com/2012/06/21/tech/web/google-fights-save-language-mashable/index.html>

# Introduction to NLP

143

Morphology and the Lexicon

# Mental Lexicon

- What is the meaning of cat?
  - Its pronunciation?
  - Part of speech?
- What is the meaning of wug?
- What is the meaning of cluvious?

“Runs”

- Two interpretations
- Affixes
  - prefixes, infixes, suffixes, circumfixes, null morpheme

Reduplication: A process of forming new words by repeating a morpheme or a part of a word. In the provided examples, "amimigo" in Pangasinan and "savavali" in Samoan are formed by repeating parts of the root words "amigo" and "savali," respectively.

## Morphological Examples

- Reduplication

- amigo = friend, amimigo = friends (in Pangasinan) [Rubino 2001]
- savali = he travels, savavali = they travel (in Samoan)

- Templatic morphology (e.g., Semitic languages):

- Imd (learn), lamad (he studied), limed (he taught), lumad (he was taught)

- Circumfixes

- spielen – gespielt (in German), light – enlighten (in English)

Pig Latin: A playful language game where the first consonant or consonant

- Pig Latin

cluster of an English word is moved to the end and followed by "-ay." For example, "happy" becomes "appyhay."

- appyhay

- Massa-freakin'-chussets

Inserting "freakin'": This is an example of playful infixation for emphasis or stylistic effect. In the word "education," you can insert "freakin'" to make it

- where can you insert "freakin'" in "education"?

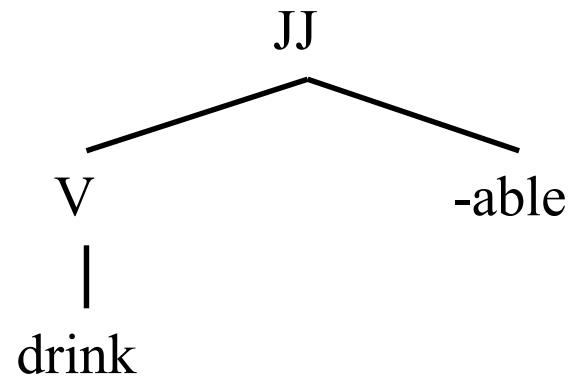
English.

# Answer

- The “freakin” infix is inserted
- ... to the left of the syllable that bears the main stress
  - edu-*freakin'*-cation
  - \* educa-*freakin'*-tion
  - \* e-*freakin'*-ducation
- though there can be exceptions

# Derivational Morphology

- Example
  - “er” (multiple interpretations)
- What do these morphemes mean?
  - prefix, stem, suffix, ending
  - ness, able, ing, re, un, er (adj)
  - JJ → V + “-able”
- Recursion:
  - unconcernednesses
- Ambiguity
  - uncloggable vs. unbelievable



# Answer to the Quiz

- **Unclogable**
  - unable to be clogged
  - able to be unclogged
- **Unbelievable**
  - unable to be believed
  - ? able to be unbelieved

# Inflectional Morphology

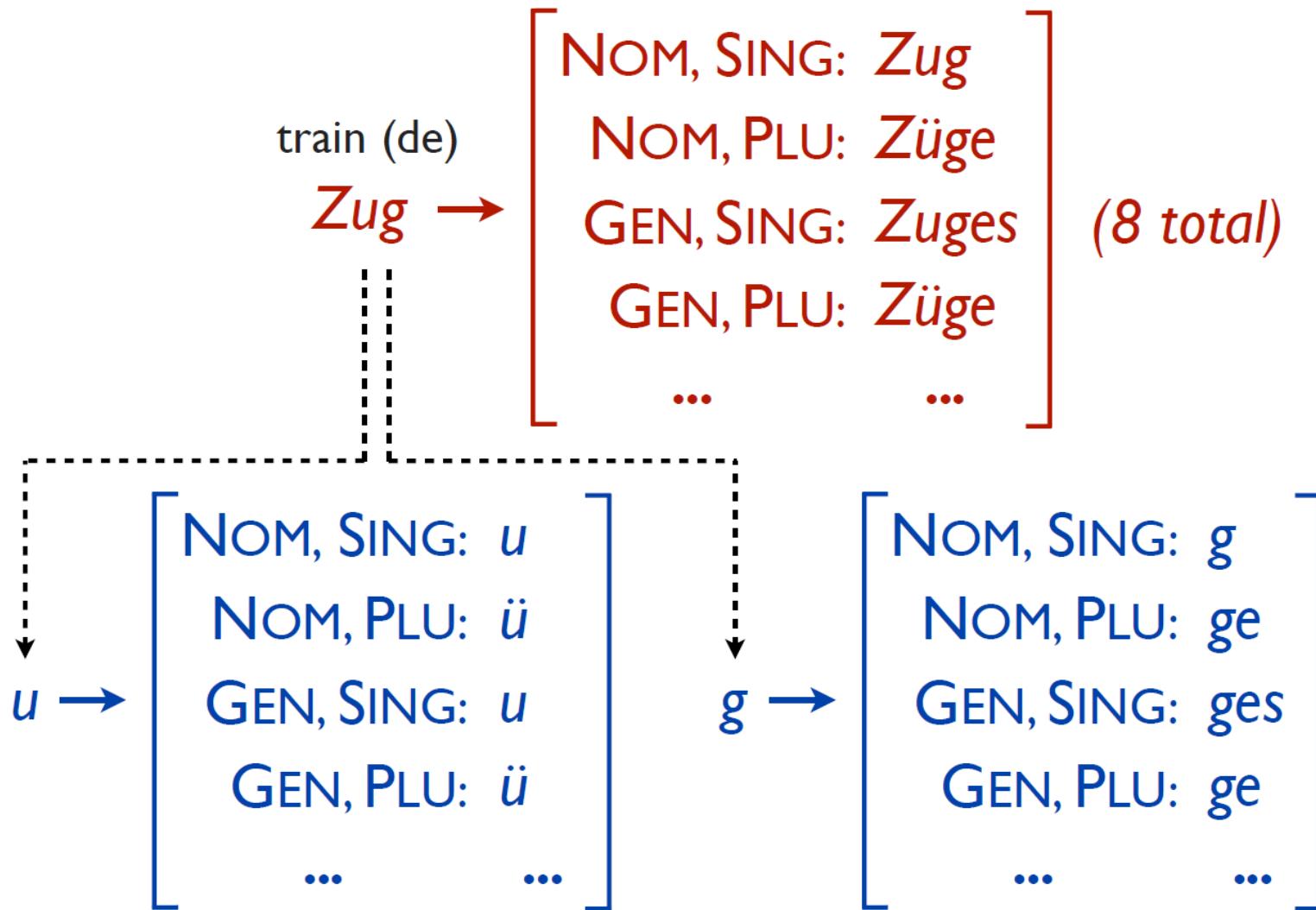
- Many forms
  - Tense, number, person, mood, aspect
  - Five verb forms in English
  - 40+ forms in French
  - Six cases in Russian:  
<http://www.departments.bucknell.edu/russian/language/case.html>
  - Up to 40,000 forms in Turkish
    - E.g., you cause X to cause Y to ... do Z)

