



香港城市大學
City University of Hong Kong



南開大學
Nankai University

‘AI’ and AI: Sharing some Progresses from EL&CL team

Qibin Ran

School of Liberal Arts, Nankai University

Jan.20, 2025



History of my academic career

- **Laboratory name: EL&CL(Experimental Linguistics and Computational Linguistics)**
- **What actually we do: Linguistics study using quantitative approaches**
- **Several turns in recent years: Linguistic distance calculation (ASJP approach)**
- **(cross-disciplinary/interdisciplinary studies based on linguistics)**
- **Deep learning/AI(applications in linguistics studies)**

Literature

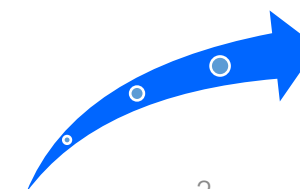
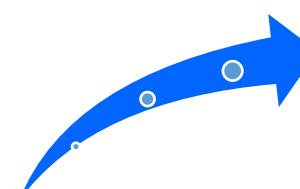
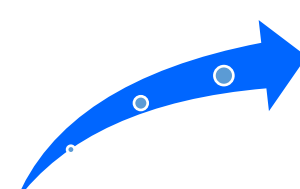
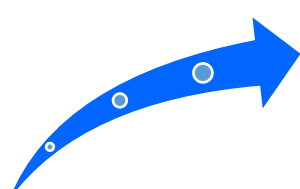
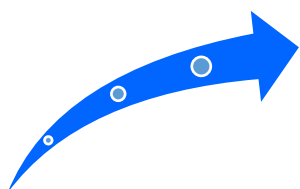
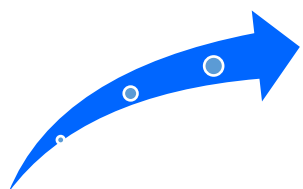
Ancient Chinese

Phonetics

Language distance
calculation

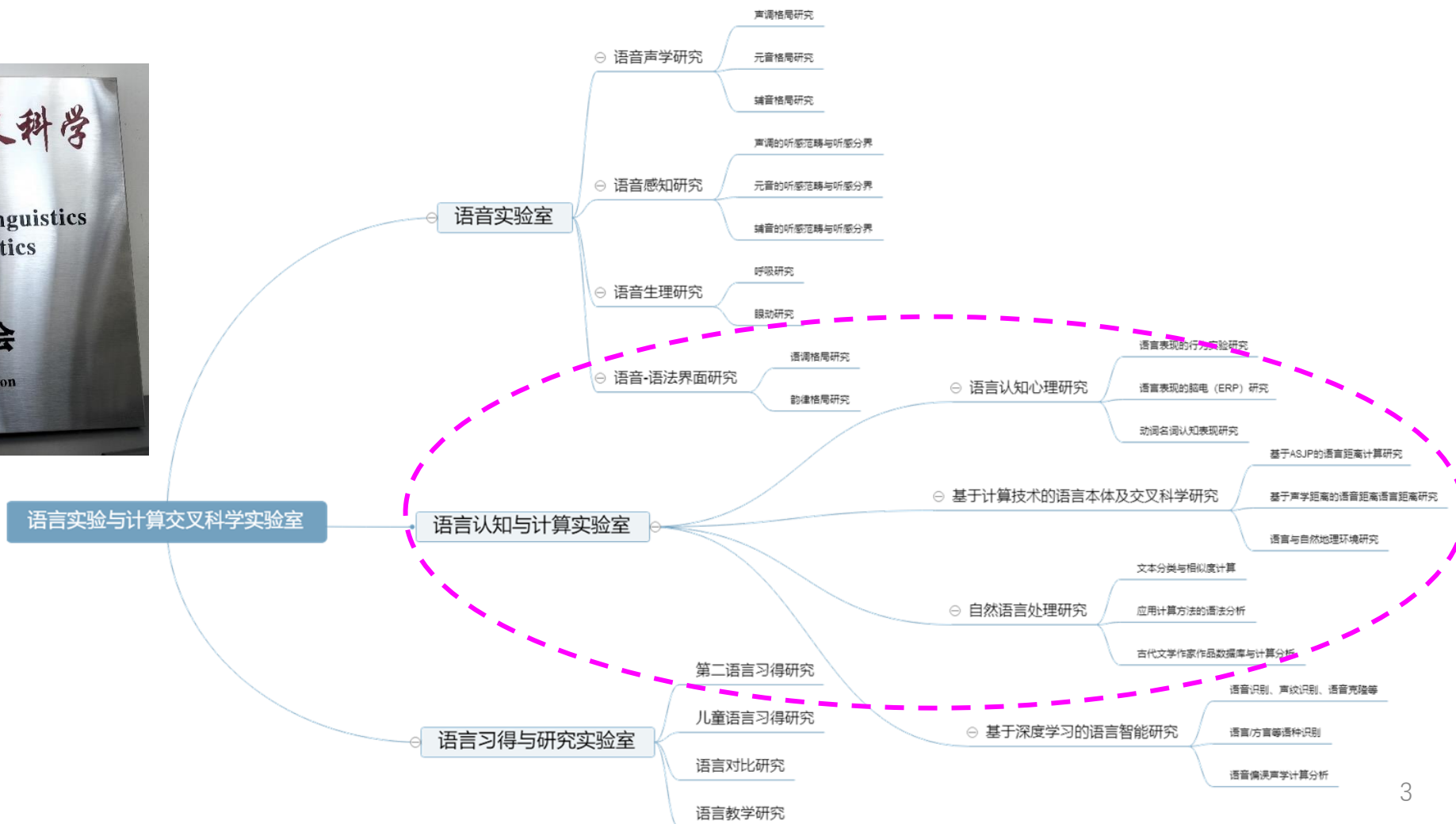
Language and
environment

Deep learning/AI





Title: Social Science Laboratory of Tianjin - Interdisciplinary Laboratory for Experimental Linguistics and Computational Linguistics





- **Social Science Laboratory of Tianjin - Interdisciplinary Laboratory for Experimental Linguistics and Computational Linguistics, Nankai University**
- **History and scope:**
- **Phonetics Laboratory(1986-), established very early in Mainland China**
- **Social Science Laboratory of Tianjin(2019-)**
- **Innovation Team of Tianjin: Language Intelligence calculation(2024-)**
- **The lab is very open and we collaborate with researchers from a variety of subjects/disciplines.**



URL: <https://github.com/EL-CL/>



Previous and ongoing researches-1

- **1-A series of studies with ASJP approach**
- **2-Studies on senses of basic kernel words in a cross-linguistic perspective**
- **3-A serial studies on differentiations between Chinese verbs and nouns**
- **4-A series of studies on acoustic distance calculation**
- **5-Serial studies on relationship between language and environment**



Previous and ongoing researches-1

- **6-Serial Acoustic voice analysis of language/dialects**
- **7-Other studies based on large database/dataset and quantitative approaches**
- **8-Serial studies based on word-vector calculation(Knowledge Graph)**
- **9-Adopting machine learning into linguistic study(classification tasks)**
- **.....**



Previous and ongoing researches-2

- **Ongoing researches using deep learning and AI**
- **1-Text-To-Speech(EL&CL speakers, specific speakers, comic characters, Chinese dialects, chanting for ancient poem...)**
- **2-ASR (IPA recognition, language recognition, foreign accent evaluation, speaker identification...)**
- **3-Machine Translation(Manchu-Mandarin automatic translation)**
- **4- digital human (Yuen-Ren Chao, Chia-ying Yeh, the first president of Nankai Univ.; Tangwangnese, Tianjin dialect)**



I-External differentiations between nouns and verbs: A myth or not?

- **Ran, Qibin (冉启斌) et al.(2023)**
Phonetic differences between nouns and verbs in their typical syntactic positions in a tonal language: Evidence from disyllabic noun–verb ambiguous words in Standard Mandarin Chinese, *Journal of Phonetics*, 98:101241.

- **夏全胜、高凯、冉启斌 (2023)**
汉语动名兼类词做名词和动词时的声学语音差异, 《语言教学与研究》第2期。

Journal of Phonetics 98 (2023) 101241



Research Article

Phonetic differences between nouns and verbs in their typical syntactic positions in a tonal language: Evidence from disyllabic noun–verb ambiguous words in Standard Mandarin Chinese



Qibin Ran ^{a,c}, Kai Gao ^{b,d}, Yuzhu Liang ^{a,c}, Quansheng Xia ^{b,c,*}, Søren Wichmann ^e

^a School of Liberal Arts, Nankai University, Tianjin, China

^b College of Chinese Language and Culture, Nankai University, Tianjin, China

^c Laboratory of Social Science of Tianjin, Tianjin, China

^d Institute of Forensic Science, Ministry of Public Security, China

^e University of Kiel, Germany

ARTICLE INFO

Article history:

Received 8 August 2021

Received in revised form 24 March 2023

Accepted 3 April 2023

Available online 24 April 2023

Keywords:

Noun

Verb

Prosody

Typical syntactic position

Phonetic difference

Information load

Number of syllables

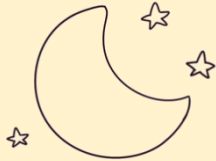

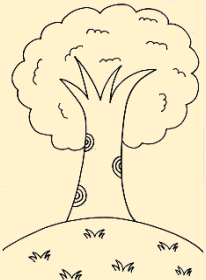
ABSTRACT

This study investigates how word categories, namely noun and verb, influence acoustic realizations (duration, F0, intensity) in Standard Mandarin Chinese, a language having phonemically distinctive tones and a simple morphological system. Noun-verb ambiguous words were selected and presented in the final positions of typical syntactic contexts in order to avoid the interference of prosodic boundary, syntactic complexity, contextual predictability, tonal environment, F0 range and syllable properties (consonant, vowel, tone, syllable length). Linear mixed models were fitted to duration, and generative additive mixed models were fitted to F0 and intensity. The results showed that phonetic differences between nouns and verbs were still evident in duration, F0 and intensity after lexical frequency, speech rate and some other related factors were taken into consideration in the models. The second syllables of nouns were longer than those of verbs, and both syllables of nouns were higher in F0 and greater in intensity than those of verbs. Since the prosodic boundary, frequency and other factors were controlled for, the phonetic differences between nouns and verbs might be attributed to their differences in information load and number of syllables. This study provided evidence that phonetic differences between nouns and verbs might be driven by the grammatical classes themselves and is not an epiphenomenon of other processes.

© 2023 Elsevier Ltd. All rights reserved.

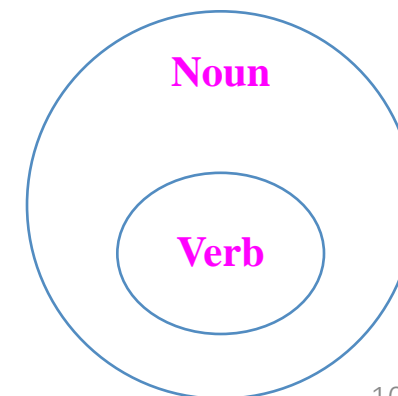
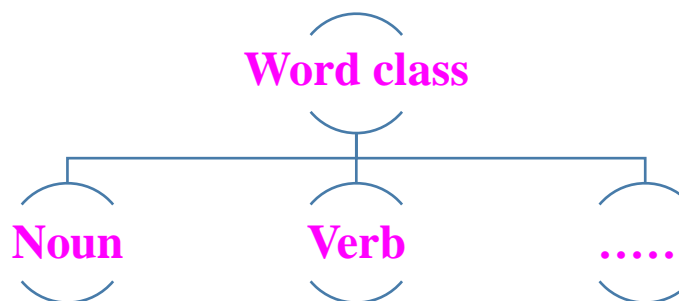
- **Nouns and verbs are different in external form/representations**
- **English:**
- **Suffix:**
- **V.→N.: -ness, -ing, -er...**
- **N.→V.: -ize, -en, -ate...**
- **Stress:**
- **re'cord (V.) : 'record(N.)**
- **??:**
- **lock(V.) : lock(N.)**



noun	verb
	
	
	



- **Most previous studies have shown that nouns tend to be longer in duration than verbs in sentences and discourse (Lohmann & Conwell, 2020; Sorensen, Cooper, & Paccia, 1978; Strunk et al., 2020).**
- **Strunk et al. (2020) showed that, among the ten different languages investigated, the grammatical class effect on duration is evident in eight non-tonal languages while it disappears in the two tonal languages, Bora and Nlɿng. These findings seem to indicate the noun–verb difference in speech is less likely to be realized in duration in tonal languages compared with non-tonal languages.**
- **Standard Mandarin is a typical tonal language with four tones.**
- **Contour tone**
- **Duration, +F0, intensity**





- **Standard Chinese:**
- 锁suo3 lock(parallel with English), 关guan1; ancient Chinese: 言yan2, 雨yu3.....
- **Disyllabic ambiguous words:**
- 爱好 àihào hobby(N)/like(V)
- 包装 bāozhuāng package(N/V)
- 工作gong1zuo4, 暗示an4shi4, 编辑bian1ji, 报告bao4gao4, 辩论bian4lun4, 表现biao3xian4.....



- **Mini pairs:**
- **(1) Noun phrase:**
- 一个+爱好 yīgè + àihào ‘a hobby’
- 一些+爱好 yīxiē + àihào ‘some hobbies’
- 几种+爱好 jǐzhǒng + àihào ‘several types of hobbies’
- **(2) Verb phrase:**
- 能够+爱好 nénggòu + àihào ‘be able to like’
- 应该+爱好 yīnggāi + àihào ‘should like’
- 可以+爱好 kěyǐ + àihào ‘can like’



- 30 mini pairs
- 28 native Mandarin speakers (14 females, 14 males, mean age = 23 yr, SD = 2.2)
- Annotated with Praat: 4 tiers
- pitch, duration, intensity

$$St = 12 * \log_2(F0/F_{ref})$$

$$duration\ ratio = D/\mu_d$$

$$z - score = (I - \mu_{st})/\sigma_{st}$$

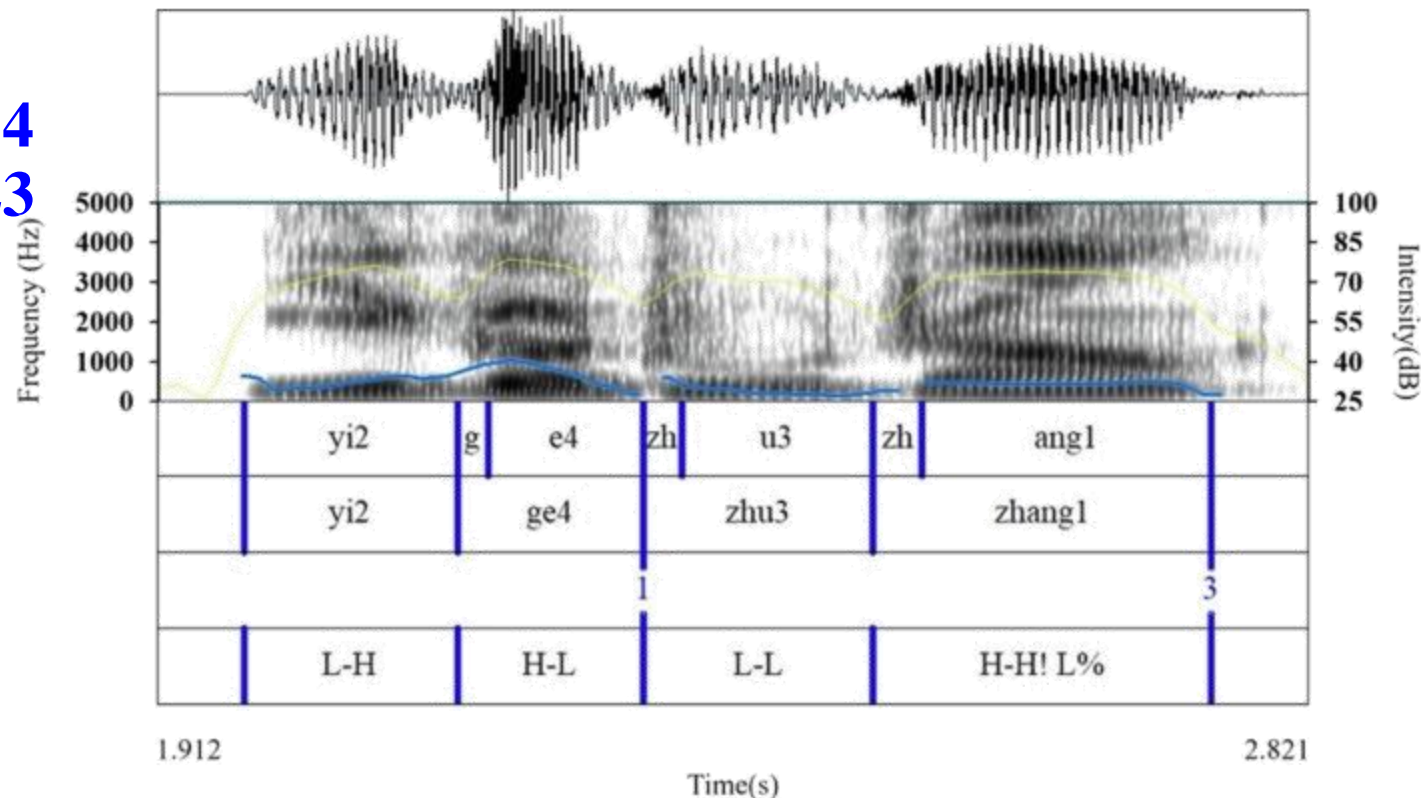


Fig. 1. Example segmentations (the blue and yellow curves displayed illustrate F0 and intensity, respectively).



• **Duration:**

Table 3

Model output of LMM model fitted to the duration (duration ratio) (N = 146160; * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$).