# Noun and Pronoun Declension in Latin

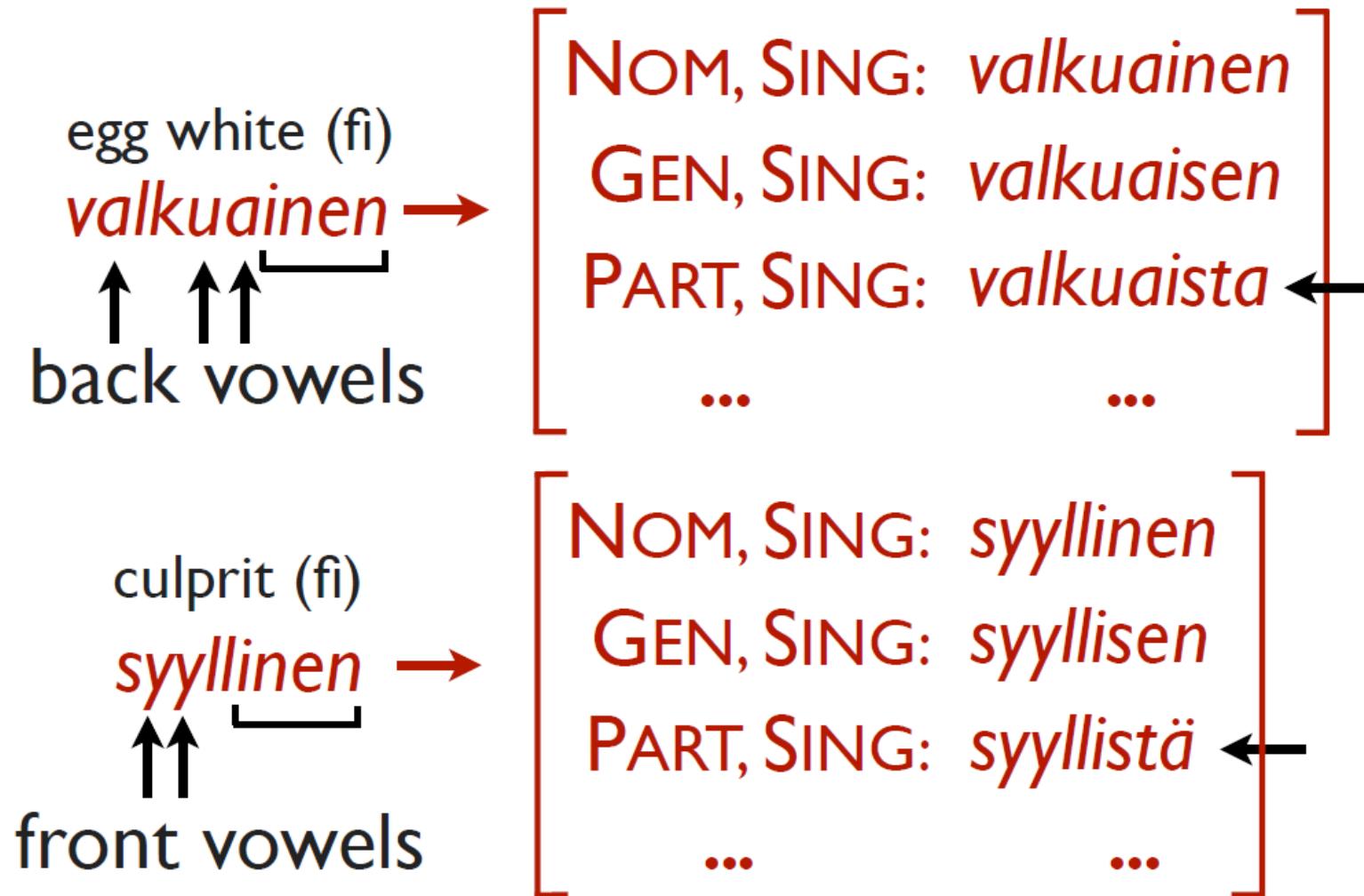
	<i>puer, puerī</i> boy m.		<i>ager, agrī</i> field m.		<i>vir, virī</i> man m.	
	Singular	Plural	Singular	Plural	Singular	Plural
<b>Nominative</b>	puer	puerī	ager	agrī	vir	virī
<b>Vocative</b>						
<b>Accusative</b>	puerum	puerōs	agrum	agrōs	virum	virōs
<b>Genitive</b>	puerī	puerōrum	agrī	agrōrum	virī	virōrum (virum)
<b>Dative</b>	puerō	puerīs	agrō	agrīs	virō	virīs
<b>Ablative</b>						

	<i>noster, nostra, nostrum</i> our, ours					
	Singular			Plural		
	Masculine	Feminine	Neuter	Masculine	Feminine	Neuter
<b>Nominative</b>	noster	nostra	nostrum	nostrī	nostrae	nostra
<b>Accusative</b>	nostrum	nostram		nostrōs	nostrās	
<b>Genitive</b>	nostrī		nostrae	nostrī	nostrōrum	nostrārum
<b>Dative</b>				nostrō	nostrīs	
<b>Ablative</b>	nostrō	nostrā				

# Automatic Inflection

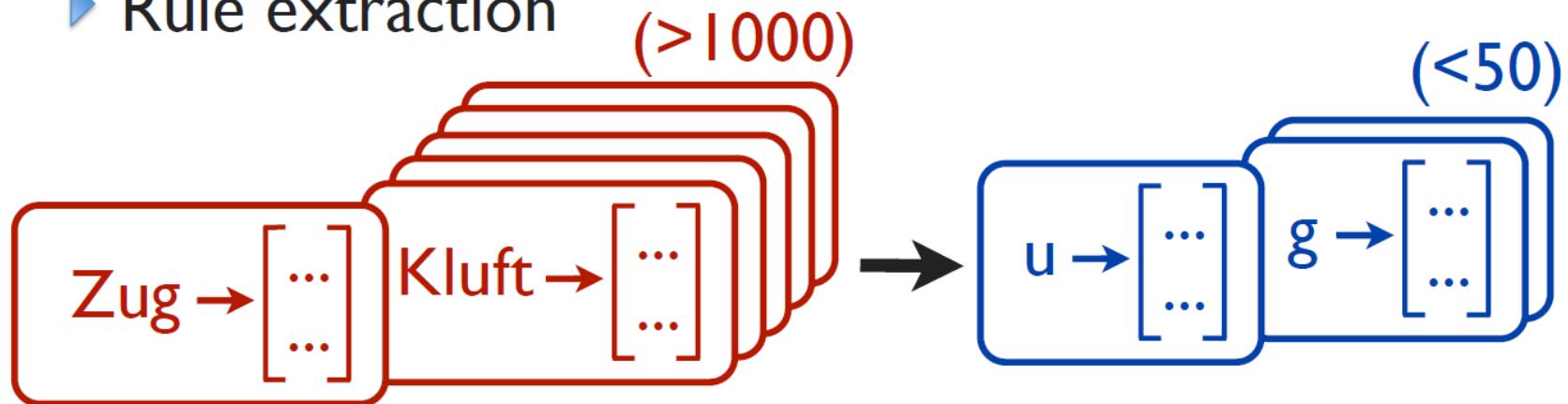


# Automatic Inflection

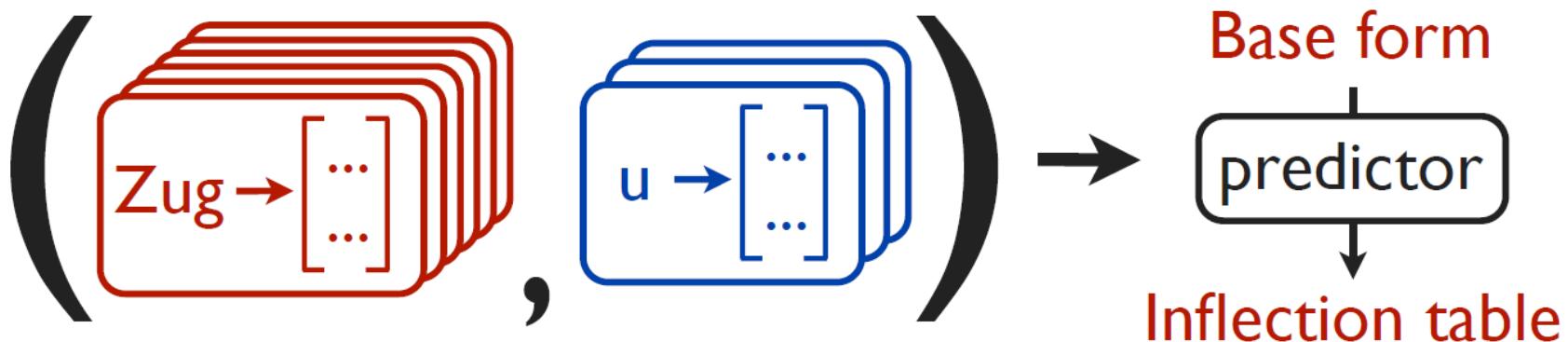


# Automatic Inflection

- ▶ Rule extraction



- ▶ Paradigm prediction



# Morphological Analysis

- sleeps = sleep + V + 3P + SG
- done = do + V + PP

Nice example from Kemal Oflazer

## Dancing in Andalusia

- A poem by the early 20th century Turkish poet Yahya Kemal Beyatlı.

## ENDÜLÜSTE RAKS

Zil, şal ve gül, bu bahçede raksın bütün hızı  
Şevk akşamında Endülüs, üç defa kırmızı  
Aşkın sihirli şarkısı, yüzlerce dildedir  
İspanya neşesiyle bu akşam bu zildedir

Yelpaze gibi çevrilir birden dönüşleri  
İşveyle devriliş, saçılış, örtünüşleri  
Her rengi istemez gözümüz şimdi aldadır  
İspanya dalga dalga bu akşam bu şaldadır

Alnında halka halkadır âşufte kâkülü  
Göğsünde yosma Gîrnata'nın en güzel gülü  
Altın kadeh her elde, güneş her gönüldedir  
İspanya varlığıyla bu akşam bu güldedir

Raks ortasında bir durup oynar, yürüür gibi  
Bir baş çevirmesiyle bakar öldürür gibi  
Gül tenli, kor dudaklı, kömür gözlü, sürmeli  
Seytan diyor ki sarmalı, yüz kerre öpmeli

Gözler kamaştıran şala, meftûn eden güle  
Her kalbi dolduran zile, her sineden ole!