R code: duration ~ category * position + frequency ratio + gender + (1 speaker) + (1 + category word)			
Random effects			
	Variance	SD	
Subject (intercept), 28	0.0054	0.074	
Word (intercept), 29	0.0015	0.038	
Category = Verb	0.0011	0.033	
Residual	0.07	0.265	
Fixed effects			
	Coefficient	Std. Error	t
Intercept	1.242e + 00	2.233e-02	55.627***
Category = Verb	-4.274e-02	7.457e-03	-5.732***
Position = 1	-1.238e-01	1.958e-03	-63.247***
frequency ratio	9.719e-04	4.220e-03	0.23
Gender of speaker (male)	-4.139e-02	2.792e-02	-1.482
Category = Verb by Position = 1	4.872e-02	2.769e-03	17.599***

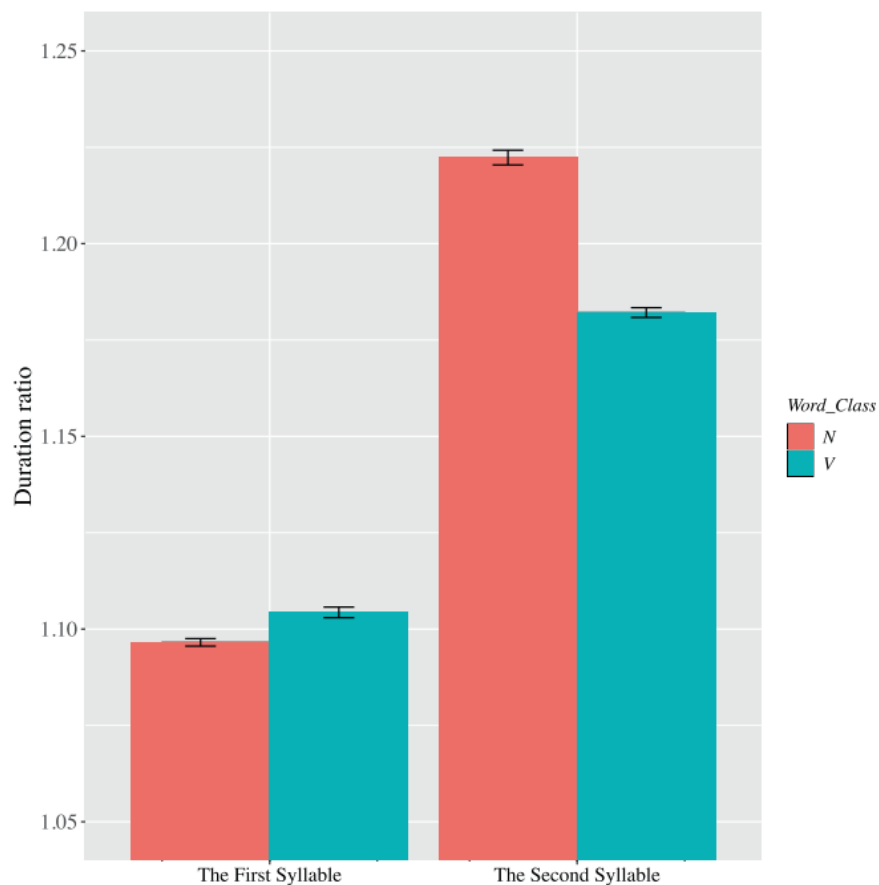


Fig. 2. Duration ratio of nouns (red) and verbs (green). Error bars enclose +/- 1 SE.

- (histogram)

The second syllables of nouns were longer than those of verbs



• Pitch:

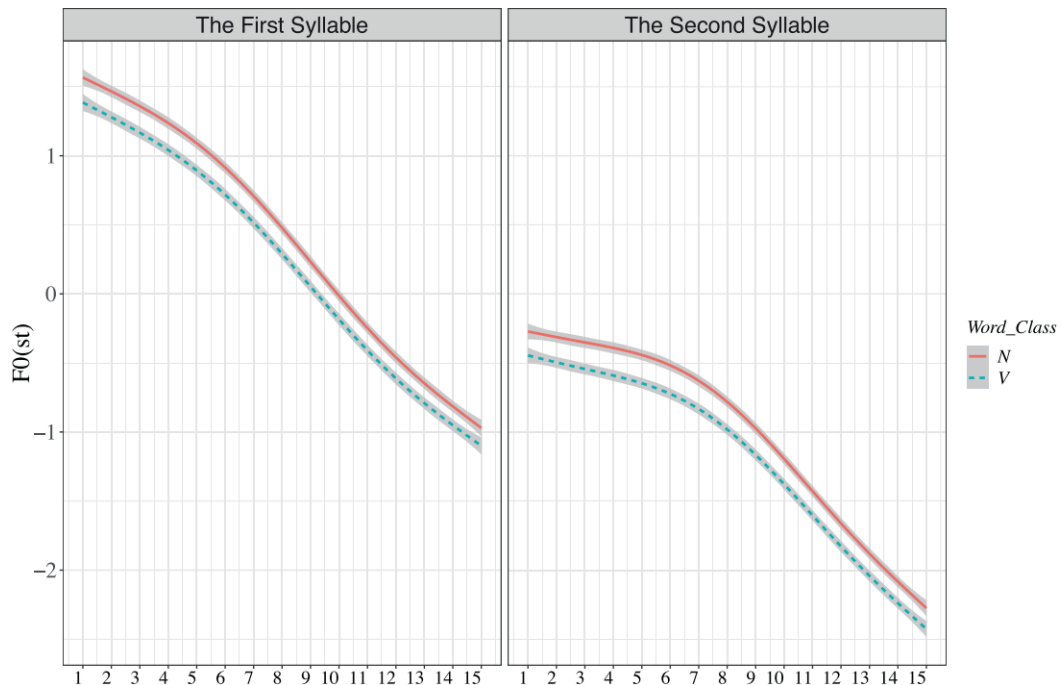


Fig. 3. The GAM smoothers for F0 of nouns (red) and verbs (green). The grey color of F0 contour shows +/- 1 SE.

- Syllables of nouns were higher in F0 than those of verbs

Table 4

Model output of GAMM model fitted to the F0 (semitone) (N = 146160; *p < 0.05, **p < 0.01, ***p < 0.001).

R code: $F0 \sim \text{category} + \text{position} + \text{frequency ratio} + \text{gender} + \text{speech rate} + \text{s}(\text{point}, \text{by} = \text{category}, \text{bs} = \text{"cr"}) + \text{s}(\text{point}, \text{by} = \text{position}, \text{bs} = \text{"cr"}) + \text{s}(\text{tone duration}, \text{bs} = \text{"cr"}) + \text{ti}(\text{point}, \text{tone duration}) + \text{s}(\text{point}, \text{subject}, \text{bs} = \text{"fs"}, \text{xt} = \text{"cr"}, \text{m} = 1, \text{k} = 15) + \text{s}(\text{subject}, \text{position}, \text{bs} = \text{"re"}) + \text{s}(\text{point}, \text{word}, \text{bs} = \text{"fs"}, \text{xt} = \text{"cr"}, \text{m} = 1, \text{k} = 15) + \text{s}(\text{word}, \text{position}, \text{bs} = \text{"re"})$

Parametric coefficients

	Estimate	Std. Error	t
(Intercept)	-0.826	0.308	-2.68**
Category = Verb	-0.063	0.021	-2.968**
Position = 1	1.438	0.382	3.768***
Frequency ratio	-0.109	0.015	-7.474***
Gender of speaker (male)	-0.024	0.109	-0.217
Speech rate	-0.038	0.16	-0.24

Smooth terms

	edf	Ref. df	F
s(Point): Category = Noun	1.808	2.249	26.576***
s(Point): Category = Verb	1.001	1.002	55.795***
s(Point): Position = 2	4.564	5.514	19.118***
s(Point): Position = 1	3.025	3.878	3.165*
s(Tone duration)	8.511	8.904	101.91***
ti(Point, Tone duration)	9.669	11.276	26.111***
s(Point, Subject)	64.909	402	10.611***
s(Subject, Position)	43.433	51	22.013***
s(Point, Word)	103.508	433	431.791
s(Word, Position)	51.181	56	231.24***



• Intensity:

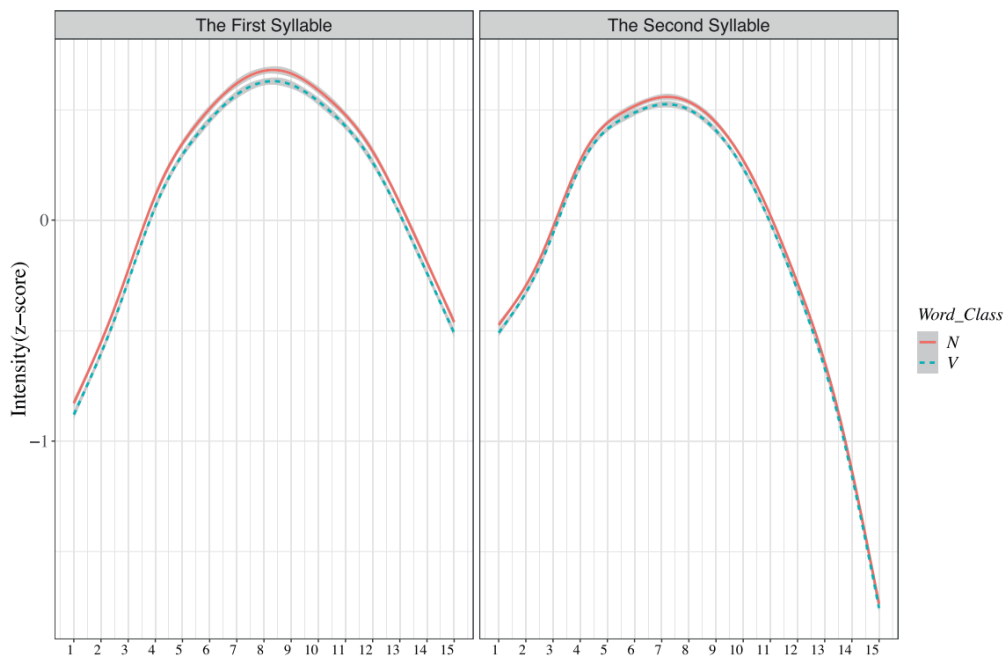


Fig. 4. The GAM smoothers for intensity of nouns (red) and verbs (green). The grey color of intensity contours shows +/- 1 SE.

Table 5

Model output of GAMM model fitted to the intensity (z-score) (N = 146160, $p < 0.05$, $** p < 0.01$, $*** p < 0.001$).

R code: `intensity ~ category + position + frequency ratio + gender + speech rate + s(point, by = position, bs = "cr") + s(duration, bs = "cr") + ti(point, duration) + s(point, subject, bs = "fs", xt = "cr", m = 1, k = 15) + s(subject, position, bs = "re") + s(point, word, bs = "fs", xt = "cr", m = 1, k = 15) + s(word, position, bs = "re")`

Parametric coefficients

	Estimate	Std. Error	t
(Intercept)	0.0331	0.0784	0.422
Category = Verb	-0.0235	0.00429	-5.473***
Position = 2	-0.228	0.0911	-2.5*
Frequency ratio	-0.159	0.00295	-5.386***
Gender of speaker (male)	-0.0113	0.0384	-0.295
Speech rate	0.2104	0.0415	5.067***

Smooth terms

	edf	Ref. df	F
s(Point): Position = 1	7.444	7.734	6.47***
s(Point): Position = 2	8.277	8.454	23.21***
s(Duration)	8.856	8.99	339.21***
ti(Point, Duration)	15.407	15.943	58.1***
s(Point, Subject)	203.312	417	219.32***
s(Subject, Position)	49.606	53	86.33***
s(Point, Word)	283.48	433	2498.28*
s(Word, Position)	47.72	56	292.35***

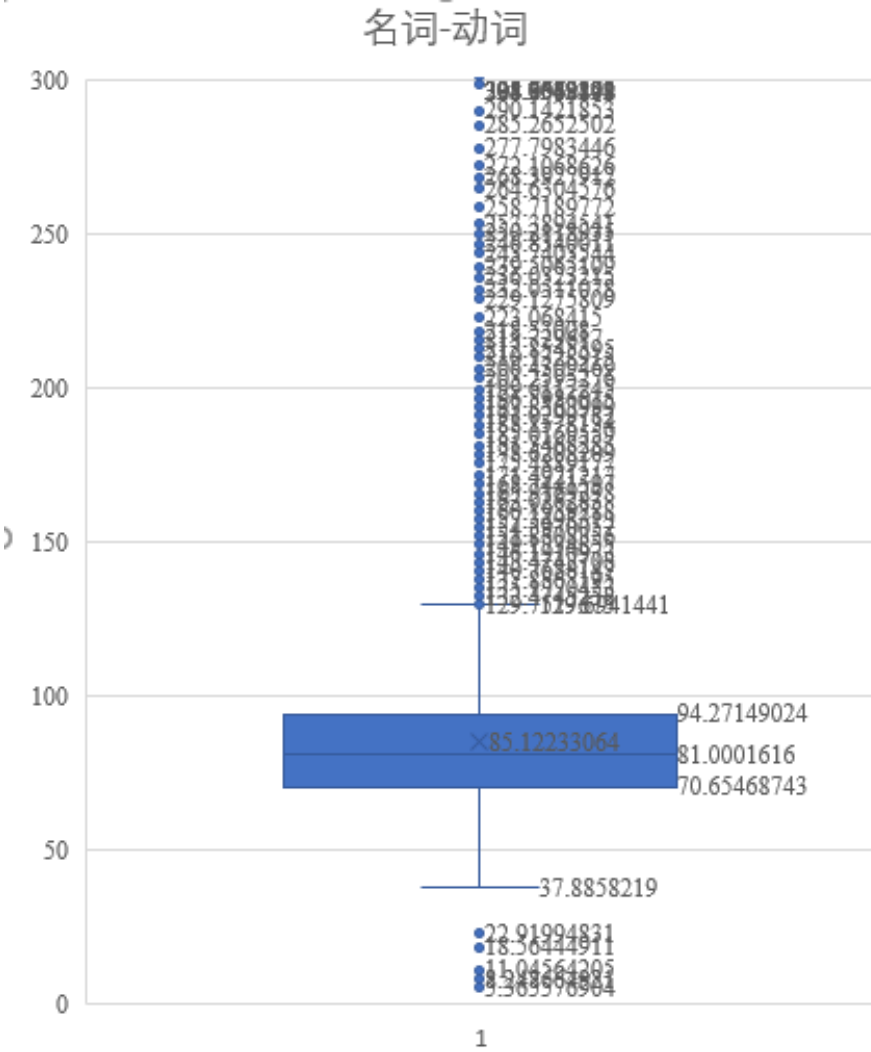
• Syllables of nouns were greater in intensity than those of verbs



- **Explanation:**
- **Strunk et al. (2020): nouns are longer than verbs in 8 non-tonal languages**
- **(1)Information load: In Mandarin Chinese, nouns are usually replaced by pronouns or omitted when mentioned a second time. Thus, nouns and verbs in Mandarin should be considered to be inherently different in information load, leading to longer duration of nouns than verbs.**
- **(2)the Smooth Signal Redundancy (SSR) hypothesis (Aylett & Turk, 2004; Turk, 2010): the more predictable or redundant a linguistic element is, the less salient its phonetic realization will be.**
- **(3)When the two grammatical classes are embedded in sentences, nouns need more planning time (Seifart et al., 2018)and are pronounced slower than verbs (Strunk et al., 2020).**
- **(4)Number of syllables: Nouns are typically disyllabic (Ke, 2012; Liu, 1996), while verbs are typically monosyllabic (Chen, 1987; Liu, 1996; Zhang, 1989).**



- Acoustic distance calculation between nouns and verbs
- audio1~audio2
- 20 words*6 times*30 speakers=3,600
- It is at a nearly random variations level.
- There isn't essential difference between nouns and verbs??



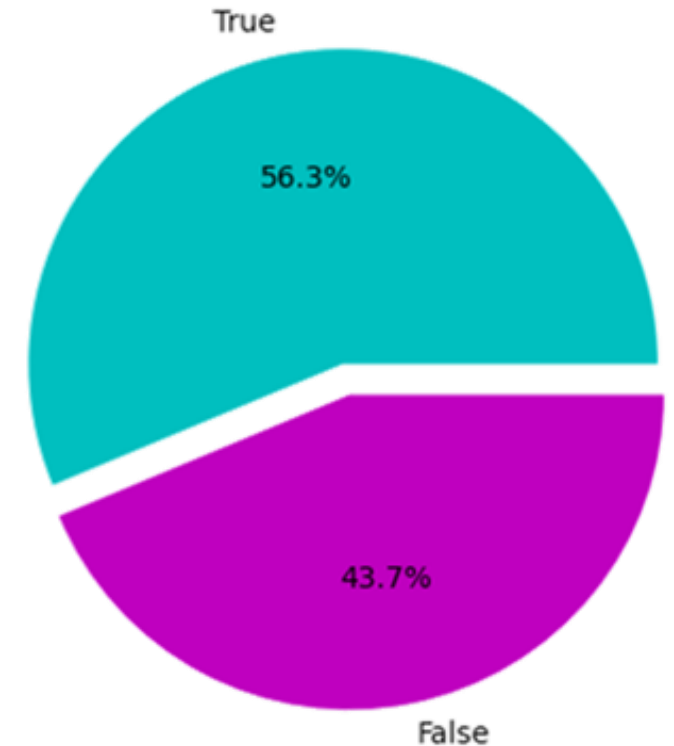


• The hierarchy of acoustic distance between different sounds (preliminary)

Category	distance (mean)	备注
Sandhi tone3/tone2	65.71	20个词对 (仁者-忍者) , 21名发音人
noun/verb	81.00 (双音节)	20对动词-名词 (爱好) , 30名发音人
4 tones	86.79	四声俱全音节, 4名发音人
n-l	88.03	55对音节, 6名发音人
Neutralized/non-N..d	109.30	14个词对 (东西) , 10名发音人
rhoticized/non-R..d	126.03 (双音节) 、 144.10 (单音节)	单、双音节各20对, 10名发音人
initials	127.57	声母均值, 4名发音人
finals	127.84	韵母均值, 4名发音人
vowel	131.05	8个单元音, 现代女发音人
inter speaker	158.08	10名发音人, 300核心词
Swadesh list1	165.87	207核心词, 69种语言
Swadesh list2	206.76	160核心词, 50种语言
between languages	215.79	44种语言, 各300核心词
between animals	334.68	175种动物, 各3种发音



- **Using machine learning to differentiate nouns and verbs**
- **3 acoustic parameters: pitch, duration, intensity**
- **151,200 samples**
- **56~57%:43~42%**
- **A little bit higher than probability level**
- **So, is there any differences between nouns and verbs??**
- **There isn't essential difference between nouns and verbs?**
- **The external differentiations between nouns and verbs: A myth or not?**





II- Studies on a dataset of cross linguistic sense for kernel words

- Yuzhu Liang (梁煜珠), Ke Xu (许可), Qibin Ran (冉启斌), (2024). Shared structure of fundamental human experience revealed by polysemy network of basic vocabularies across languages. *Scientific Reports* 14: 5877.
- 许可、冉启斌 (2024) 基本核心词词义的跨语言表现与特点——对61种语言数据的分析, 《世界汉语教学》第1期
- 许可、冉启斌、李坤怡 (2022) 基本核心词词义的数量及其内部关系——对50种语言词义数据的分析探索, 《中国语言学研究》第二辑, 北京: 中国社会科学出版社
- 许可、李坤怡、冉启斌、黄玮 (2022) 跨语言基本核心词词义分析——以24种语言基本核心词词义为例, 《语文学刊》第5期

www.nature.com/scientificreports

scientific reports

Check for updates

OPEN Shared structure of fundamental human experience revealed by polysemy network of basic vocabularies across languages

Yuzhu Liang^{1,3}, Ke Xu^{1,3} & Qibin Ran^{1,2}✉

How are concepts related to fundamental human experiences organized within the human mind? Our insights are drawn from a semantic network created using the Cross-Linguistic Database of Polysemous Basic Vocabulary, which focuses on a broad range of senses extracted from dictionary entries. The database covers 60 basic vocabularies in 61 languages, providing 11,841 senses from 3736 entries, revealing cross-linguistic semantic connections through automatically generated weighted semantic maps. The network comprises 2941 nodes connected by 3573 edges. The nodes representing body parts, motions, and features closely related to human experience occupy wide fields or serve as crucial bridges across semantic domains in the network. The polysemous network of basic vocabularies across languages represents a shared cognitive network of fundamental human experiences, as these semantic connections should be conceived as generally independent of any specific language and are driven by universal characteristics of the real world as perceived by the human mind. The database holds the potential to contribute to research aimed at unraveling the nature of cognitive proximity.