## ENDÜLÜSTE RAKS

Zil, şal ve gül, bu bahçede raksın bütün hızı  
Şevk akşamında Endülüs, üç defa kırmızı  
Aşkın sihirli şarkısı, yüzlerce dildedir  
İspanya neşesiyle bu akşam bu **zildedir**

Yelpaze gibi çevrilir birden dönüşleri  
İşveyle devriliş, saçılış, örtünüşleri  
Her rengi **istemez** gözümüz şimdi aldadır  
İspanya dalga dalga bu akşam bu şaldadır

Alnında halka halkadır âşufte kâkülü  
Göğsünde yosma Gîrnata'nın en güzel gülü  
Altın kadeh her elde, güneş her gönüldedir  
İspanya **varlığıyla** bu akşam bu güldedir

Raks ortasında bir durup oynar, yürüür gibi  
Bir baş çevirmesiyle bakar öldürür gibi  
Gül tenli, kor dudaklı, kömür gözlü, surmeli  
Şeytan diyor ki sarmalı, yüz kerre öpmeli

Gözler **kamaştıran** şala, meftûn eden güle  
Her kalbi dolduran zile, her sineden ole!

**zildedir:** a verb derived from the locative case of the noun “zil” (castanet)  
“is at the castanet”

**dönüşleri:** plural infinitive and possessive form of the verb “dön” (rotate)

“their (act of) rotating”

**istemez:** negative present form of the verb “iste” (want)  
“it does not want”

**varlığıyla:** singular possessive instrumental-case of the noun “varlık” (wealth)  
“with its wealth”

**kamaştıran:** present participle of the verb “kamaş” (blind)  
“that which blinds....”

# Aligned Verses

Zil, şal ve gül, **bu** bahçede raksın bütün hızı

Şevk akşamında Endülüs, üç defa kırmızı

**Aşkın sihirli şarkısı**, yüzlerce **dildedir**

İspanya neş'esiyle bu akşam bu **zildedir**

Castañuela, mantilla y rosa. El baile veloz llena **el jardín**...

**En esta noche de jarana**, Andalucía se ve tres veces carmesí...

Cientas de **bocas recitan** la canción mágica del amor.

La alegría española esta noche, **está en las castañuelas**.

Castanets, shawl and rose. Here's the fervour of dance,

Andalusia is threefold red **in this evening of trance**.

Hundreds of **tongues utter love's magic refrain**,

**In these castanets** to-night survives the gay Spain,

Zimbel, Schal und Rose- Tanz **in diesem Garten** loht.

**In der Nacht der Lust** ist Andalusien dreifach rot!

Und in tausend **Zungen Liebeszauberlied erwacht-**

Spaniens Frohsinn **lebt in diesen Zimbeln** heute Nacht!

# Agglutinative Languages

- How does English become Turkish?

if we will be able to make ... become strong

if we will be able to make ... become strong

... strong become to make be able will if we

... sağlam +laş +tır +abil +ecek +se +k

↓  
... sağlamlaştıabileceksek

## ■ Aymara

- ch'uñuwinkaskiriyätwa
- ch'uñu +: +wi +na -ka +si -ka -iri +: +ya:t(a) +wa
- I was (one who was) always at the place for making ch'uñu'

ch'uñu	N		'freeze-dried potatoes'
+:		N>V	be/make ...
+wi		V>N	place-of
+na			in (location)
-ka		N>V	be-in (location)
+si			continuative
-ka			imperfect
-iri		V>N	one who
+:		N>V	be
+ya:ta		1P	recent past
+wa			affirmative sentencial

Example Courtesy of Ken Beesley

## ■ Finnish Numerals

- ☐ Finnish numerals are written as one word and all components inflect and agree in all aspects

## Kahdensienkymmenensienkahdeksansien

second                    tenth                            eighth    (28th)  
kaksi+Ord+Pl+Gen kymmenen+Ord+Pl+Gen kahdeksan+Ord+Pl+Gen  
**kahde** ns i en **kymmene** ns i en **kahdeksa** ns i en

## ■ Swahili

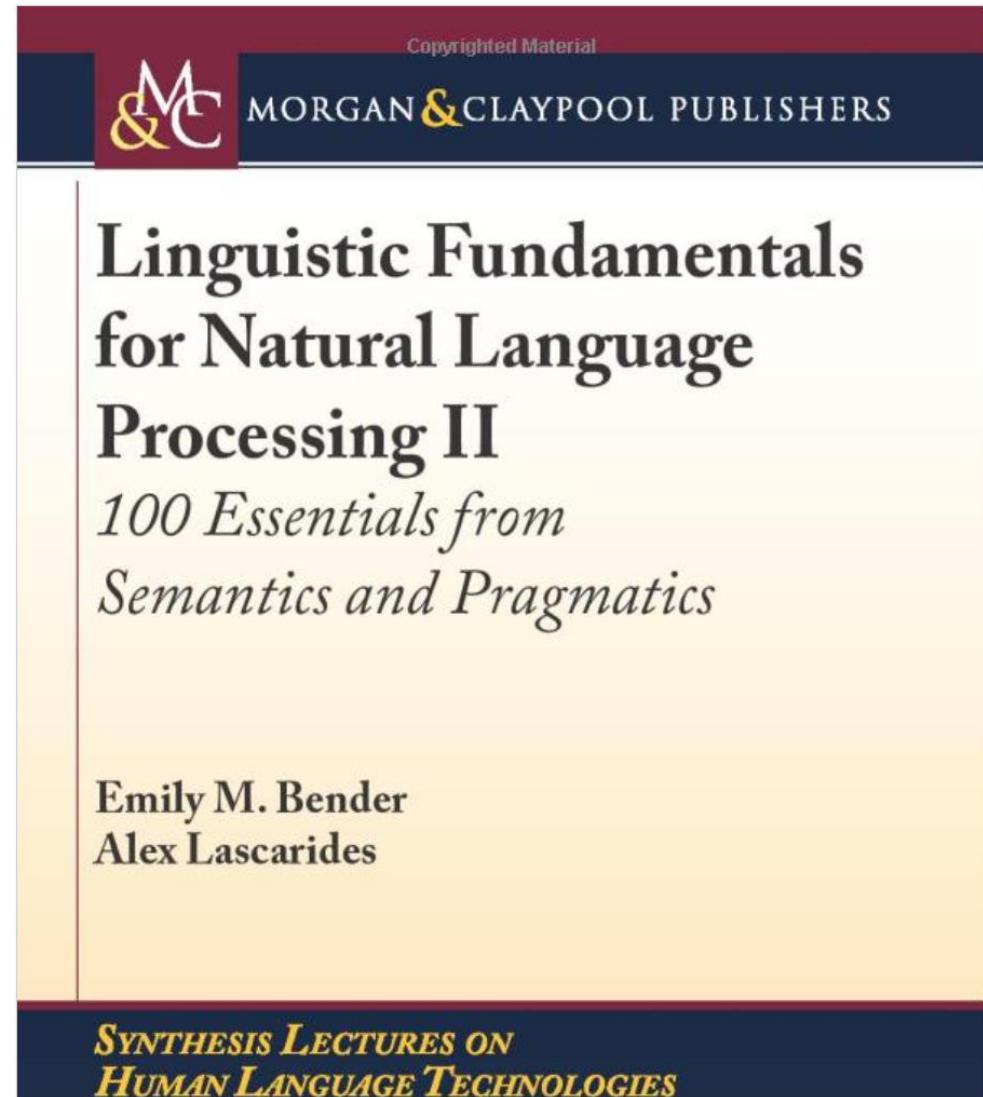
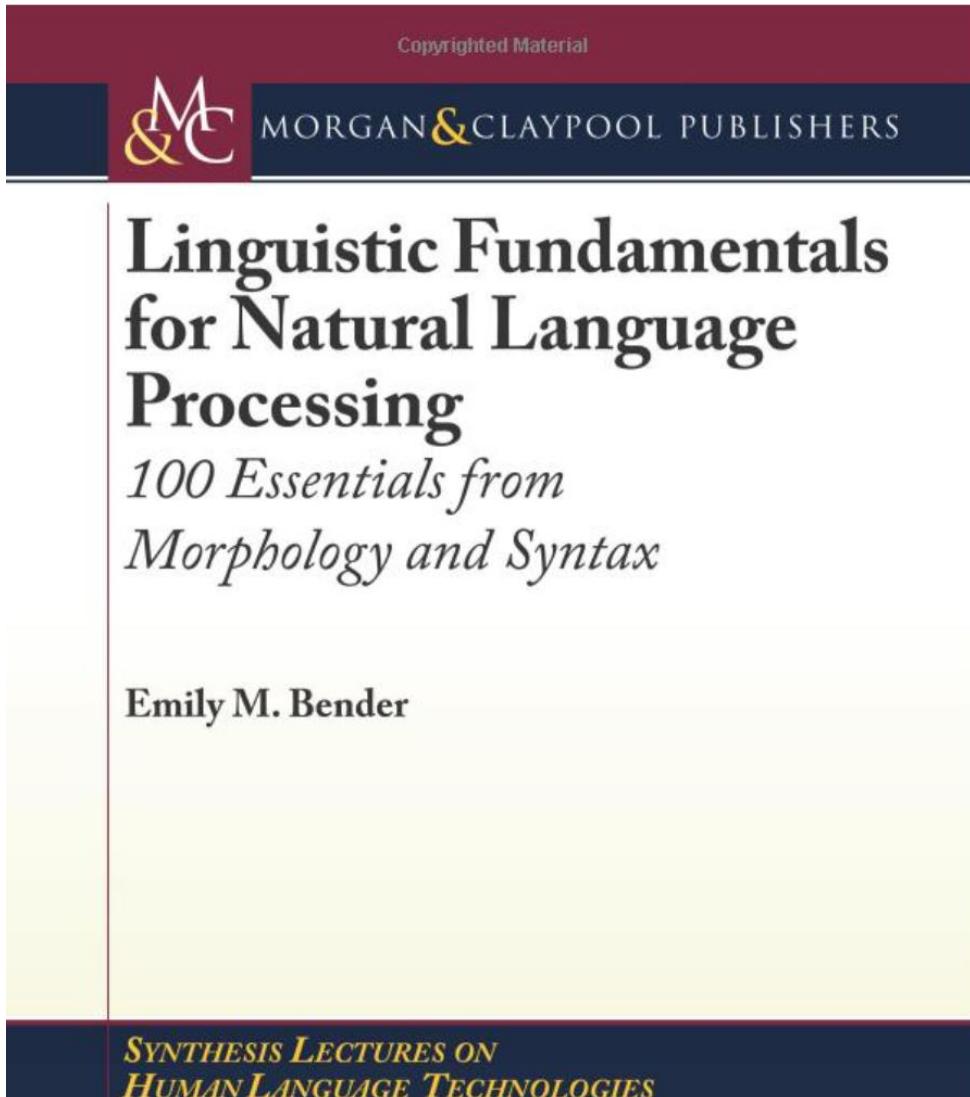
- walichotusomea = wa[Subject Pref]+li[Past]+cho[Rel Prefix]+tu[Obj Prefix 1PL]+**som**[read/Verb]+e[Prep Form]+a[]
    - that (thing) which they read for us
  - tulifika=tu[we]+li[Past]+**fik**[arrive/Verb]+a[]
    - We arrived
  - ninafika=ni[I]+na[Present]+**fik**[arrive/Verb]+a[]
    - I am arriving

# Turkish Vowel Harmony

	Front		Back	
	Unrounded	Rounded	Unrounded	Rounded
High	i	ü	ı	u
Low	e	ö	a	o

- Back vowels
  - in the room → oda~~da~~
  - at the door → kapı~~da~~
- Front vowels
  - at home → ev~~de~~
  - at the lake → göl~~de~~
  - on the bridge → köprü~~de~~