• Basic concepts and their derived senses

huǒ (火)

火 huǒ ① (～儿) 名 物体燃烧时所发的光和焰：～光|～花|灯～|点～|～越烧越旺。② 指枪炮弹药：～器|～力|～网|军～|走～。③ 名 火气③：上～|败～。④ 形容红色：～鸡|～腿。⑤ 比喻紧急：～速|～急。⑥ (～儿) 名 怒气：冒～|心头～起。⑦ (～儿) 动 比喻发怒：～性|他～儿了。⑧ 〈口〉 形 兴旺；兴隆：买卖很～。⑨ 同“伙¹、伙²”。⑩ (Huǒ) 名 姓。

fire  (fɪr)

n.

1.
 - a. A rapid, persistent chemical change that releases heat and light and is accompanied by flame, especially the exothermic oxidation of a combustible substance: *destruction by fire*.
 - b. A specific instance of this change that destroys something: *a house fire*.
 - c. A burning fuel: *a cooking fire*.
2. Burning intensity of feeling; ardor or enthusiasm: *a musical performance that had fire*. See Synonyms at **passion**.
3. Luminosity or brilliance, as of a cut and polished gemstone.
4. Liveliness and vivacity of imagination; brilliance: *the fire of an artistic genius*.
5. A severe test; a trial or torment: *went through fire to become a leader*.
6. A fever or bodily inflammation: *tormented by the fire in an infected toe*.
7.
 - a. The discharge of firearms or artillery: *heard the fire of cannon*.
 - b. The launching of a missile, rocket, or similar ballistic body.
 - c. Discharged bullets or other projectiles: *subjected enemy positions to heavy mortar fire; struck by rifle fire*.
8. Intense, repeated attack or criticism: *answered the fire from her political critics*.

v. fired, fir-ing, fires



一 1 yī ① 一 最小的正整数。参看 1212 页
【数字】。② 一 表示同一：咱们是～家人|
你们～路走|这不是～码事。③ 一 表示另一：
番茄～名西红柿。④ 一 表示整个；全：～冬|
～生|～路平安|～屋子人|～身的汗。⑤ 表示
专一：～心～意。⑥ 一 表示动作是一次，或
表示动作是短暂的，或表示动作是试试的。a)
用在重叠的动词(多为单音)中间：歇～歇|笑
～笑|让我闻～闻。b)用在动词之后，动量词
之前：笑～声|看～眼|让我们商量～下。⑦ 一
用在动词或动量词前面，表示先做某个动作
(下文说明动作结果)：～跳跳了过去|～脚把
它踢开|他在旁边～站，再也不说什么。⑧ 一
与“就”配合，表示两个动作紧接着发生：～请
就来|～说就明白了。⑨ 一旦；一经：～失足
成千古恨。⑩ 〈书〉一 用在某些词前加强语
气：～何速也|为害之甚，～至于此！ || 注意
“一”字单用或在一词一句末尾念阴平，如“十
一、一一得一”，在去声字前念阳平，如“一半、
一共”，在阴平、阳平、上声字前念去声，如“一
天、一年、一点”。本词典为简便起见，条目中
的“一”字，都注阴平。

one (wǔn)

adj.

1. Being a single entity, unit, object, or living being: *I ate one peach.*
2. Characterized by unity; undivided: *They spoke with one voice.*
3.
 - a. Of the same kind or quality: *two animals of one species.*
 - b. Forming a single entity of two or more components: *three chemicals combining into one solution.*
4. Being a single member or element of a group, category, or kind: *I'm just one player on the team.*
5. Being a single thing in contrast with or relation to another or others of its kind: *One day is just like the next.*
6. Occurring or existing as something indefinite, as in time or position: *He will come one day.*
7. Occurring or existing as something particular but unspecified, as in time past: *late one evening.*
8. Informal Used as an intensive: *That is one fine dog.*
9. Being the only individual of a specified or implied kind: *the one person I could marry; the one horse that can win this race.*

n.

1. The cardinal number, represented by the symbol 1, designating the first unit in a series.
2. A single person or thing; a unit: *This is the one I like best. Of her many books, the best ones are the last two.*
3. A one-dollar bill.

pron.

1. An indefinitely specified individual: *She visited one of her cousins.*
2. An unspecified individual; anyone: *"The older one grows the more one likes indecency" (Virginia Woolf).*



• The differences and similarities in the derived senses of basic concepts

stone	shi2 (石)
stone	stone
a piece of stone shaped for a purpose	stone inscription
a hard seed in some fruits	
a unit of weight	
a natural shade of whitish or brownish grey	

go	zou3(走)
English	Chinese
go	go
leave	run
intend to be or do something	move or work
pass into a specified state	become
turn out	leave
be harmonious, complementary, or matching	die
function	visit
be regularly kept or put in a particular place	pass through
use a toilet	leak
	change

listen	ting1(听)
English	Chinese
listen	listen
accept one's opinion	accept one's opinion
pay attention to	at one's convenience
	govern or judge

- **Deriving directions and sense network**



- **If we collect enough senses from a vast amount of languages, we can build a large cross-linguistic sense network.**



- **Investigations:**
- **How many derived senses for different basic concepts in a cross-linguistic perspective?**
- **What are the universals and differences of deriving directions in a large dataset?**
- **Do the derived senses reflect the cognitive significance of different languages?**
- **Are there differences of deriving depth in different languages?**
- **What are the most prominent sense units(including basic senses and derived senses from a cross-linguistic perspective?)**
- **What are the most prominent connections in the big network of sense units?**
- **.....**



- **Before we started the research, we reviewed a vast amount of literature.**
- **Talmy (1985,2000); Anderson(1978)words for body parts; Comrie(2011) numeral words**
- **伍铁平 (1989) 、蒋栋元 (2002) 、温凌云 (2007) 、蒋绍愚 (2008) 、黄树先 (2012)**
- **Max-Plank Institute: Wilkins(1996), Fortson IV(2004), Goddard & Wierzbicka(2008, Newman(2009, Urban(2010, Pericliev(2015).....**
- **Miller et al., WordNet (Miller 1995, Fellbaum 2012)**
- **<https://wordnet.princeton.edu/>; <http://lope.linguistics.ntu.edu.tw/cwn/>;
<http://compling.hss.ntu.edu.sg/cow/>**
- **Chinese Concept Dictionary, CCD, 2001**
- **List, Mayer, Terhalle & Urban, CLICS (Rzymiski et al. 2020): colexification: <https://www.clips.com/>**
- **These studies can not fully meet our requirements.**



60 basic kernel words(selected from Swadesh list)

Part of speech	Sense category	word
noun (25)	Natural objects	water、fire、stone、tree、earth/soil、rain、sky、sun、moon、night
	Living being	human、bird、dog、person、man、woman
	Body parts	head、eye、ear、mouth、hand、foot、heart、bone、blood
verb (11)	Action/motion	eat、drink、say、see/look、hear/listen、stand、sleep、give、go、come、die
adj. (14↑)	Quantity/size	many、long、big、small
	Property/quality	new、old、good、bad
	Color	black、white
	Shape	square、round
	Temperature	cold、hot
pro-noun (6↑)	Personal pronoun	I、you
	demonstrative pronoun	this、that
	interrogative pronoun	who、what
num. (3)	Numeral	one、two、three
adv. (1)	negator	not



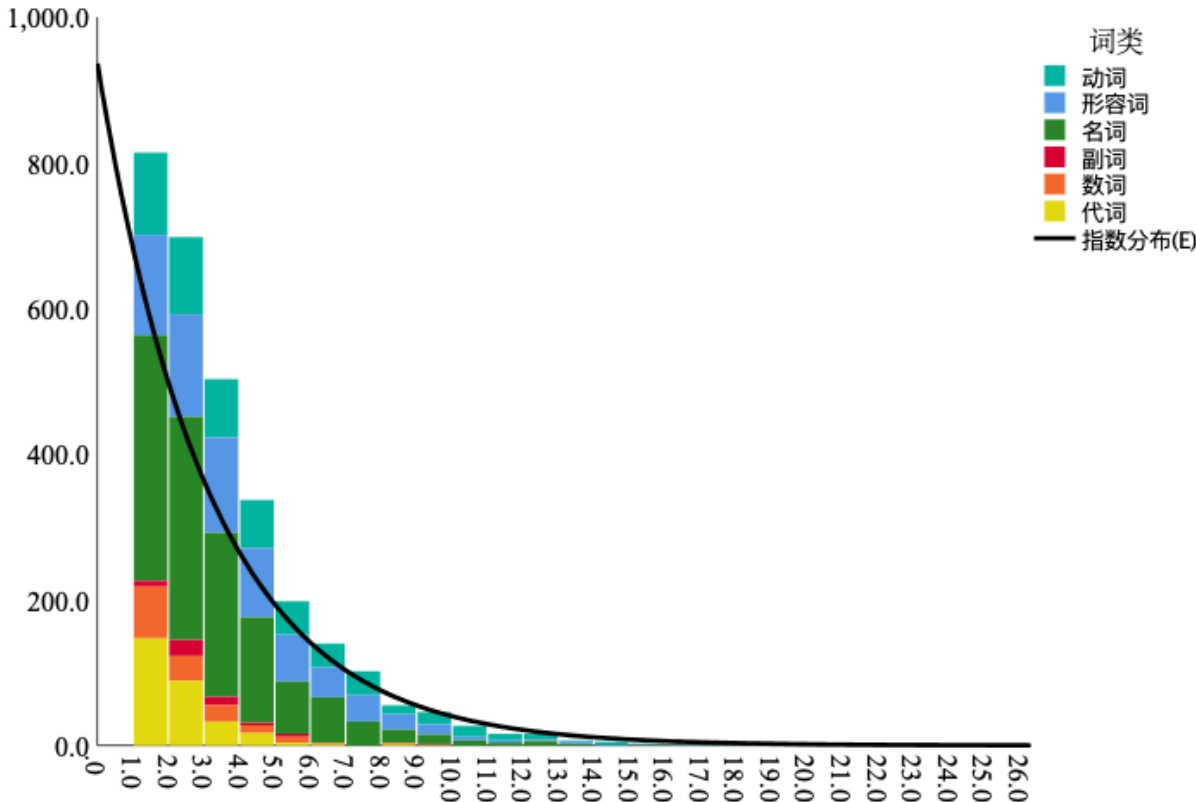
- about 5 years
- 61 languages
- Medium sized authoritative dictionaries
- URL:
https://github.com/ELL-CL/CLD_Polysemous_Basic_Vocabulary-

A	B	C	D	E	F	G	H
序号	整理时间	language	语言 (汉语拼音首序)	词典	作者	出版社	出版时间
1	20190614	Andric	阿拉伯语	高级阿拉伯词典			
2	20190614	Irish	爱尔兰语	Irish-English Dictionary with	O'Reilly, Edward, and John O'Donnell	James Duffy	
3	20190614	Polish	波兰语	Collins Polish-English Dictio	www.collinslanguage.com		
4	20190614	Persian	波斯语				
5	20190614	German	德语	Oxford Duden German Diction	J.B.Sykes, Werner Scholze-Stubrecht		2008
6	20190614	Russian	俄语	现代俄汉双解词典		外语教学与研究出版社	1992
7	20190614	French	法语	FRENCH-ENGLISH Dictionary	Barron's Educational Series Inc., U.S.		2016
8	20190614	Finnish	芬兰语	bab.la Finnish-English diction			
9	20190614	Korean	韩语	DONG-A'S PRIME KOREAN- 영어 (斗山东亚)			1998
10	20190614	Chinese	汉语	现代汉语词典 (第7版)	商务印书馆		2016
11	20190614	Czech	捷克语	Velký česko-anglický slovník (英文翻译: Big Czech-English Dicti		Leda	2013
12	20190614	Latin	拉丁语	Collins Latin Dictionary & Gr	世界图书出版公司北京公司		2013
13	20190614	Manchu	满语				
14	20190614	Burmese	缅甸语	缅甸词典	北京大学东方语言文学系缅甸语教	商务印书馆	
15	20190616	Norwegian	挪威语	Norwegian-English dictionary		Oslo : Universitetsforlaget	1990
16	20190614	Portuguese	葡萄牙语	Oxford Portuguese Dictionary	Sinda López Fuentes & Ana Franke	Oxford University Press	
17	20190614	Japanese	日语	A STANDARD JAPANESE - English Dictionary	TSUNETA TAKEHARA M. A. Ph. D.	东京大阪 广文堂	
18	20190614	Swedish	瑞典语	The Standard Swedish-English	Cassell Publishers Ltd Villiers Hou	首次出版者: Holt Rinehart & Winston	
19	20190614	Slovak	斯洛伐克语				
20	20190615	Swahili	斯瓦西里语	A standard Swahili-English dic	Inter-territorial Language (Swahili) C	Nairobi, Kenya : Oxford	1939
21	20190614	Tamil	泰米尔	Tamil-English, English-Tamil (Victor	Hippocrene Books	2004
22	20190614	Thai	泰语	THAI-ENGLISH DICTIONARY			2003
23	20190614	Uighur	维吾尔语				
24	20190614	Hungarian	匈牙利语	Hungarian-English Dictionary		匈牙利科学院	
25	20190614	Spanish	西班牙语	Oxford Spanish Dictionary	Galimberti Janner, Roy Russell, Car	OLIP Oxford, 4	2008
26	20190615	Hebrew	希伯来语				
27	20190614	Greek	希腊语	Collins Greek-English Diction	HaperCollins Publishers		2003
28	20190614	Hawaiian	夏威夷语	Hawaiian Dictionary	Mary Kawena Pakui, Samuel H. Elbr	University of Hawai'i Pre	1986
29	20190614	Italian	意大利语	Oxford-Paravia Italian Dictionary		Oxford University Press	2010
30	20190614	Hindi	印地语	Hindi-English English-Hindi Dictionary		HIPPOCRENE BOOKS, INC.171 Madison	
31	20190614	Indonesian	印度尼西亚语	A Comprehensive INDONESIA- English Dictionary	Alan M. Stevens and A. Ed. Schmid	Ohio University Press	
32	20190614	English	英语	Oxford Intermediate Learner's	商务印书馆		2010
33	20190614	Vietnamese	越南语	现代越汉词典 (第2版)	曹航 李宝红 主编	外语教学与研究出版社	
34	20190820	Scottish Gaelic	苏格兰盖尔语	The Gaelic-English dictionary	Colin Mack	London ; New York : Ri	2004
35	20200903	Turkish	土耳其语	土耳其语汉语词典	周正清 周廷立 主编	商务印书馆	2008
36	20200903	Mongolian	蒙古语	新蒙汉词典	《新蒙汉词典》编委会编	商务印书馆	1999
37	20200926	Malay	马来语	Kamus umum bahasa malaysia	Yang Kail Yee,Chan Meow Wah	The World Book Co.Ssd	1988
38	20201020	Urdu	乌尔都语	乌尔都语汉语词典	孔菊兰 主编	高等教育出版社	2014
39	20201009	Lingao	菲商话	菲商汉语词典	中国社会科学院民族研究所主编 刘	四川民族出版社	2000
40	20201020	Zhuang Language	壮语	壮汉词汇	广西壮族自治区少数民族语言文字	广西民族出版社	1984
41	20201027	Bouyei	布依语	布依汉语词典	中国社会科学院民族研究所主编 吴	民族出版社	2002
42	20201104	Li	黎语	黎汉词典	中国社会科学院民族研究所主编 郑	四川民族出版社	1993
43	20201208	Khmer	高棉语	Cambodian-English Dictionary	Robert K. Headley, Ruth Chim, and	Darwoody Press	1997
44	20201216	Lao	老挝语	Lao-English Dictionary	William L. Patterson and Mario E. S	Darwoody Press	1995
45	20201212	Lithuanian	立陶宛语	Lithuanian dictionary	Bronius Picasasas, Bronius Svecevi	Routledge	1995
46	20201220	Dutch	荷兰语	Prisma Handwoordenboek Eng	Pruc Gargano, Fokko Veldman	Uitgeverij Unieboek	2010
47	20201220	Hausa	豪萨语	A Hausa-English Dictionary	Paul Newman	Yale University Press	2007
48	20210323	Malagasy	马尔加什语	A new Malagasy-English dictio	J. Richardson	The London Missionary	1885
49	20210107	Danish	丹麦语	Danish Dictionary-Danish-Eng	Arna Guld	Routledge	1995
50	20210602	Kyrgyz	柯尔克孜语	Kyrgyz-English Dictionary	Karl A. Krippes	Darwoody Press	1998
51	20210616	Basque	巴斯克语	CBS-Momis English-Basque/ E	Momis Academy Press	Mikel Momis	2010
52	20210609	Teluga	泰卢固语	A TELUGU-ENGLISH DI	New York: Oxford University Press	Gwynn, J. P. L. (John P	1991
53	20210719	icelandic	冰岛语	Íslensk-ensk orðabók / Con	Sverrir Hólmarsósson; Sanders, Ch	Bunn	1989
54	20210610	Amuzian	亚美尼亚语	A Comprehensive Dictionary /	Mesrob G. Kouyoumdjian	Atlas Press	1970
55	20210613	Croatian	克罗地亚语	Veliki hrvatsko-engliski rječn	Z.Bujas	Globus	2008
56	20210726	Bengali	孟加拉语	Samsad Bengali-English dic	Biwas, Sallendra	Sahitya Samsad	2000
57	20210628	maori	毛利语	A Dictionary of the Maori Lan	Herbert W. Williams	Government Printer	1957
58	20210906	Pashto	普什图语	A Dictionary of the Pashto	Raverty, H. G.	Literary Licensing, LLC	1867
59	20210906	Shan	掸族语	http://scslang.org/shan/dictionary.htm			
60	20210906	Pali	巴利语	http://www.ahandiko.files.wordpress.com/documents/Concise%20Pali%20English%20Dictionary_Buddhadatta.pdf			
61	20210906	Marshallese	马绍尔语	Marshallese-English Dictionary	Takaji Aho, Byron W. Bender, Alfr	University of Hawai'i Pre	1977