# Books by Emily Bender



# Introduction to NLP

144

Word Distributions

# Word Distributions

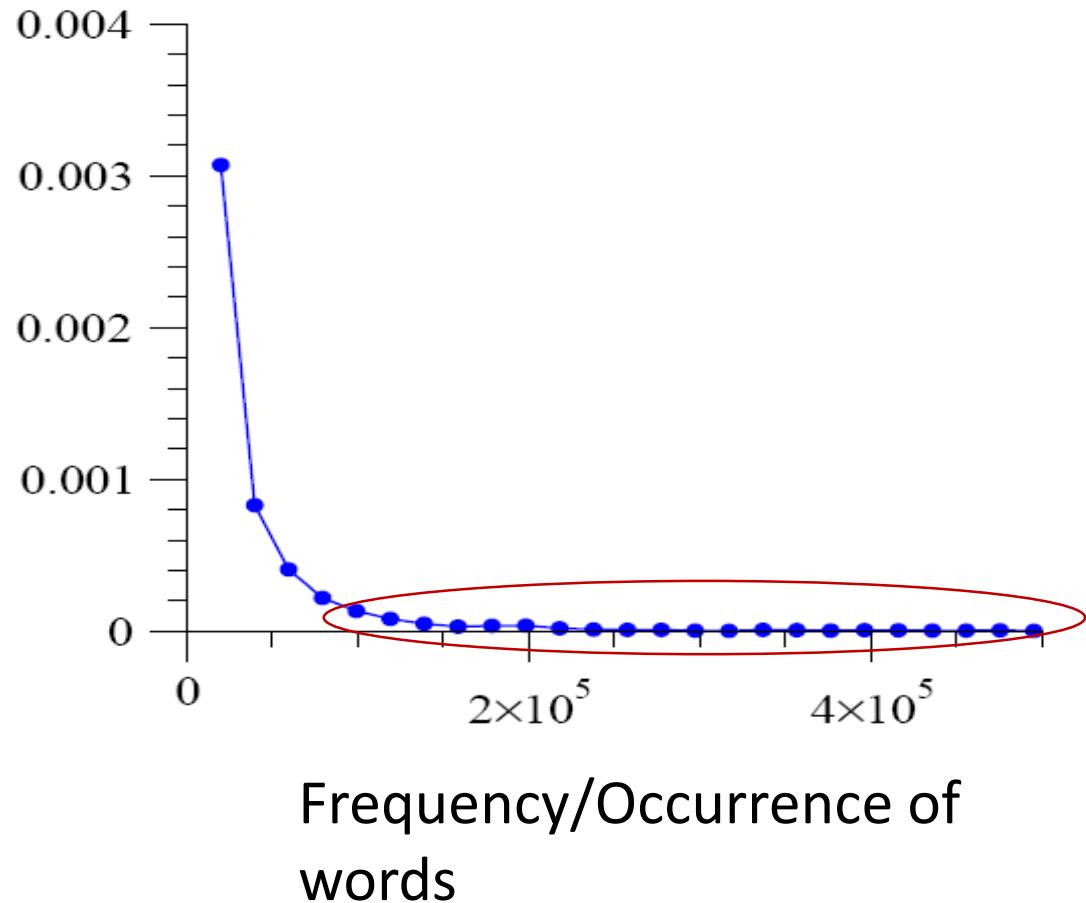
- Words are not distributed evenly!
  - Same goes for letters of the alphabet (ETAOIN SHRDLU), city sizes, wealth, etc.
- Usually, the 80/20 rule applies
  - 80% of the wealth goes to 20% of the people or it takes 80% of the effort to build the easier 20% of the system
  - more examples coming up...

# Stop Words

- Fact:
  - 250-300 most common words in English account for 50% or more of a given text.
- Example:
  - “the” and “of” represent 10% of tokens. “and”, “to”, “a”, and “in” - another 10%. Next 12 words - another 10%.
- Moby Dick Ch.1:
  - 859 unique words (types), 2256 word occurrences (tokens). Top 65 types cover 1132 tokens (> 50%).
- Token/type ratio:
  - $2256/859 = 2.63$