Statistics of the dataset

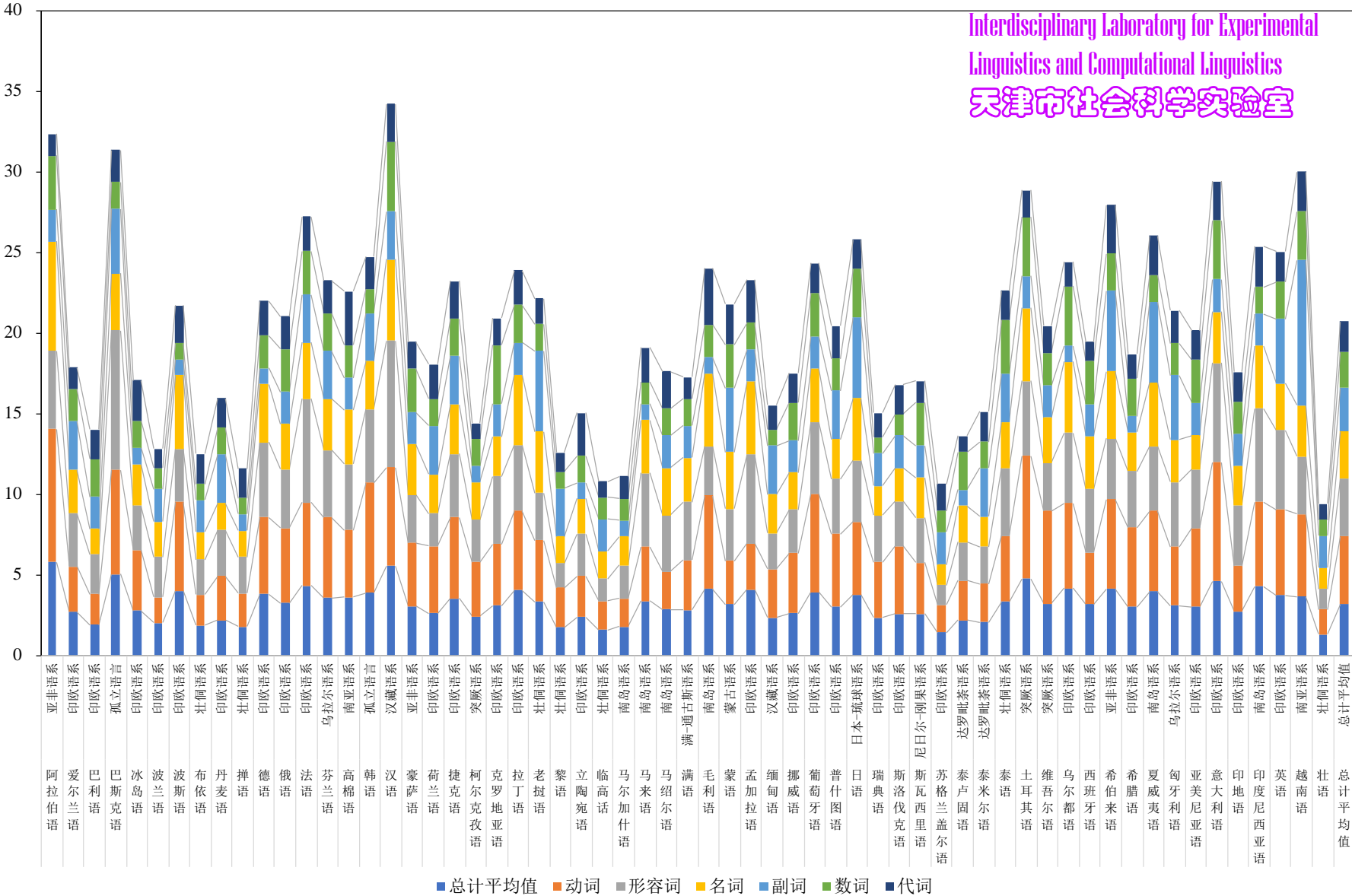
- The minimum number of senses: 1;
the maximum number of senses:
Arabic: eye, 22
- the maximum of sense in different
word classes:
- Verbs: 20, adjective:15, adv.(not):
9, pronoun:9, numeral:8
- 义项数量分布广，但多数集中在
较低值上
- exponential distribution



Stacked histogram of sense number for different word classes



词类、语言对义项数量均值有显著影响；不同语言中不同词类的义项数量大小可能不同 (F=69.702, df=5, p<0.001; F=3.022, df=60, p<0.001; 交互作用: F=1.32, df=300, p<0.001)



Stacked histogram of sense number for different word classes



- **General descending order: verb > adjective > noun > adverb > numeral > pronoun**
- **Chinese: adjective > verb > noun > other classes. (including Indonesian, Hindi, etc.)**
- **Korean: verb > adjective > noun > other classes. (Japanese, Greece, etc.)**
- **Hawaiian, Latin: verb > noun > adjective > other classes.**
- **Thai, Burmese, Irish, Tamil: almost equally distributed in different word classes.**
- **Pronoun**
- **Lithuanian: I:1, you:1, this:1, that:2, who:8, what:3**
- **Who** → “whoever, anybody (无论谁, 任何人)”、“someone (某人)”、“something (某事)”、“anything (任何事)”、what、that、“the thing or things that (.....的事)”
- **negative adverbs**
- **Japanese, Hawaiian, Vietnamese: not → ‘without’ ‘never’ ‘avoid’ ‘empty’ ‘gratis’**

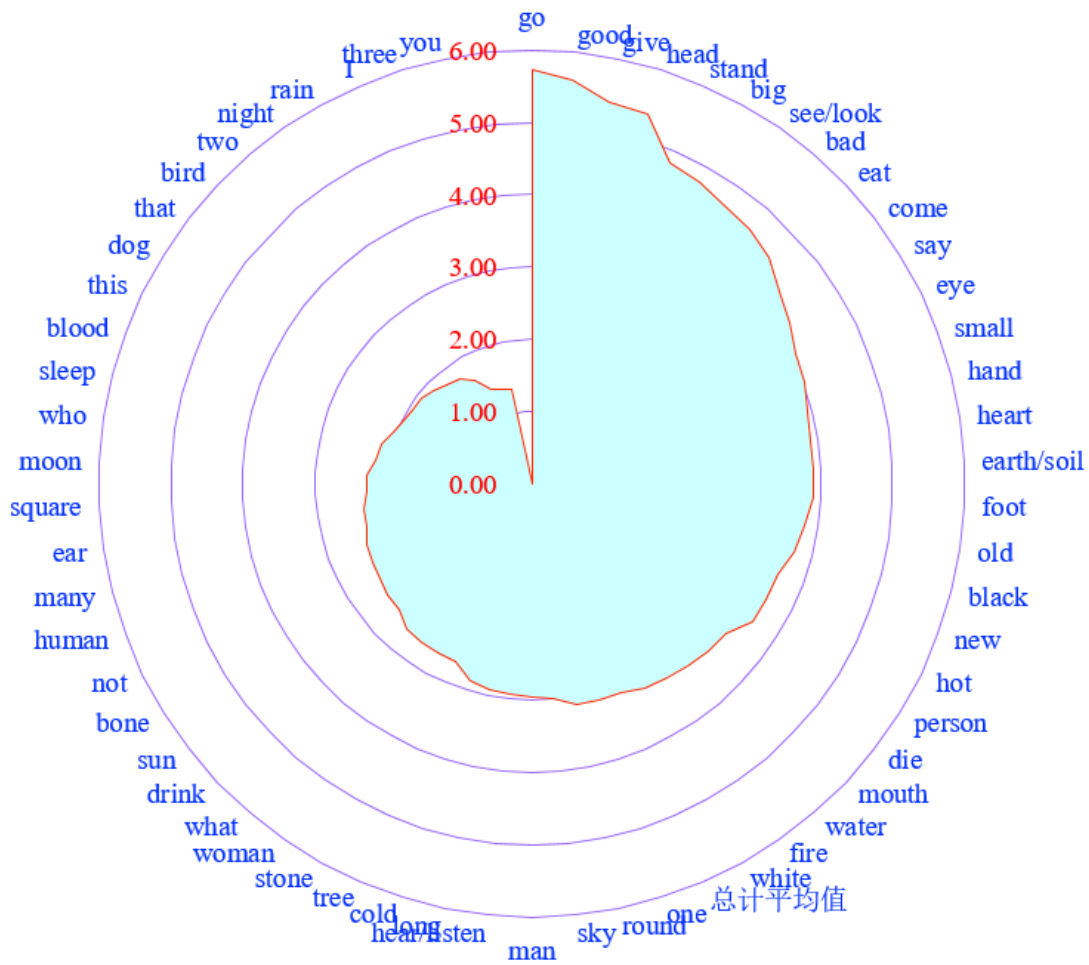


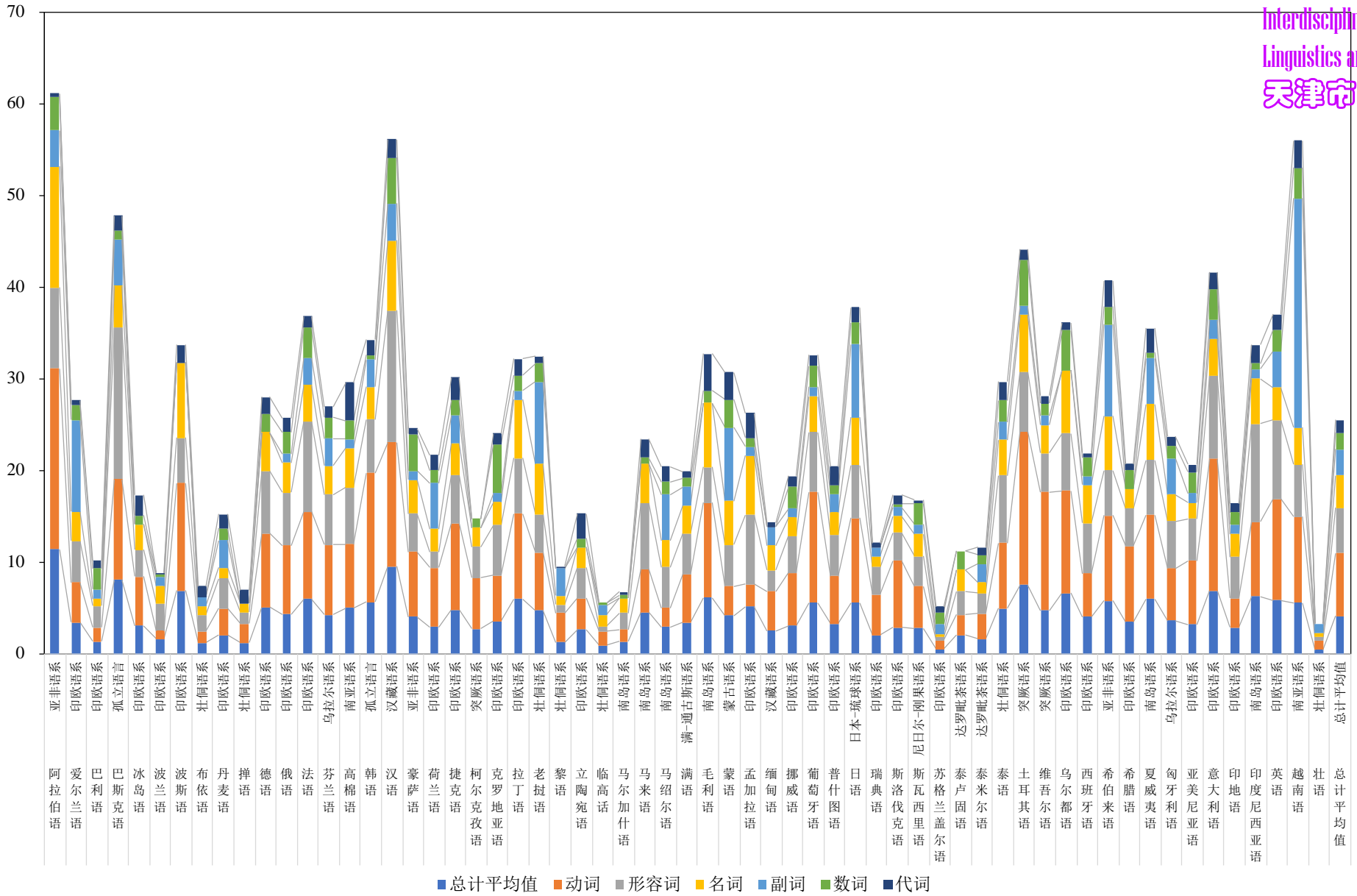
- The degree of dispersion in the distribution of the number of senses**

	Sum.	verb	adj.	noun	adv.	numeral	pronoun
Mean	3.17	4.14	3.56	2.99	2.58	2.12	1.86
S. D.	2.45	3.24	2.39	2.18	1.43	1.44	1.21
Maximum	21.00	20.00	15.00	21.00	9.00	8.00	9.00
Maximum estimated within the general range	10.51	13.87	10.73	9.53	6.87	6.45	5.49
Minimum	1.00	1.00	1.00	1.00	1.00	1.00	1.00



- Radar chart of sense mean values of the 60 words
- Cross-linguistic data
- 从基本概念来看
- 最容易产生出引申意义的基本概念
- 最容易发生认知延伸的基本概念
- 具体性? “tree”、“square”义项数量较少。“head”义项数量多
- 基本性? eat与drink
- 笼统性? cold等在不同语言中用多种不同级别的概念表示不





• Stacked histogram of average extension ability in 61 languages



- There is extremely rich information in the network atlas.
- The 10 most prominent basic concepts:
- **Go, good, give, head, stand, big, see/look, bad, eat, come**, are important basis for cognitive derivation in kernel words, it frequently undertakes the starting point and foundation of derived senses.
- The basic concepts that are most prone to cognitive derivation are not very “substantive” concepts, but rather those with “abstract” senses.
- “substantive” word: head (indicating that “head” is very important in linguistic cognition of human?)
- good: Pollyanna Hypothesis; give:??; stand.....



• 20 most frequent connections in deriving senses:

- not → no
- big → great, superior
- sky → paradise, heaven
- head → chief, leader
- foot → leg
- small → little
- head → top
- woman → female
- new → fresh
- earth/soil → ground
- come → reach, arrive
- say → talk, speak
- say → tell
- earth/soil → land
- sun → sunshine, sunlight
- black → dark
- bad → debased, evil, wicked
- hand → arm
- man → husband
- moon → month

• 20 most rare connections in deriving senses:

- person → capital towns
- person → chief, leader
- person → civil rights
- person → flesh
- see/look → have a look
- person → head
- person → life
- see/look → look into
- person → mask Latin
- person → principle
- person → principle money Latin
- person → source of rivers Latin
- see/look → spot
- person → substance, chapter
- person → the world Latin
- person → top Latin
- die → be exhausted
- die → breathe Icelandic
- die → crack
- die → explode Icelandic
- die → split



- The 60 kernel words can be clustered into 10 groups.
- (different clustering methods were used)

Module	Basic kernel words
1	stand、 give、 go、 hear/listen、 see/look、 come、 say、 sleep、 die
2	rain、 sky、 eat、 drink、 water、 moon
3	hand、 round、 square、 ear、 earth/soil、 white
4	head、 mouth、 eye
5	stone、 tree、 heart、 foot、 cold、 bone
6	blood、 bird、 man、 person、 you、 I、 human、 woman
7	big、 many、 long、 good、 old、 not、 three
8	dog、 bad、 black、 small、 night
9	one、 two、 that、 what、 this、 who
10	sun、 hot、 fire、 new



• **III- Acoustic distance analysis on linguistic elements/language/dialect**

• **Jun Ding^{1†}, Tianheng Wang^{1†}, Ke Xu¹, Aleksandr Mitkov¹, Lining Wang^{2*}, Gang Peng^{3*}, Qibin Ran^{1,4*}, Acoustic distances of 300 core words imply Indo-European phylogeny and chronology.**

Submitted Manuscript: Confidential
Template revised November 2023

Title: Acoustic distances of 300 core words imply Indo-European phylogeny and chronology

Authors: Jun Ding^{1†}, Tianheng Wang^{1†}, Ke Xu¹, Aleksandr Mitkov¹, Lining Wang^{2*}, Gang Peng^{3*}, Qibin Ran^{1,4*}

5 **Affiliations:**

- ¹School of Liberal Arts, Nankai University; Tianjin, 300071, China.
 - ²Center for the Protection and Research of Language Resources of China, Beijing Language and Culture University; Beijing, 100083, China.
 - ³Research Centre for Language, Cognition, and Neuroscience, Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University; Kowloon, Hong Kong SAR.
 - ⁴Laboratory of Social Science of Tianjin, Nankai University; Tianjin, 300071, China.
- *Corresponding author. Email: wangln@bclu.edu.cn (L.W.); gpeng@polyu.edu.hk (G.P.); ranqibin@126.com (Q.R.)

†These authors contributed equally to this work.

15 **Abstract:** Language offers crucial clues to human history; yet can these clues be directly gleaned from speech sounds? Unlike previous investigations into language chronology and phylogeny that primarily rely on historical comparisons, this study explores an acoustic approach to examine the history of Indo-European languages by leveraging synthesized sound samples of 300 core concepts in 42 languages. The acoustic distance-based clustering of these languages aligns well with established classifications, and the acoustic distances effectively reflect language divergence ages. Acoustic results imply that the origin of Indo-European languages (except Anatolian and Tocharian) dates back to 5700–5000 BP. This groundbreaking acoustic approach not only shows potential to unravel language chronology and phylogeny, but also provides fresh methods and perspectives on human prehistory.

25 **One-Sentence Summary:** Synthesized sounds from 42 Indo-European languages disclose the age, classification, and origin date of these languages.



THE INTERNATIONAL PHONETIC ALPHABET (revised to 2018)

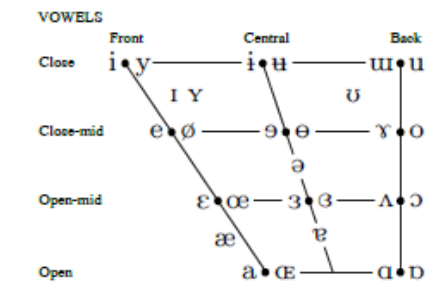
CONSONANTS (PULMONIC)

	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retroflex	Palatal	Velar	Lingual	Pharyngeal	Glottal
Plosive	p b			t d							
Nasal	m	ɱ		n		ɳ	ɲ	ŋ			
Trill				r						ʀ	
Tap or Flap				ɾ		ɽ					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	h ɦ
Lateral fricative				ɬ ɮ							
Approximant		ʋ		ɹ		ɻ	j	ɰ			
Lateral approximant				l		ɭ	ʎ	ʟ			

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
◌ Bilabial	ɓ Bilabial	ʼ Examples:
Dental	ɗ Dental/alveolar	pʼ Bilabial
! (Post)alveolar	ɟ Palatal	tʼ Dental/alveolar
‡ Palatoalveolar	ɡʷ Velar	kʼ Velar
Alveolar lateral	ɠ Uvular	sʼ Alveolar fricative



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

OTHER SYMBOLS

ʌ Voiceless labial-velar fricative ʎ ʐ Alveolo-palatal fricatives
 W Voiced labial-velar approximant ɺ Voiced alveolar lateral flap
 ɥ Voiced labial-palatal approximant ɧ Simultaneous ʃ and x
 H Voiceless epiglottal fricative
 ʕ Voiced epiglottal fricative Affricates and double articulations can be represented by two symbols joined by a tie bar if necessary. ts̺ k̟p̟
 ʡ Epiglottal plosive

DIACRITICS Some diacritics may be placed above a symbol with a descender, e.g. ɹ̥̄

◌ Voiceless	◌̥ ◌̜	◌ Breathily voiced	◌̚ ◌̝	◌ Dental	◌̪ ◌̫
◌ Voiced	◌̩ ◌̬	◌ Creaky voiced	◌̰ ◌̱	◌ Apical	◌̽ ◌̿
◌ Aspirated	◌ʰ ◌ʰ̚	◌ Linguolabial	◌̍ ◌̎	◌ Laminal	◌̥ ◌̦
◌ More rounded	◌̙	◌ Labialized	◌̙ ◌̙̟	◌ Nasalized	◌̃
◌ Less rounded	◌̚	◌ Palatalized	◌̟ ◌̟̟	◌ Nasal release	◌̚̚
◌ Advanced	◌̟	◌ Velarized	◌̙ ◌̙̟	◌ Lateral release	◌̚̚̚
◌ Retracted	◌̚	◌ Pharyngealized	◌̙̟ ◌̙̟̟	◌ No audible release	◌̚̚̚̚
◌ Centralized	◌̚̚	◌ Velarized or pharyngealized	◌̙̟̟		
◌ Mid-centralized	◌̚̚̚	◌ Raised	◌̚̚̚̚ (◌̚̚̚̚ = voiced alveolar fricative)		
◌ Syllabic	◌̚̚̚̚	◌ Lowered	◌̚̚̚̚̚ (◌̚̚̚̚̚ = voiced bilabial approximant)		
◌ Non-syllabic	◌̚̚̚̚̚	◌ Advanced Tongue Root	◌̚̚̚̚̚̚		
◌ Rhoticity	◌̚̚̚̚̚̚ ◌̚̚̚̚̚̚̚	◌ Retracted Tongue Root	◌̚̚̚̚̚̚̚		

SUPRASEMENTALS

◌ Primary stress ˈ founəˈtʃən
 ◌ Secondary stress ˌ
 ◌ Long ː
 ◌ Half-long eˑ
 ◌ Extra-short ɛ̥
 ◌ Minor (foot) group ◌̥
 ◌ Major (intonation) group ◌̥
 ◌ Syllable break .i.sɛkt
 ◌ Linking (absence of a break) ◌̥

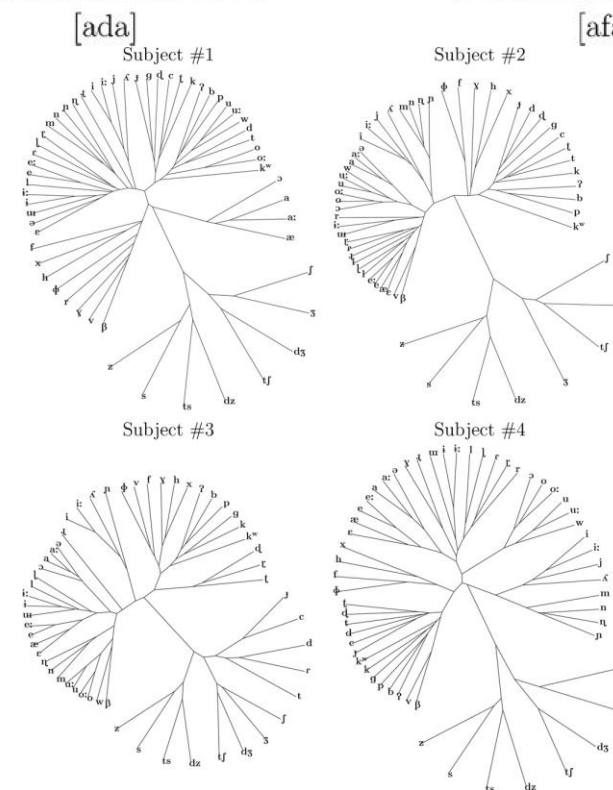
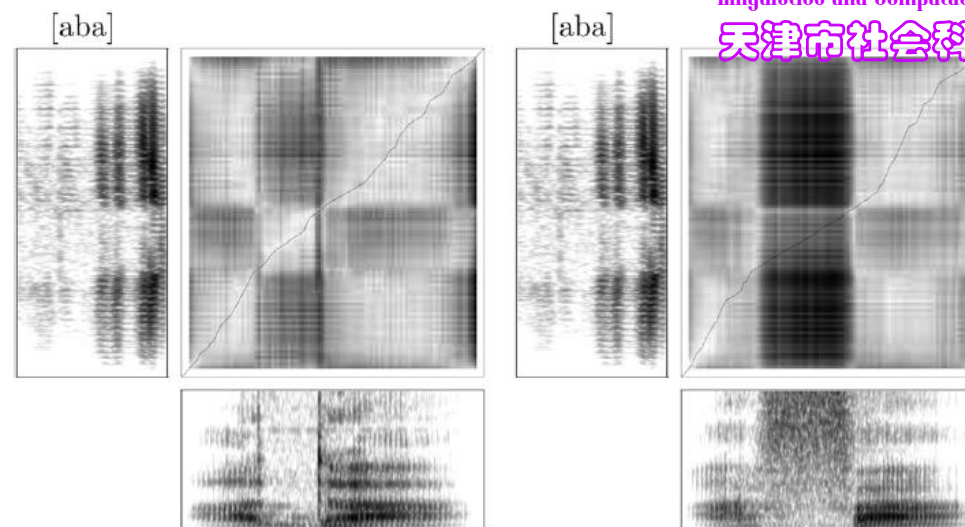
TONES AND WORD ACCENTS

LEVEL CONTOUR

◌̥ or ˥ Extra high ◌̥ or ˨ Rising
 ◌̥ High ◌̥ Falling
 ◌̥ Mid ◌̥ High rising
 ◌̥ Low ◌̥ Low rising
 ◌̥ Extra low ◌̥ Rising-falling
 ↓ Downstep ↗ Global rise
 ↑ Upstep ↘ Global fall

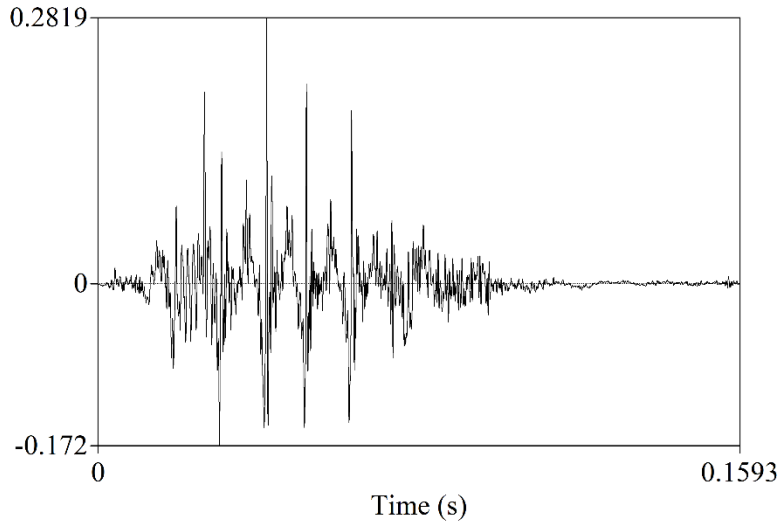
- Difference between sounds
- Distances between sounds are equal or not?
- D (s,m) > D (ɛ,æ)
- D (p,s) > D (n,l)
-

- Itakura(板仓1965): DTW
- Holmes and Holmes (2001) : (DTW, dynamic time warping algorithm)
- Mielke(2012): VXV
- J. Mielke(2012)A phonetically based metric of sound similarity, *Lingua*122(2).
- Mielke(2012): 58 phonemes (vowel: 17, consonant: 41)
- Phylip

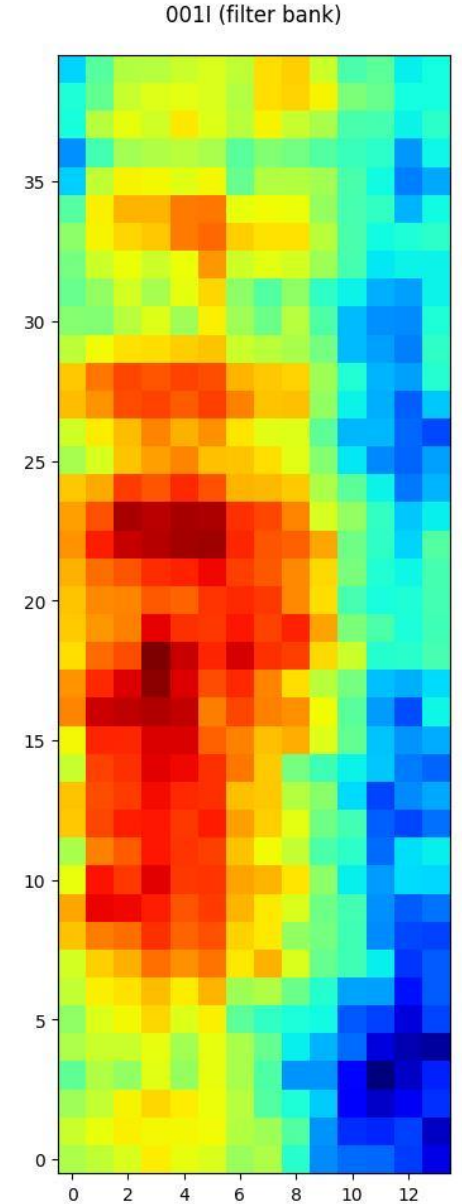
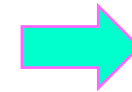




- Changing a sound file into logbank spectrogram



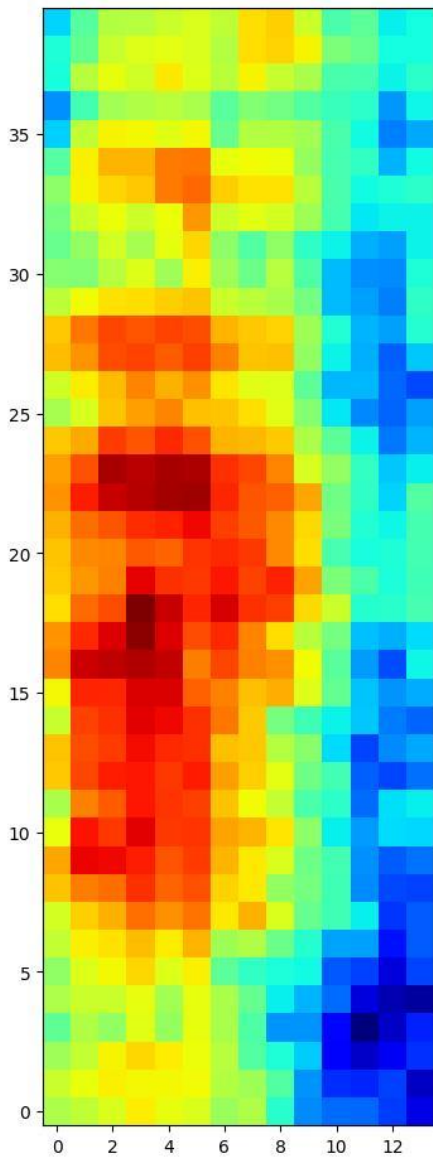
```
AfrikaansMale.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
AfrikaansMale
001I
[[16.09318872 16.38537773 16.83003579 17.47683358 17.08911863 16.85441997
15.73627862 16.11543865 13.77538309 11.80183848 11.36335178 11.2798829
10.67291908 9.34844731]
[16.58517382 17.10813229 17.43071866 17.28786053 17.31027178 17.24261797
16.17248206 15.83117861 14.50531471 11.92550699 10.45546413 10.33278323
10.72898563 8.90095918]
[15.82529966 16.4834277 17.36693708 17.78335786 17.4806111 17.19161539
16.26481013 14.52458704 13.62485972 12.69860732 9.75052798 8.96553931
9.48600653 10.49473751]
[14.72298115 16.11393345 15.61291548 16.98935248 15.69078029 17.08828438
16.03738268 14.45094048 11.92037299 11.91672489 9.88481705 8.02566548
8.96278404 10.29343605]
[16.05661021 16.55791366 16.60123298 17.07844943 15.91260572 17.0938056
16.01016829 14.8321726 13.2252715 12.39695091 11.37621658 9.24268275
8.65956295 8.43519 ]
[15.54487616 16.89257079 17.32946719 17.8346006 16.9248057 17.44017453
14.64323941 13.95905894 13.53141577 13.37355062 11.09456808 10.76666433
9.2163525 10.79778705]
[16.4820098 17.419895 17.61326176 18.17816624 17.46027706 18.33582932
15.82245147 16.11213465 14.95319686 13.72892263 12.13627781 12.08509168
10.0697676 11.15493427]
[16.7365217 17.91804913 18.37319209 19.36032961 18.84584334 19.3042299
17.52697959 18.37171566 16.86118357 14.83554119 14.15208485 13.23441042
10.56745502 11.15495653]
[18.11195695 19.15372323 19.40038847 20.30760187 19.57353913 19.80908273
17.86239659 17.50256681 15.62080003 15.09018312 14.1734914 11.8049879
10.84228499 10.76922011]
[18.54152722 21.01034634 20.95344386 20.6278906 19.74407471 20.17624332
18.30591591 17.51937492 16.83569349 15.01115479 14.24508105 11.83882843
11.12643407 11.47943924]
[17.15057314 20.77350211 20.22063685 21.09120356 20.2463769 20.27976977
18.58319707 18.30533252 17.60033072 15.50000647 13.19898895 12.06067805
12.95104901 12.84717946]
[16.00552 19.12182162 19.67924697 20.71441759 20.29035037 20.11900052
18.07272306 17.23956444 16.46996405 14.38013725 13.84908192 11.37722387
12.96338229 13.21392332]
[18.03535828 19.96956623 20.59378352 20.70333884 20.23486475 20.64173168
18.66050151 17.94574558 17.02550565 14.98169702 14.20834192 11.21812716
10.8239827 11.39608847]
[18.08515385 19.87385218 20.16557974 20.86785916 20.44791867 20.30797087
18.16905411 18.03310602 16.16399936 15.43304232 12.92637755 10.78404985
11.78492508 12.20868754]
[16.53737842 20.07301907 20.31719053 21.12820701 20.94713208 20.32528508
```



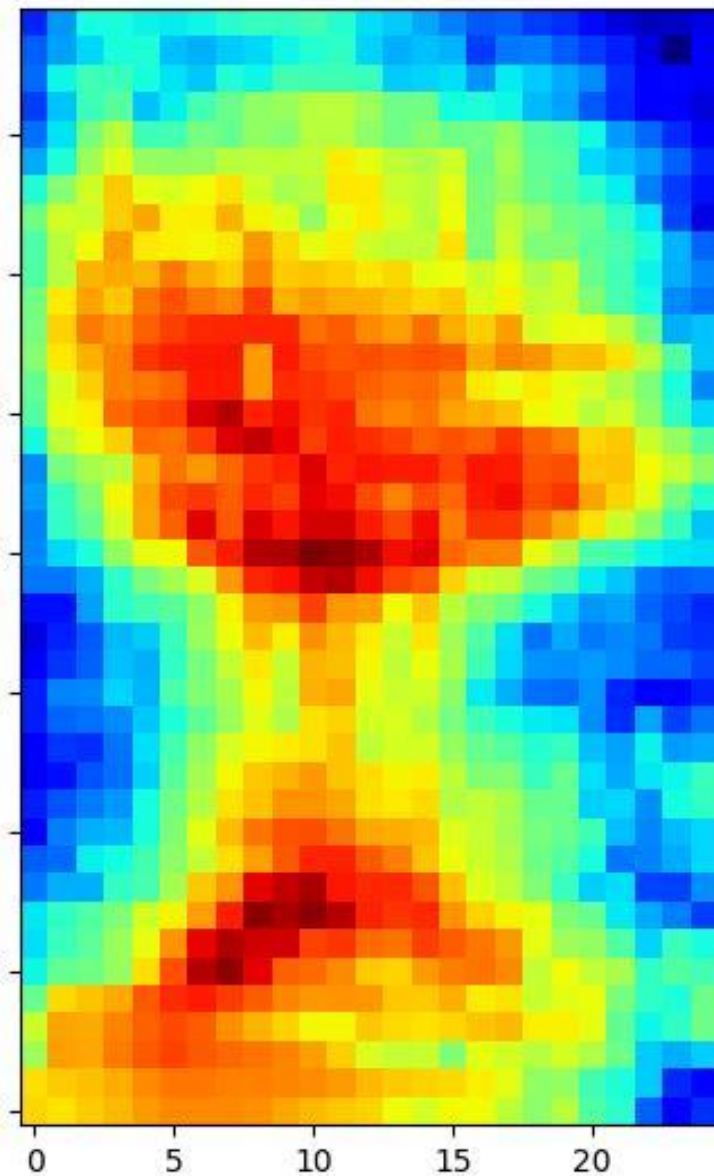


• Afrikaans

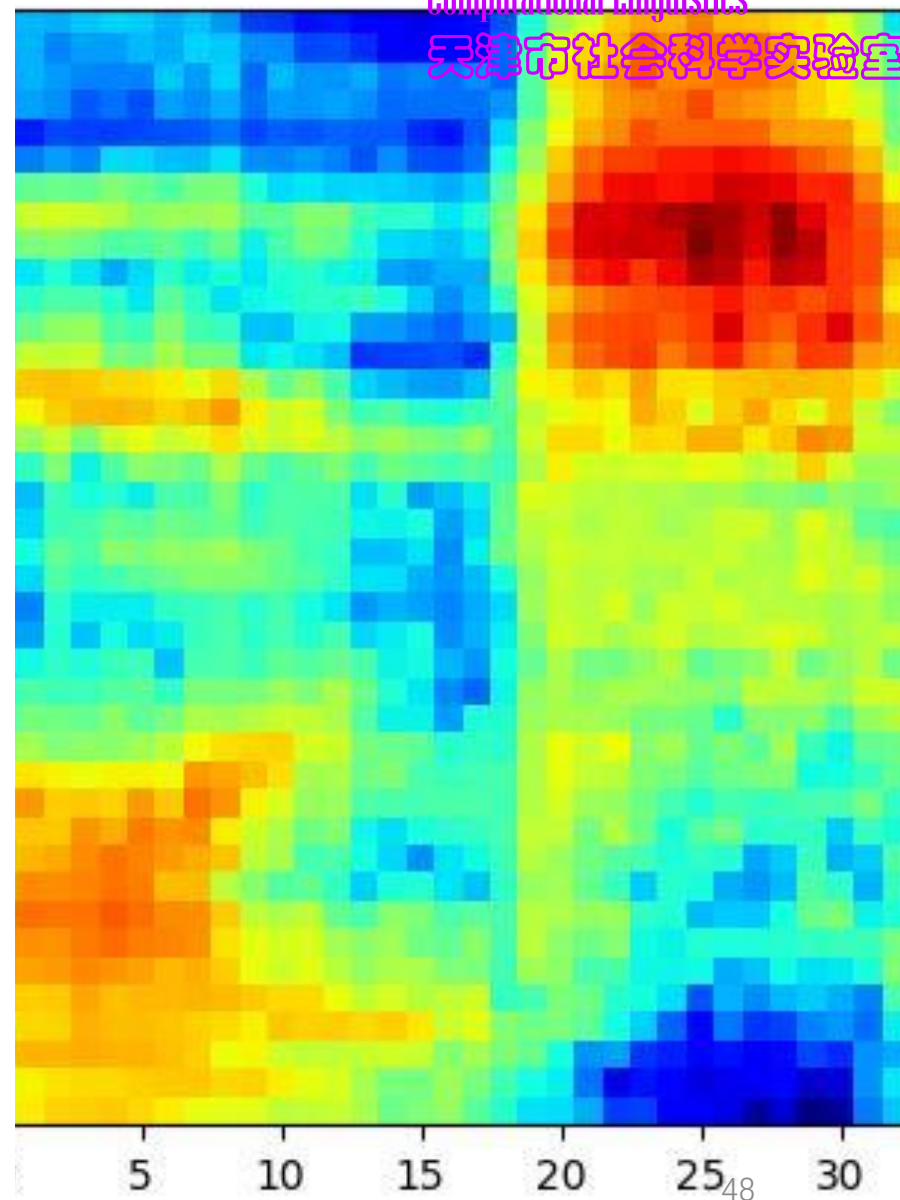
001i (filter bank)