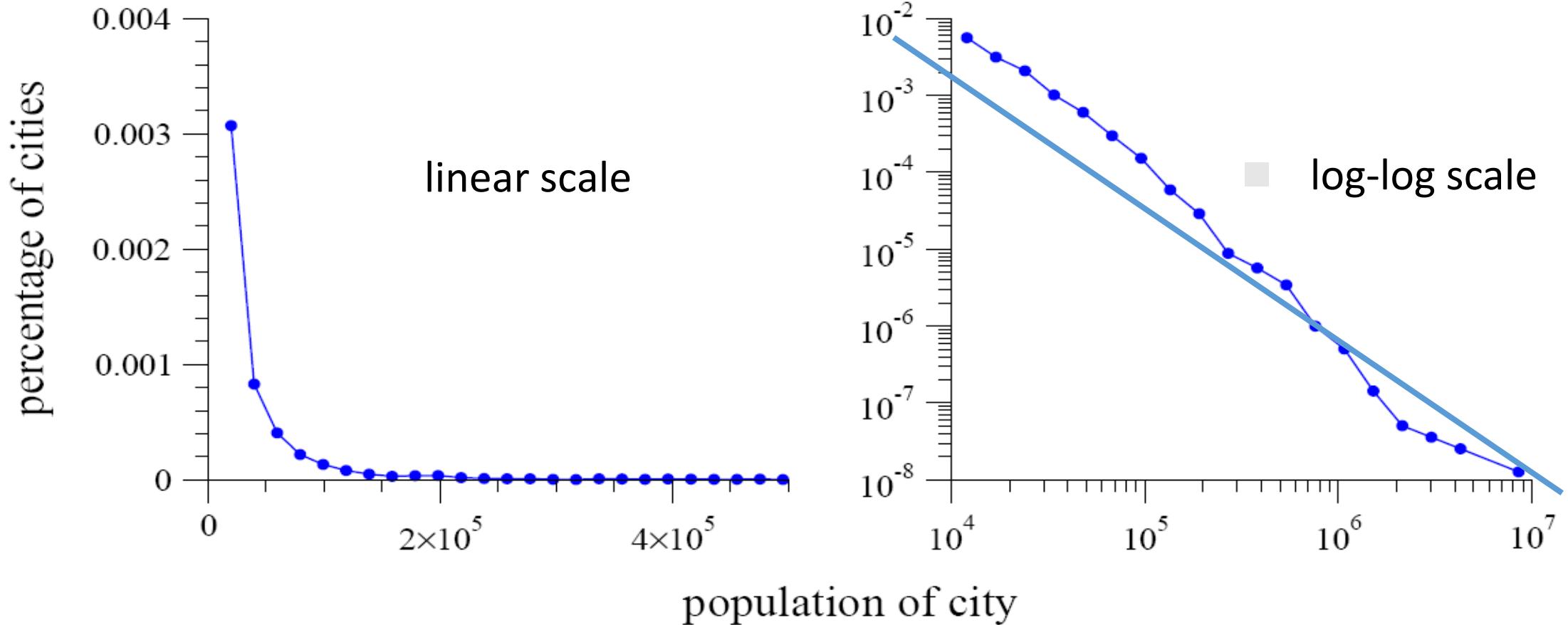
# Power-law Distribution

Percentage of words



- Power-law
  - Many words with a small frequency of occurrence
  - A few words with a very large frequency
  - High skew (asymmetry)
- Compared to a normal distribution:
  - Many people of a medium height
  - Almost nobody of a very high or very low height
  - Symmetry

# Scaling the Axes



■ Long-tail on a linear scale - straight line on a log-log plot

# Power Law Distribution

- The probability of observing an item of size  $x$  is given by

$$p(x) = cx^{-\alpha}$$

normalization constant (probabilities over all  $x$  must sum to 1)

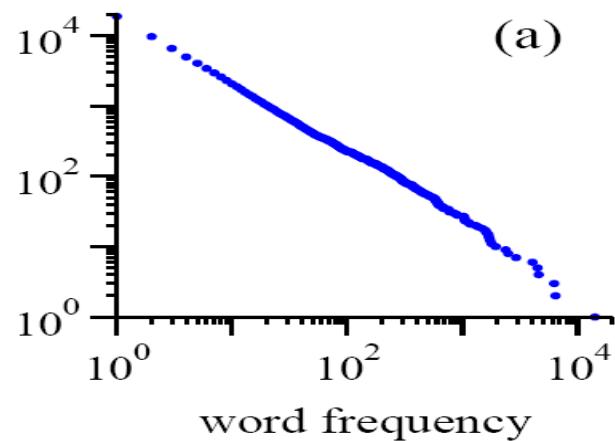
$\alpha$  : scaling exponent, or power law exponent

- Straight line on a log-log plot

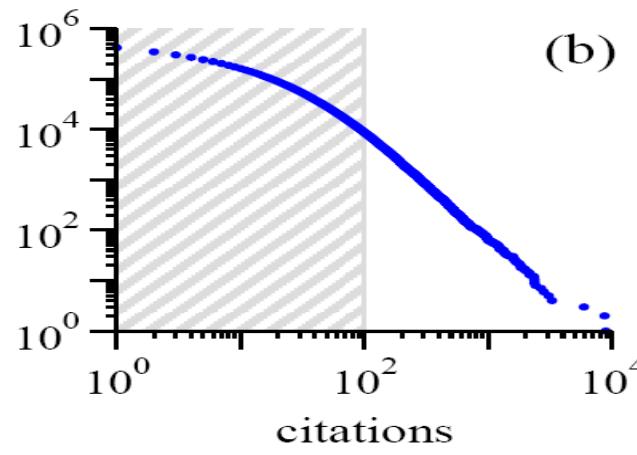
$$\ln(p(x)) = c - \alpha \ln(x)$$

# Power Laws Are Seemingly Everywhere

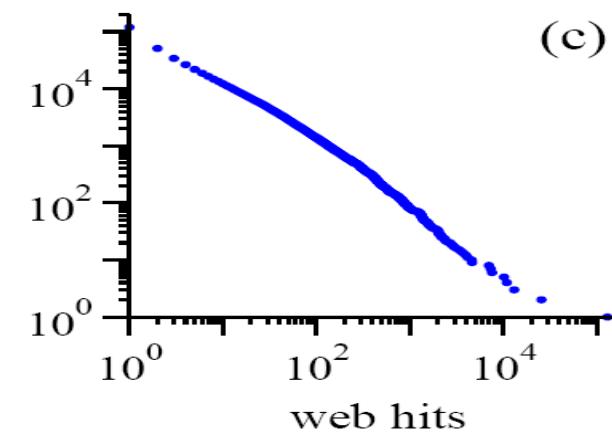
note: these are cumulative distributions



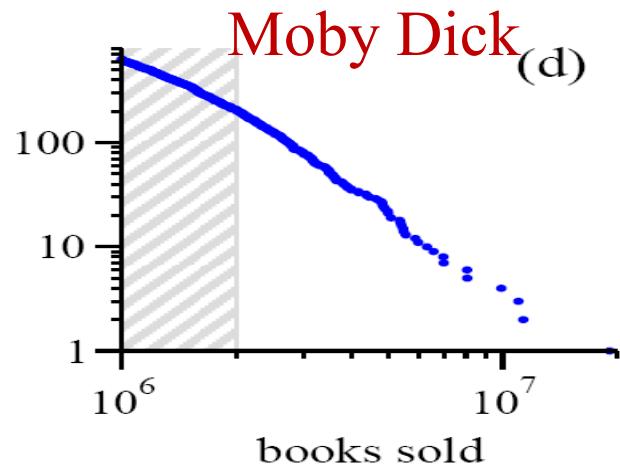
(a)



(b)

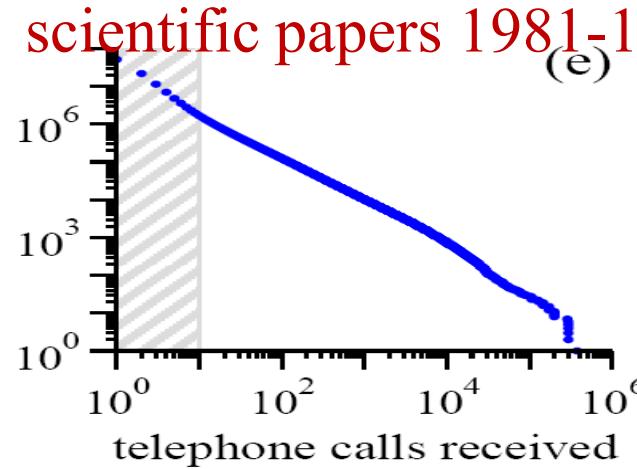


(c)