002you-a (filter bank)



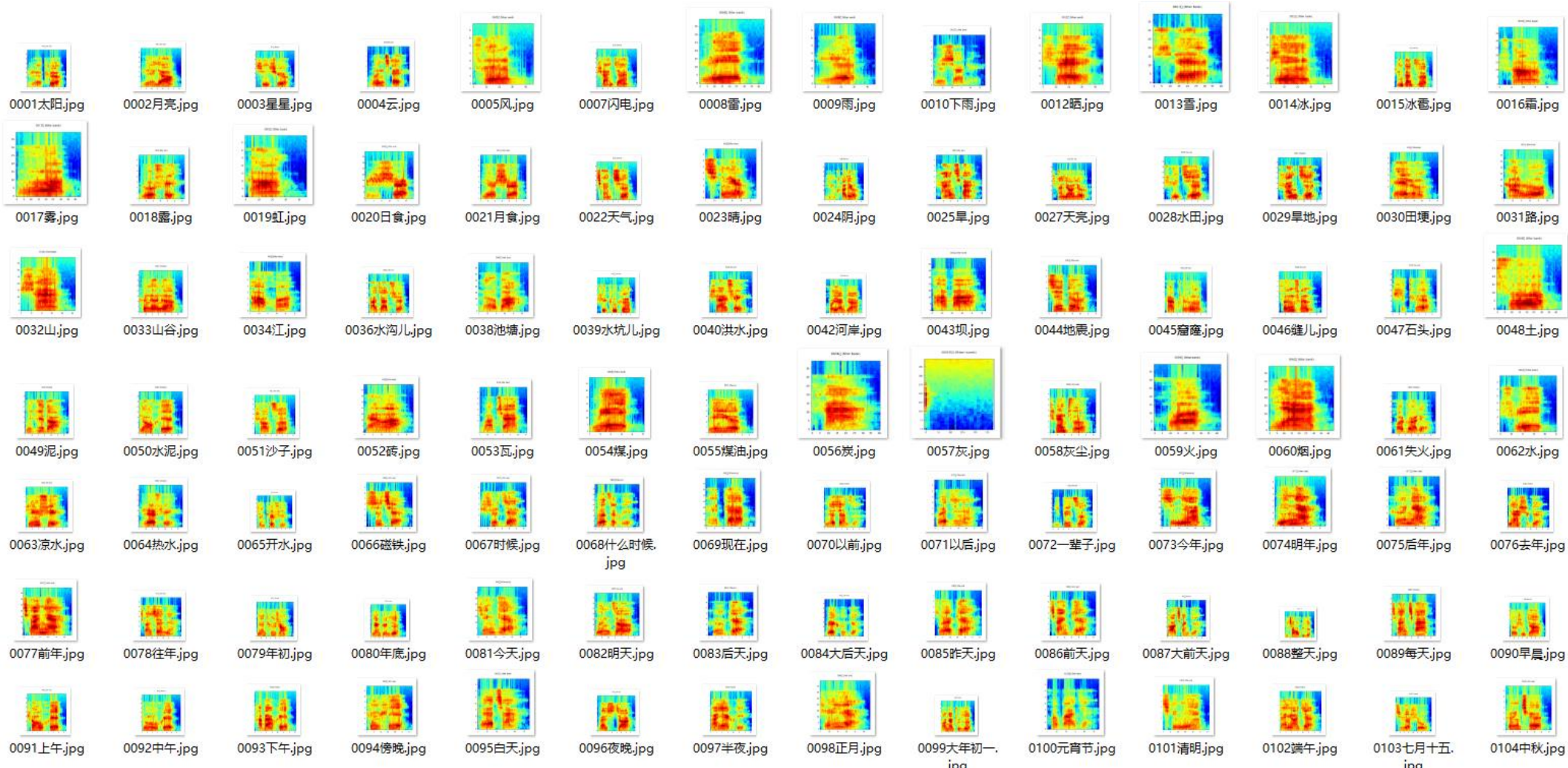
003we (filter bank)



Interdisciplinary Laboratory for
Experimental Linguistics and
Computational Linguistics

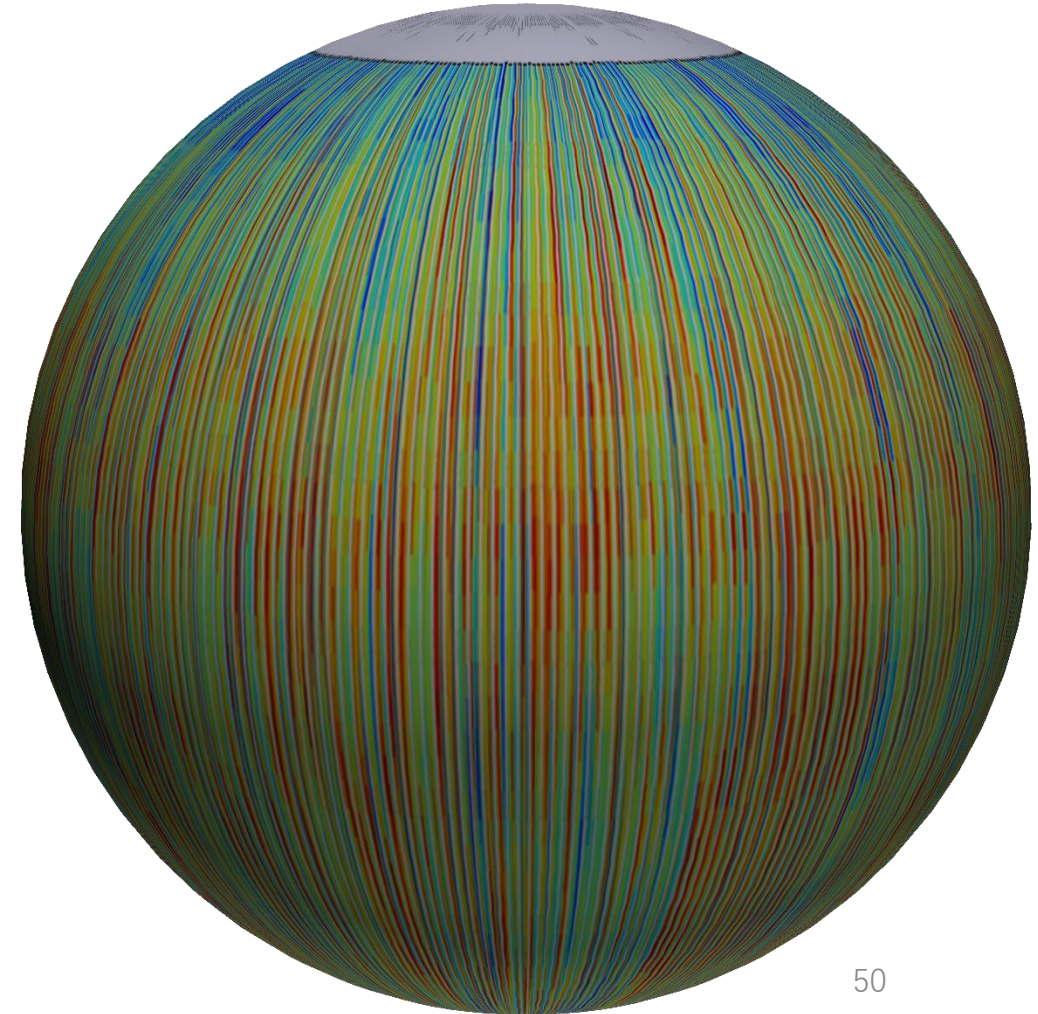
天津市社会科学实验室

- Microsoft TTS engine: Azure
- Afrikaans.txt (98/300)



- **Acoustic graph of Affrikaans words(300 kernel words)**

- [Afrikaans.tif](#)
- **Acoustic sphere of Affrikaans(planet/star)**
- **Digital acoustic map/atlas of kernel words?**
- **Sound-Meaning spheres of a language**

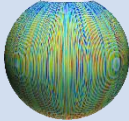
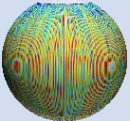
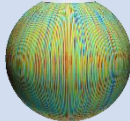
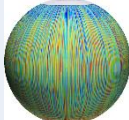
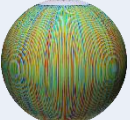
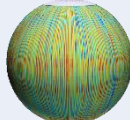
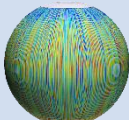
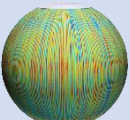
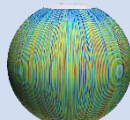
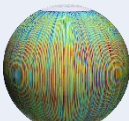
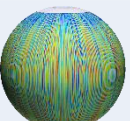
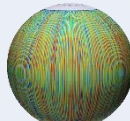
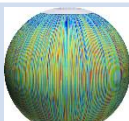
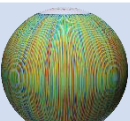
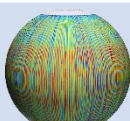
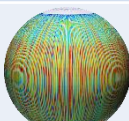
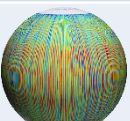
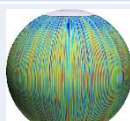
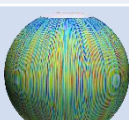
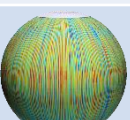
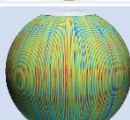




- Sound-Meaning spheres of 300 kernel concepts of IE languages (1)
- 21/42 IE languages

language	S-M sphere	language	S-M sphere	language	S-M sphere
Afrikaans		Croatian		German	
Albanian		Czech		Greek	
Armenian		Danish		Gujarati	
Bengali		Dutch		Hindi	
Bosnian		English		Icelandic	
Bulgarian		French		Irish	
Catalan		Galician		Italian	

• Sound-Meaning spheres of 300 kernel concepts of IE languages (2)

language	S-M sphere	language	S-M sphere	language	S-M sphere
Latvian		Persian		Slovak	
Lithuanian		Polish		Slovenian	
Macedonian		Portuguese		Spanish	
Marathi		Romanian		Swedish	
Nepali		Russian		Ukrainian	
Norwegian		Serbian		Urdu	
Pashto		Sinhala		Welsh	

• Using DTW algorithm to calculate distance between any 2 languages

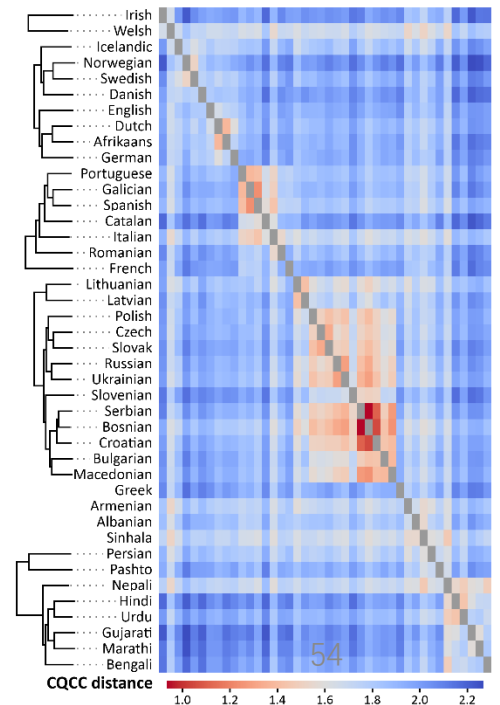
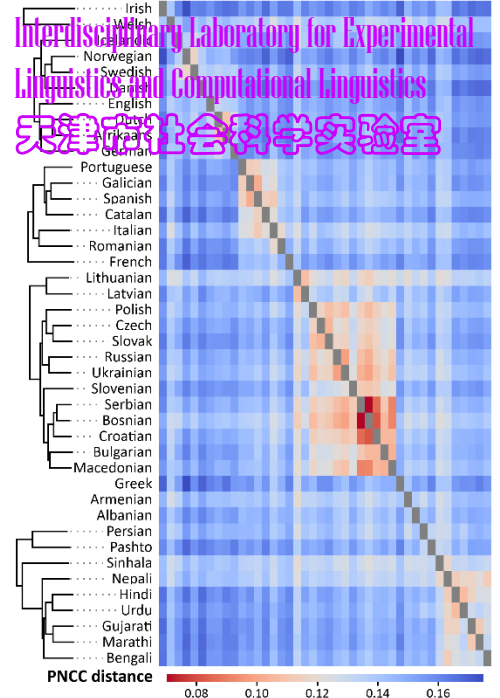
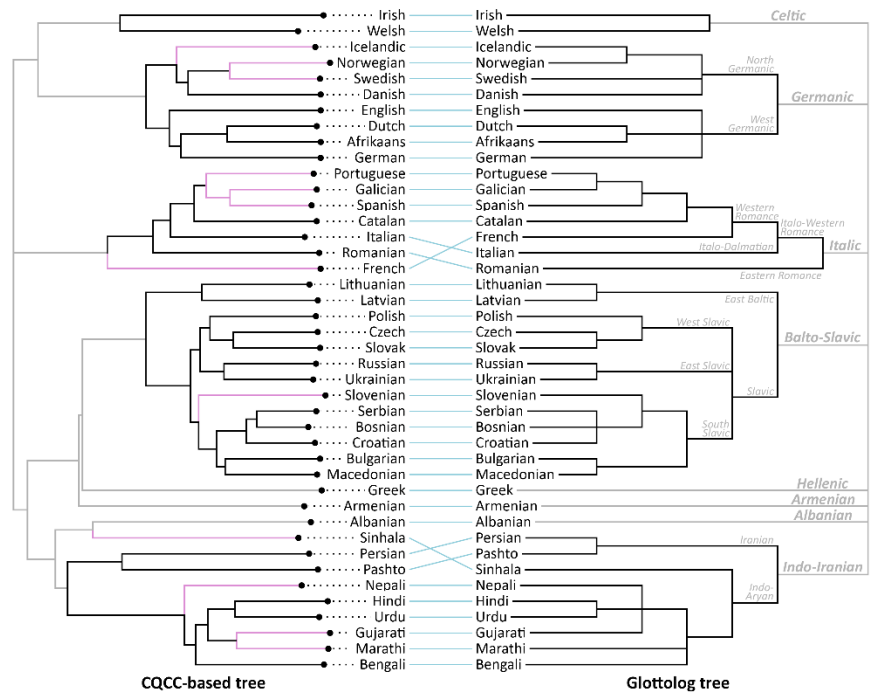
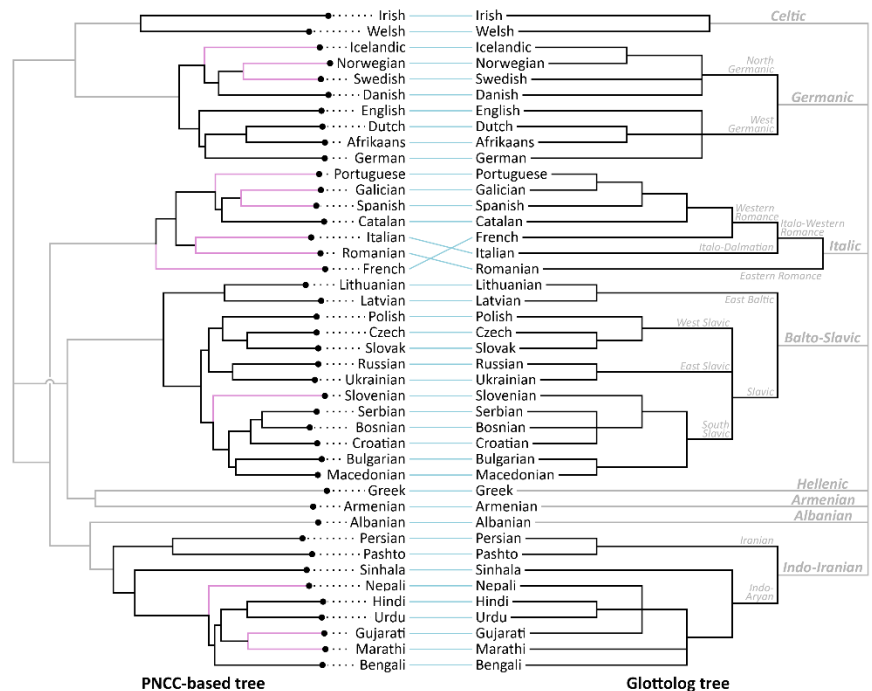


• Matrix:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ											
1	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ											
2	Afrikaans	Albanian	Armenian	Bengali	Bosnian	Bulgarian	Catalan	Croatian	Czech	Danish	Dutch	English	French	Galician	German	Greek	Gujarati	Hindi	Icelandic	Irish	Italian	Latvian	Lithuanian	Macedonian	Marathi	Nepali	Norwegian	Pashto	Persian	Polish	Portuguese	Romanian	Russian	Serbian	Sinhala	Slovak	Slovenian	Spanish	Swedish	Ukrainian	Urdu	Welsh												
3	0	0.89998	0.28939	0.31269	0.27288	0.21656	0.42708	0.311769	0.12706	0.12448	0.29954	0.29295	0.20711	0.26355	0.26842	0.2929	0.24463	0.20028	0.29874	0.24221	0.20624	0.20781	0.25989	0.29199	0.26182	0.11684	0.29482	0.200787	0.2711	0.22228	0.22487	0.23396	0.23655	0.22228	0.22487	0.23396	0.23655	0.22228	0.22487	0.23396	0.23655	0.22228	0.22487	0.23396	0.23655									
4	0.28998	0	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626	0.274626									
5	0.28939	0.274626	0	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023							
6	0.31269	0.274626	0.274626	0	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023	0.21023						
7	0.27288	0.21656	0.27288	0.21656	0	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399	0.246399						
8	0.21656	0.274626	0.274626	0.21656	0.246399	0	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225	0.236225					
9	0.27288	0.21656	0.27288	0.21656	0.246399	0.236225	0	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526	0.217526					
10	0.21656	0.274626	0.274626	0.21656	0.246399	0.236225	0.217526	0	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715	0.217715				
11	0.27288	0.21656	0.27288	0.21656	0.246399	0.236225	0.217526	0.217715	0	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235	0.203235			
12	0.21656	0.274626	0.274626	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171	0.293171				
13	0.27288	0.21656	0.27288	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0.293171	0	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189	0.218189		
14	0.21656	0.274626	0.274626	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0.293171	0.218189	0	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099	0.20099		
15	0.27288	0.21656	0.27288	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0.293171	0.218189	0.20099	0	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	0.299777	
16	0.21656	0.274626	0.274626	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0.293171	0.218189	0.20099	0.299777	0	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	0.28296	
17	0.27288	0.21656	0.27288	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0.293171	0.218189	0.20099	0.299777	0.28296	0	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296	0.292296
18	0.21656	0.274626	0.274626	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0.293171	0.218189	0.20099	0.299777	0.28296	0.292296	0	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366	0.27366		
19	0.27288	0.21656	0.27288	0.21656	0.246399	0.236225	0.217526	0.217715	0.203235	0.293171	0.218189	0.20099	0.299777	0.28296	0.292296	0.27366	0	0.273869	0.2																																			



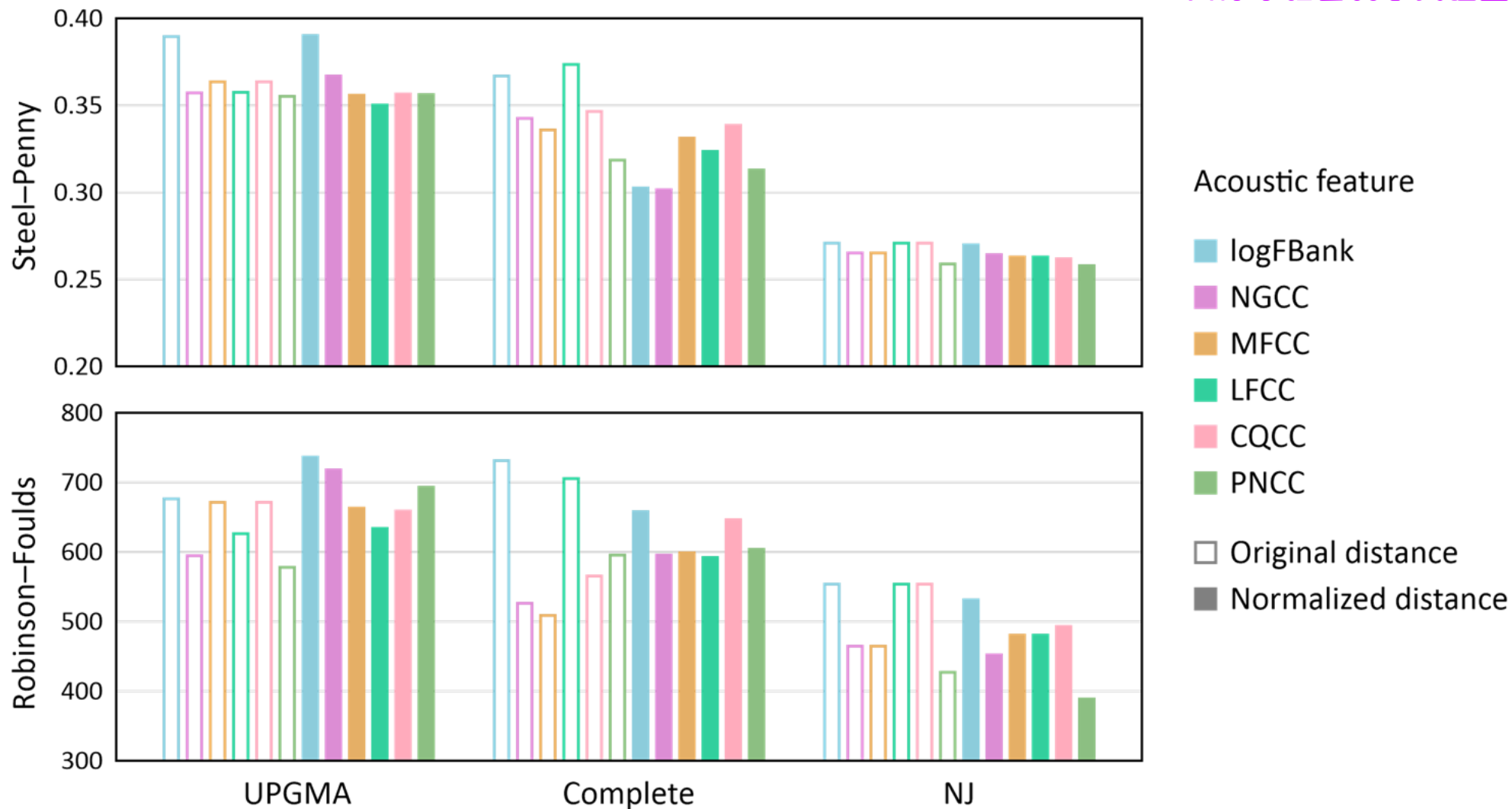
- New phylogenetic tree is very similar to the traditional classification.
- This is really interesting.
- We tested more than 10 acoustic features.
- ✕ writing system
- ✕ transcription
- Language history is hided/embedded in speech sounds, or in everyone's mouth!



Interdisciplinary Laboratory for Experimental Linguistics and Computational Linguistics
天津社会科学实验室



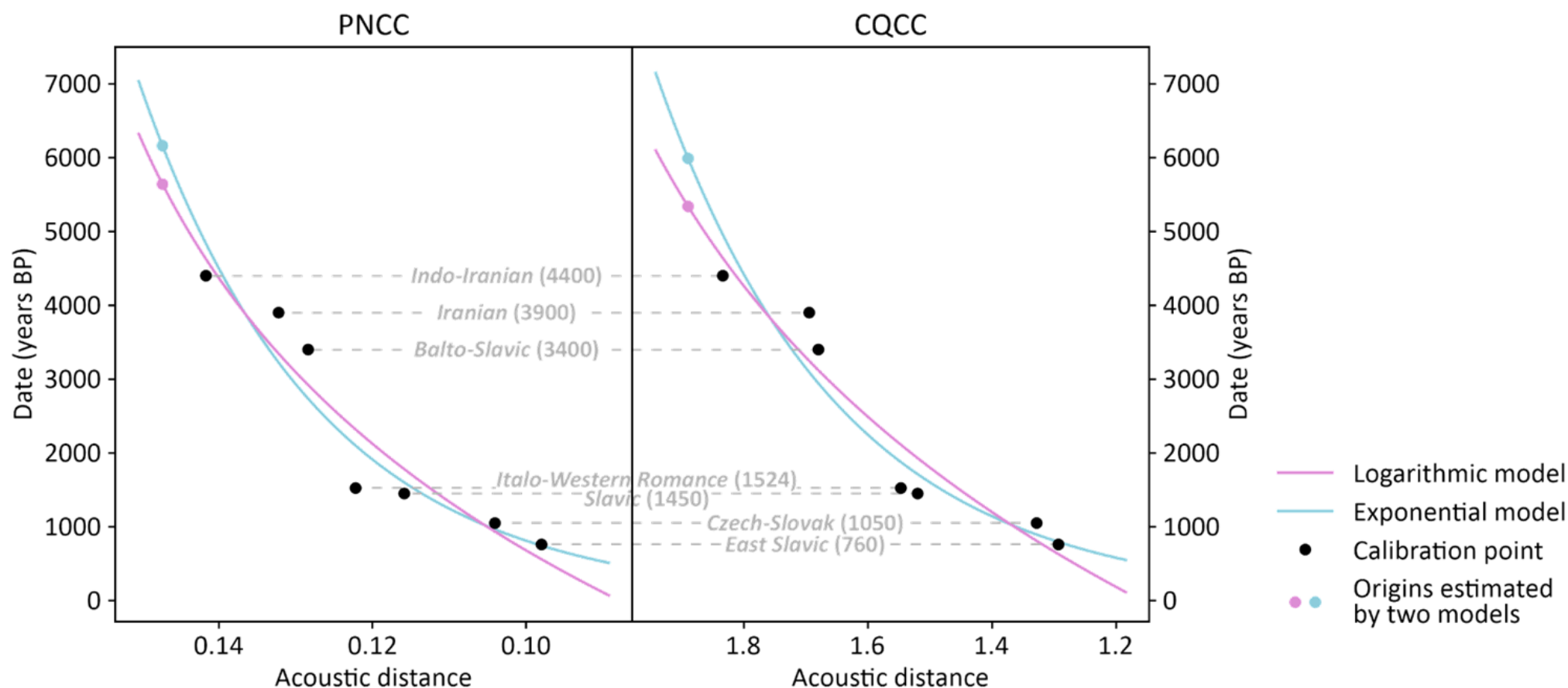
- Different acoustic features



- Shorter the column is, better the acoustic feature/method is.



- **Fitting chronology!**
- **The acoustic distance matrix, can not only be used to construct the phylogenetic tree, but also can be used to fit I-E language chronology.**





• The acoustic distance-based clustering of these languages aligns well with established classifications, and the acoustic distances effectively reflect language divergence ages.

• Acoustic results imply that the origin of Indo-European languages (except Anatolian and Tocharian) dates back to **5,700–5,000 BP**.

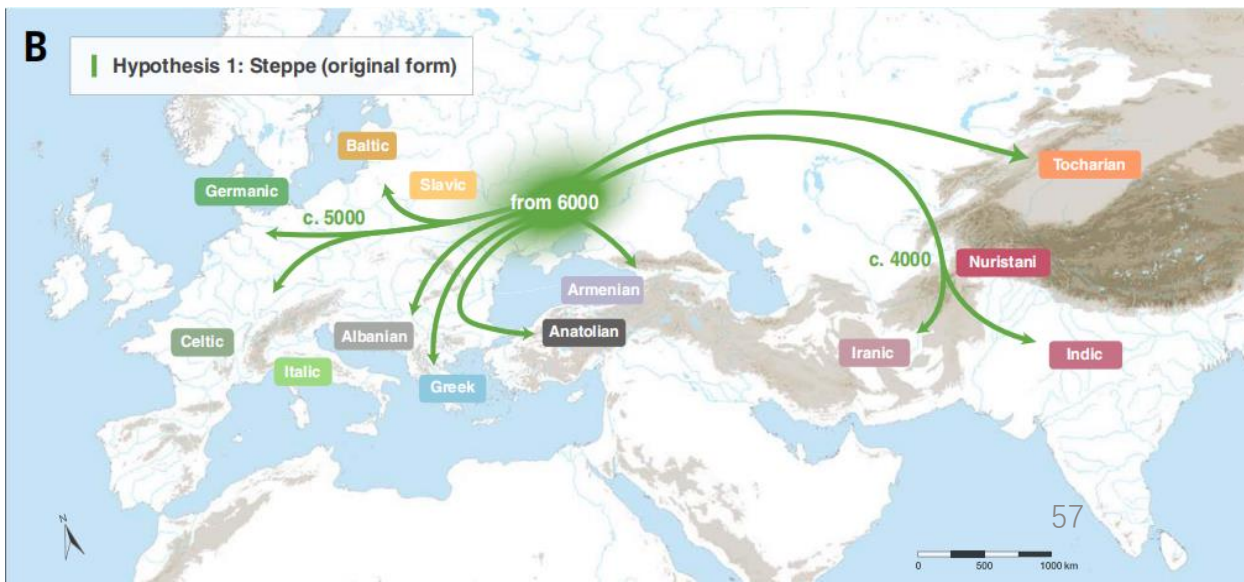
• 2 hypothesis:

• Anatolian hypothesis: 8,500BP, agricultural diffusion

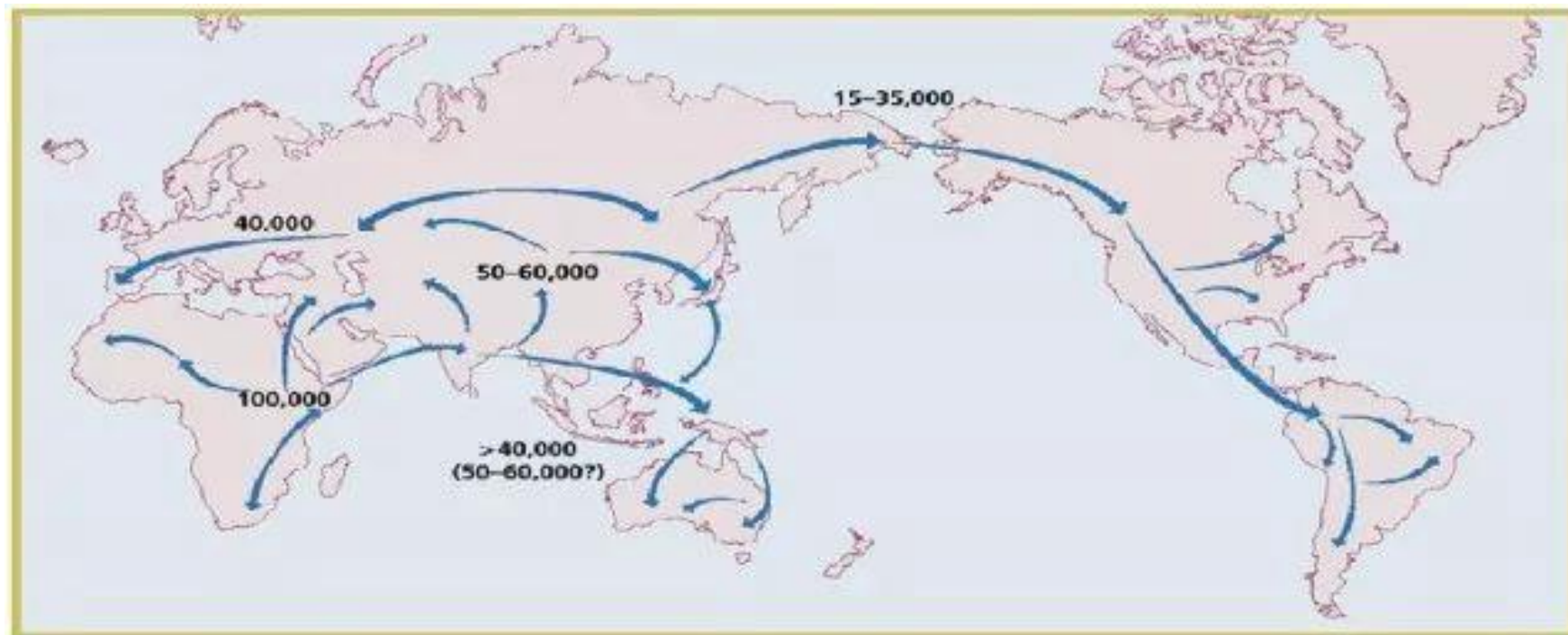
• ✓ Kurgan hypothesis (grassland hypothesis): 6,500BP-4,500BP

• This is really interesting and meaningful.

• It is beyond expectations.

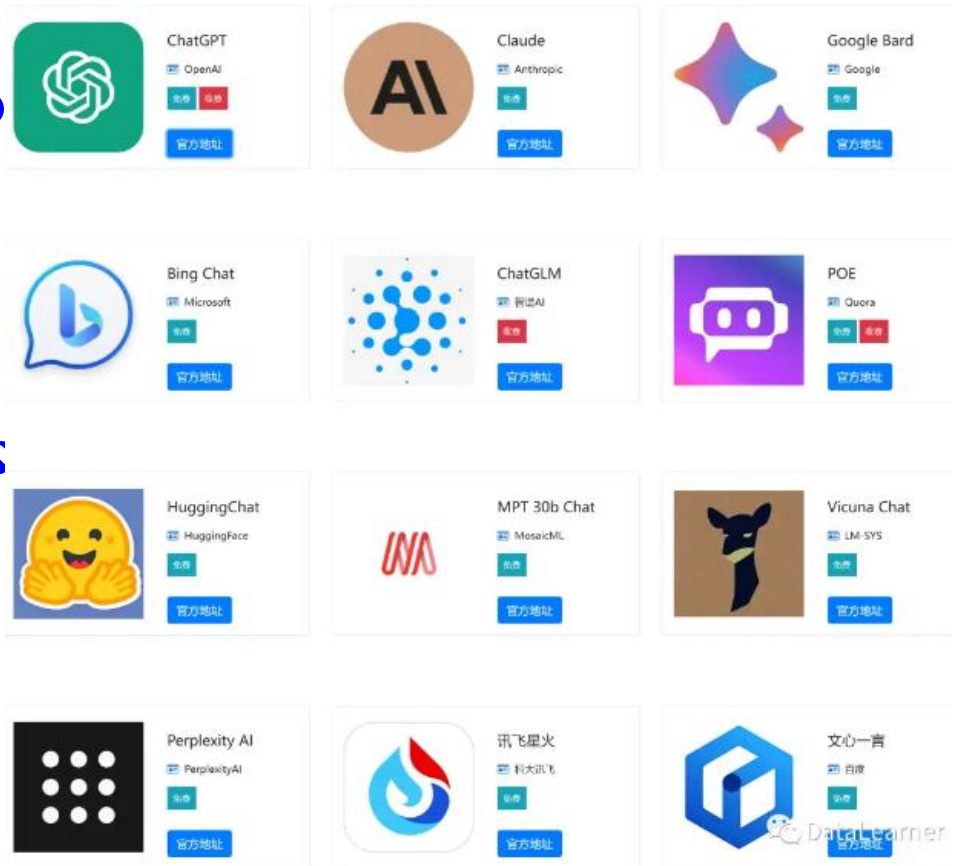


- **Future work:**
- **Simulating and reconstructing the migration route of ancestors of modern human from the ‘homeland’ in Africa**
- **Global Recordings Network(<https://globalrecordings.net/en/>)**
- **6,527 language varieties**



The age of Artificial Intelligence

- We are in the age of AI developing rapidly.
- Accurate automatic speech recognition, Text-to-Speech(speech synthesis).....
- Image recognition, Face recognition.....
- Prompt Engineering:
- Large Language Model(LLM, ChatGPT, 文心一言等).....
- Generate images, Generate videos.....
- transportation、shopping、financial service、education, health care.....