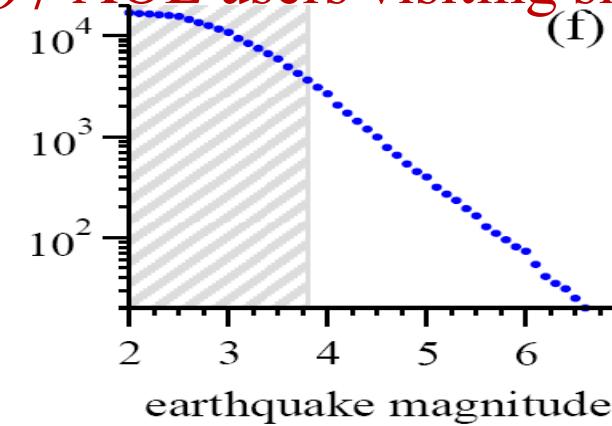


Moby Dick

(d)



(e)



(f)

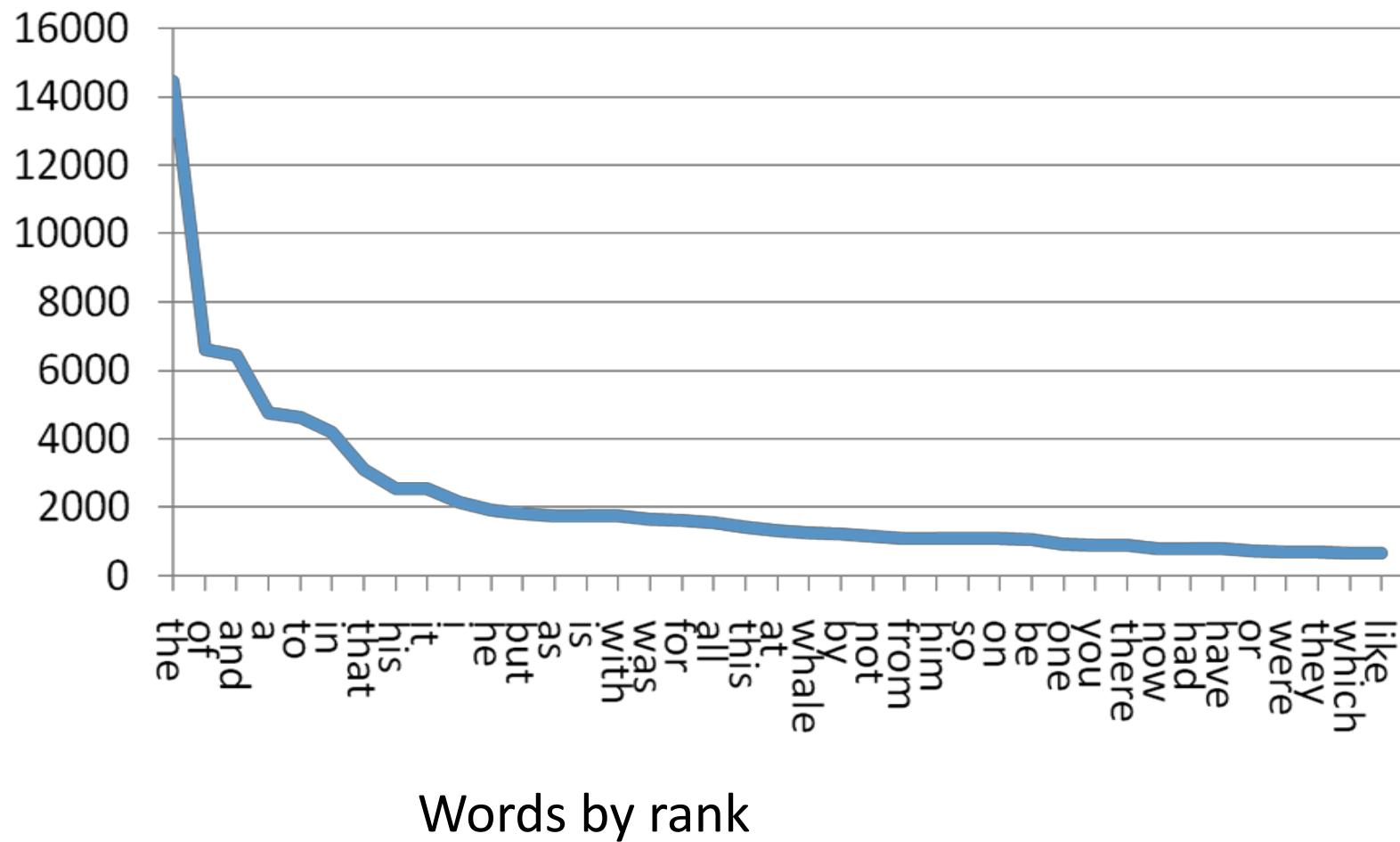
bestsellers 1895-1965

AT&T customers on 1 day

California 1910-1992

# Zipf's Distribution

Word frequency



$$p(k) \sim k^{-\alpha}$$

# Zipf's Law in Natural Language

Rank x Frequency  $\approx$  Constant

- Constant  $\approx 0.1 \times$  Length of collection (in words)
- Not accurate at the tails, but accurate enough for our purposes

Rank	Term	Freq.	Z	Rank	Term	Freq.	Z
1	the	69,9	0.07	6	in	21,3	0.12
2	of	36,4	0.07	7	that	10,5	0.07
3	and	28,8	0.08	8	is	10,0	0.08
4	to	26,1	0.10	9	was	9,81	0.08
5	a	23,2	0.11	10	he	9,54	0.09