IV-Engineering attempts using deep learning and AI technology in our lab

- 1-Speech synthesis(Text-To-Speech)
- LingSound speakers: Lingsheng, Lingzhi, Junsheng, Junzhi
- specific speakers, comic figures, Chinese dialects, chanting for ancient poem...
- Pingwu(Sichuan) dialect; Tianjin dialect;

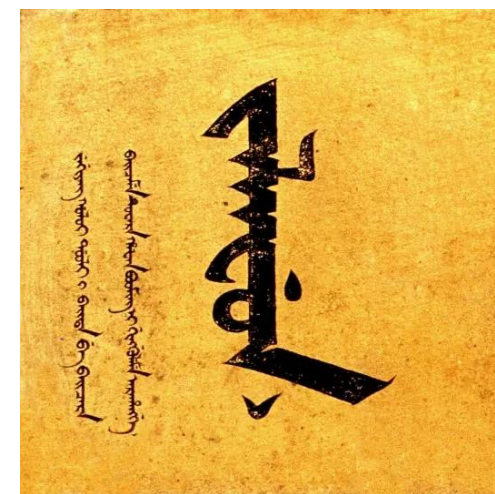


项目	链接	简介
语音合成	语音合成	支持包括实验室自研男声、女声；名人声音；天津话、四川话等方言的合成
国际音标识别	国际音标识别	支持国际音标的流式识别
语种识别	语种识别	支持汉语普通话、英语、法语、德语、西班牙语、日语、加泰罗尼亚语、卢旺达语等8种语言的识别。
汉语学习者发音评分	汉语学习者发音评分	对汉语学习者的口语进行评估
声纹识别	声纹识别	识别已注册声音
满汉互译系统	满汉互译系统	支持满语、汉语双向翻译（初步测试）



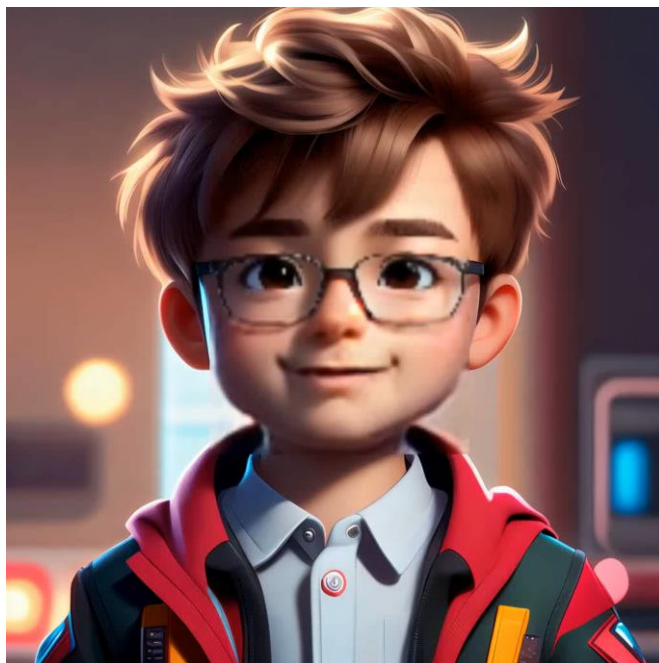
- **2-Automatic Speech Recognition(IPA ASR, language recognition, foreign accent evaluation, speaker identification...)**
- **demo**
- **<https://www.ranqibin.com/language-identification>**

- **3-Machine Translation(Manchu-Mandarin bidirectionally automatic translation)**
- **Endangered language, Paul Georg von Möllendorff's transcription**
- **<https://manchumt.app.ranqibin.com/>**



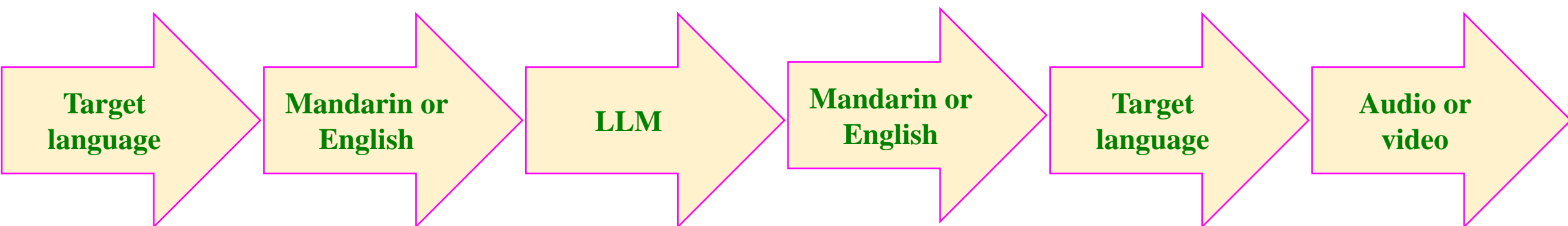


- **Engineering attempts using deep learning and AI in our lab**
- **4-digital human (Yuen-Ren Chao, Chia-ying Yeh, the first president of Nankai Univ.; Tangwangnese, Tianjin dialect)**



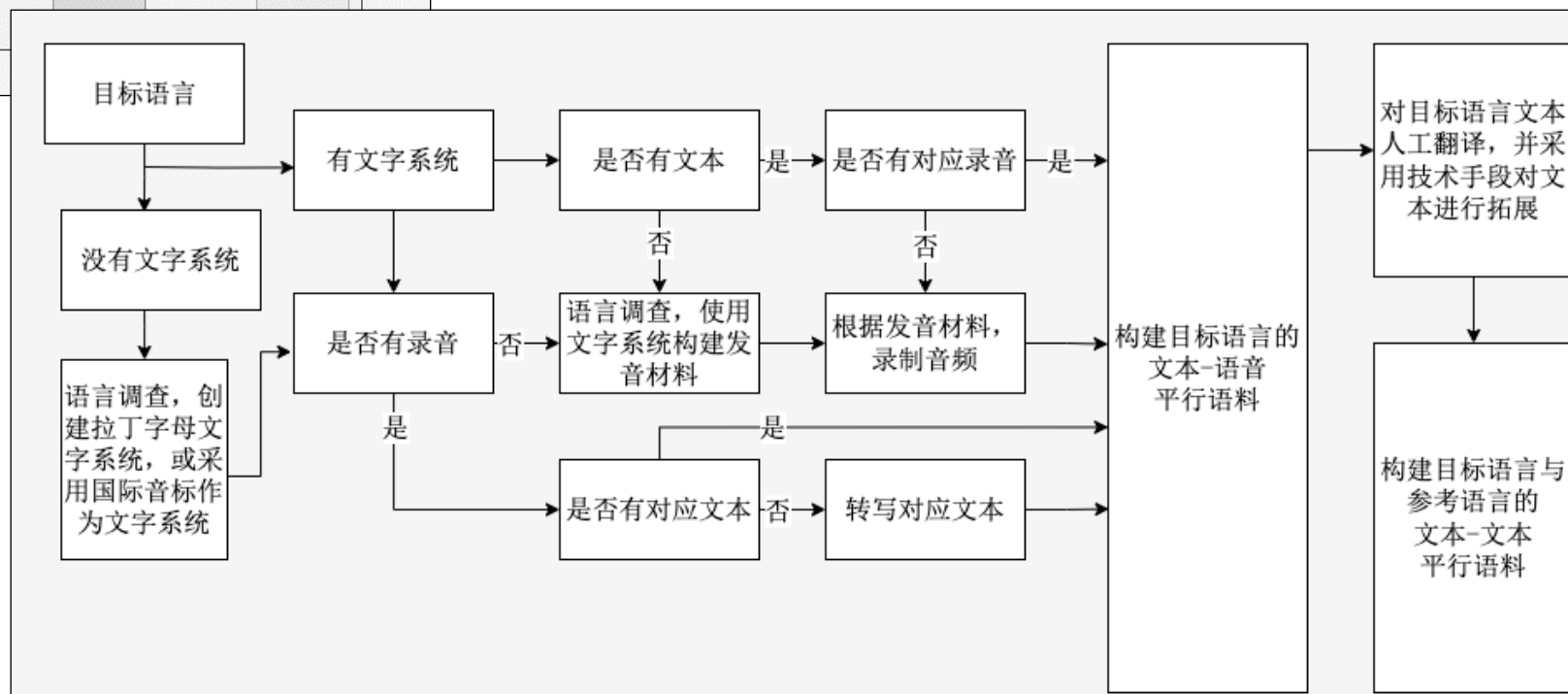
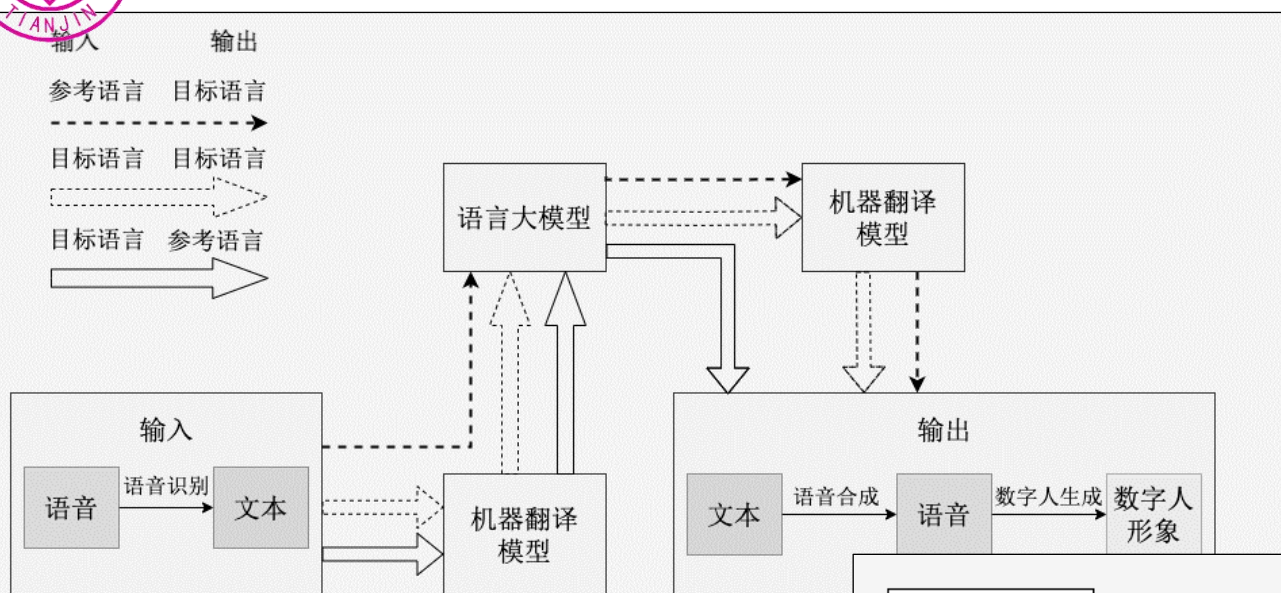


- **Based on the AI technologies displayed above, my team is engaged in 2 large projects:**
- **1-AI-driven Language preservation(especially for endangered languages or dialects)**
- **A huge project: full-chain integrated solution**





A variety of methods and approaches



- 2-AI-driven Reconstruction of history/language/oral culture
- Technology of digital human



- **Synthesizing oral culture audios/videos**
- **Chanting of Weidong Zhang (张卫东)**
- **Chanting of official Mandarin (国子监官韵吟诵)**
- **ancient Chinese prose (古文)**
- **Chinese poetry (诗词)**



- **《论语·学而》：子曰：“学而时习之，不亦说乎？有朋自远方来，不亦乐乎？人不知而不愠，不亦君子乎？”**
- **The Analects of Confucius; Mencius**
- **Original sound and synthesized sound**



- ☆ **Reconstructing historic audio and video**
- **Hu Shih (胡适, 1891-1962), Yuan-Ren Chao(赵元任, 1892-1982)**
- **New Year's speech, 1935**
- 新年前的两日，我正在作长途的旅行。寂寞的旅途是我最欢迎的，因为平常某日有应作的事，有不能不见的客，很少有整天可以自由用来胡思乱想的。这两天在火车上，因为要替《大公报》写新年的第一篇星期论文。我一路上想的是“我盼望我们这个国家在这新开始的一年里可以做到的几件什么事?”.....
- **Yuan-Ren Chao:**
- **Widen history window.**
- **Ethical standards for artificial intelligence issued by Ministry of Science and Technology.**
(https://www.most.gov.cn/kjbgz/202109/t20210926_177063.html)

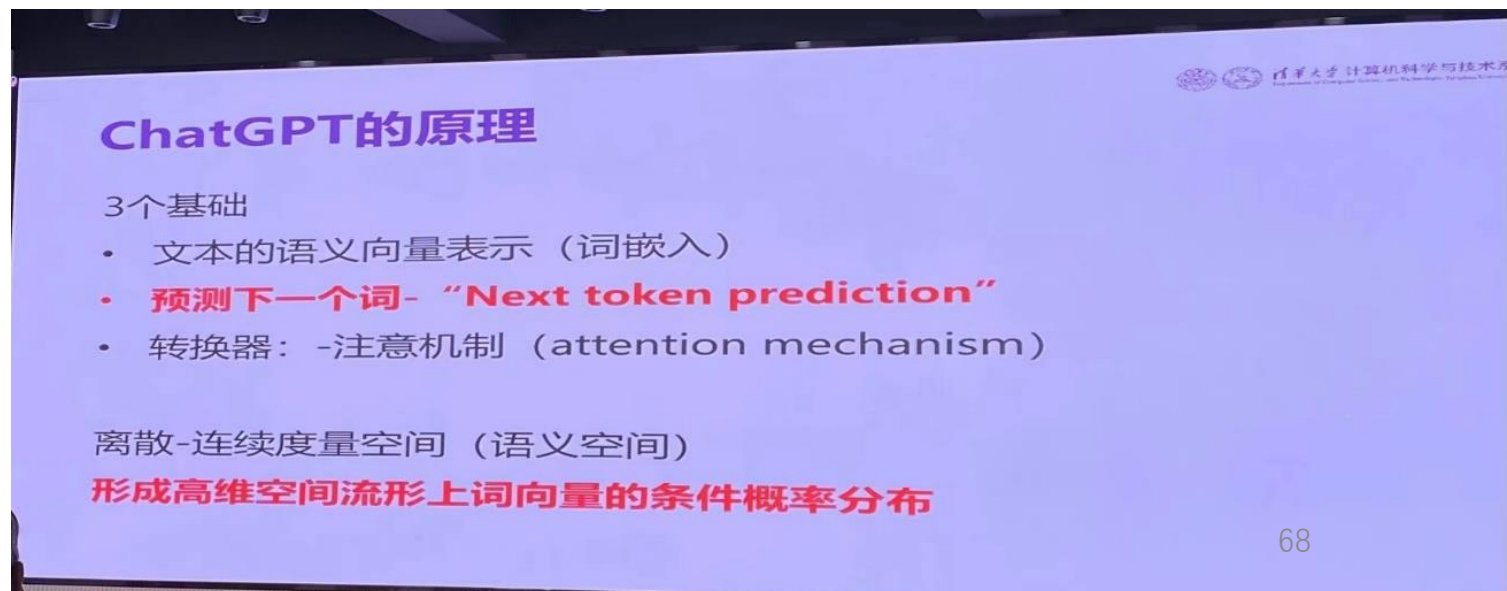




- **Essential mechanism of Large Language Model:**
- **3 basic technology:**
- **(1) Semantic Vector Representation of Text (Word embedding)**
- **(2) Next token prediction**
- **(3) Transformer (attention mechanism)**

- **Discrete Continuous Metric Space (Semantic Space)**

- **Conditional probability distribution of word vectors on high-dimensional space manifolds**





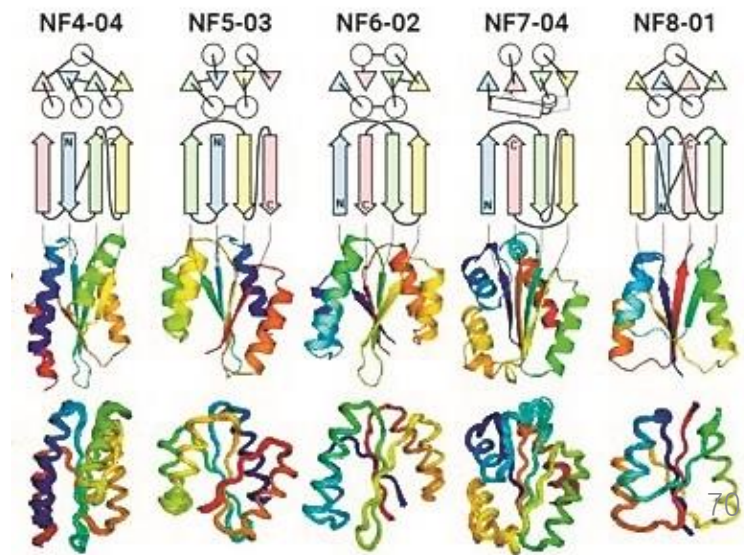
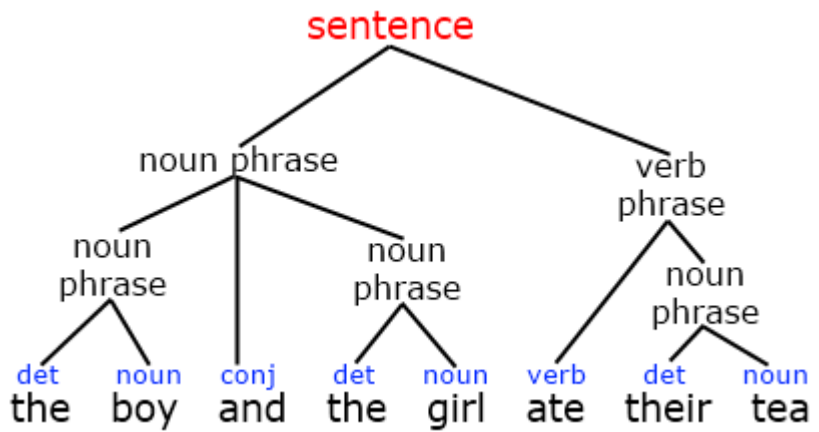
- **Chomsky vs. Hinton**
- **“basically high-tech plagiarism”**



- **AIGC(Artificial Intelligence Generated Content)**
- **Generative linguistics (Transformational-generative grammar)**



- **AlphaFold 3**
- “AlphaFold is a model that is capable of high-accuracy prediction of complexes containing nearly all molecular types present in the Protein Data Bank.”
- **amino acids: 20**
- **protein: Proteins have complex and diverse spatial conformations. Their primary structure is a linear arrangement of amino acids, while secondary structures such as alpha helices and beta folds are further folded to form tertiary structures. Multiple subunit combinations form quaternary structures. This precise structure endows proteins with specific functions.**
- **AlphaFold inspired Large Language Model**
- **Interesting!**





香港城市大學
City University of Hong Kong



南開大學
Nankai University

Thanks for listening.

Looking forward to your comments and suggestions.

Looking forward to your comments and suggestions.

Qibin Ran
1@ranqibin.com