Language and cognitive diversity and global wellbeing The Behemoth Language Bias Damián Blasi









The Branco Weiss Fellowship Society in Science



I research **language diversity** over broad scales of time and space through a multidisciplinary lens



appointed ner to tollow next ourse, as the rest had done bef began in this manner. nom , to acquaint e mockery, d (not in ieft, e Gentlewodeuoute reliwill yeelde for afure and reecular vnderdy he or thee commonly ns (at least) like notoinucnters of mes, as thinore wife and many other; worthorva every beft of



Why should anyone <u>besides us</u> care or invest into language diversity research?



ARGUMENT I **THE ROMANTIC SPIEL**

Language diversity is a rightful object of human curiosity and awe

15 MIN READ JUNE 1, 2023

This Ancient Language Has the Only Grammar Based Entirely on the Human Body

An endangered language family suggests that early humans used their bodies as a model for reality



20 amazing words that don't exist in English but really should

Shelby Slauer Updated Oct 9, 2019, 5:51 PM GMT+2

"Schadenfreude" — German



> <

 \bigcirc \bigcirc \bigcirc

Botswana has a language made entirely of clicks. **Click languages or Khoisan** languages use solely clicking sounds to communicate. Typically, they are made up of four clicks, but the Southern languages use a fifth ("kiss" click) as well.



Ide" once in a while. Alan Crowhurst/Getty Images



ARGUMENT II THE STUDY OF LANGUAGED HUMANS

Some aspects of human behavior and cognition are influenced by the language(s) used by the individual. No true science of the languaged human without sampling across languages

Trends in **Cognitive Sciences**

Feature Review

Over-reliance on English hinders cognitive science

Damián E. Blasi ^{(1,2,3,*,@} Joseph Henrich, ¹ Evangelia Adamou, ⁴ David Kemmerer, ^{5,6} and Asifa Majid ^{7,*,@}



COMPUTATIONS AND REPRESENTATIONS





SOCIAL COGNITION

CAUSAL COGNITION



VISUAL AND AUDITORY ATTENTION











ARGUMENT III THE ETHICAL IMPERATIVE Languages become dormant (i.e. lose all active users) at a concerning rate, and documenting this diversity should be at the fore of our efforts concerning language

Former rainforest environment of a Bagyeli population in Mimbosso, Centre-South province, Cameroon (June 2024)

-





THE ROMANTIC SPIEL

Members of affluent societies who are capable of consuming language diversity as entertainment (<16%*)





* % of population with Development Index 2024



THE ROMANTIC SPIEL

Members of affluent societies who are capable of consuming language diversity as entertainment (<16%*)

THE STUDY OF LANGUAGED HUMANS

Members of affluent societies who are capable of reaping the benefits of psychology and cognitive science (<16%*)





* % of population with Development Index 2024



THE ROMANTIC SPIEL

Members of affluent societies who are capable of consuming language diversity as entertainment (<16%*)

THE STUDY OF LANGUAGED HUMANS

Members of affluent societies who are capable of reaping the benefits of psychology and cognitive science (<16%*)

THE ETHICAL IMPERATIVE

Users of endangered languages (<20%[†])



* % of population with Development Index 2024

⁺ approx. % of population who do not speak the top 100 largest languages by Ethnologue



MY ARGUMENT HERE **THE BEHEMOTH LANGUAGE BIAS**

English and a few other large and influential languages ("behemoth languages") serve, in technology, medicine, education, law, etc. as prototypes models for language in general, *leading to a measurably poorer performance* for users of other languages

MY ARGUMENT HERE **THE BEHEMOTH LANGUAGE BIAS**

English and a few other large and influential languages ("behemoth languages") serve, in technology, medicine, education, law, etc. as prototypes models for language in general, *leading to a measurably poorer performance* for users of other languages

This bias affects over 60% of the world's population, involving large vital languages with hundreds of millions of users



Useful speech cues for detecting Parkinson's in English lead to misdiagnose in vowel-rich and tonal languages Pinto et al. (2017)



Direct translation of questions in the PISA evaluations is fair in most large languages but in many minority languages leads to more complex sentences and lower scores

El Masri, Baird & Graesser (2016)



Language models are fundamental for human-like cognition in AI but they work better in morphology-poor and fixed-word order languages like English Schwartz et al. (2020), Park et al. (2020), Gerz et al. (2020)



Public health communication during COVID was clear in English but the same guidelines applied to other languages yield difficult to understand advice Blasi et al. (2021)

If a resource is language-based, then very likely it's <u>not language-independent</u>

The world's population is aging, which increases the burden of neurodegenerative diseases (Alzheimer's, Parkinson's,

dementia)



ur World in Data	
Villion	
Median age in 2100: 41.6 years Median age in 2075: 39 years	
Median age in 2018:	
Median age in 1950:	
23.6 years	
205 207	
bor May Reser	
nor wax noser.	

The world's population is aging, which increases the burden of neurodegenerative diseases (Alzheimer's, Parkinson's, dementia)





There is a race for developing fast, cheap, and reliable markers for detecting these diseases early



Linguistic biomarkers are extremely promising and developing at a breakneck speed



PROBABILISTIC DIAGNOSIS



In English, individuals afflicted by Alzheimer's disease (AD) use more pronouns in contrast to healthy individuals

> Healthy individual: the dog is chasing a cat

AD patient: he is chasing it



In English, individuals afflicted by Alzheimer's disease (AD) use more pronouns in contrast to healthy individuals

> Healthy individual: the dog is chasing a cat

AD patient: he is chasing it

...maybe because AD patients have problems with *perspective taking* (Bittner et al. 2022)





Convenient basis for a <u>universal biomarker</u>: pronoun frequency

In English, individuals afflicted by Alzheimer's disease (AD) use more pronouns in contrast to healthy individuals

> Healthy individual: the dog is chasing a cat

AD patient: he is chasing it

...maybe because AD patients have problems with *perspective taking* (Bittner et al. 2022)



However, the **exact opposite pattern** is found in speakers of Bengali (বাংলা): they use less, no more, pronouns

However, the **exact opposite pattern** is found in speakers of Bengali (বাংলা): they use less, no more, pronouns

English personal pronouns (NOM)

			Nominative
First-person	First-person Singular		
	Ρ	we	
	Singular	Standard	you
Second-person		Poetic/dialectal	thou
	Ρ	you	
Third-person	Singular	Masculine	he
		Feminine	she
		Neuter	it
		Epicene	they
	Plural		they

However, the **exact opposite pattern** is found in speakers of Bengali (বাংলা): they use less, no more, pronouns

English personal pronouns (NOM)

			Nominative
First-person	First-person Singular		
	Ρ	we	
	Singular	Standard	you
Second-person		Poetic/dialectal	thou
	Р	you	
Third-person	Singular	Masculine	he
		Feminine	she
		Neuter	it
		Epicene	they
	Plural		they

Bengali personal pronouns (NOM)

Subject	Proximity	Honor	Singular	Plural
1		আমি (<i>ami</i> , I)	আমরা (<i>amra</i> , we)	
VF		তুই (<i>tui</i> , you)	তোরা (<i>tora,</i> you)	
	2		তুমি (<i>tumi</i> , you)	তোমরা (<i>tomra,</i> you)
F		Р	আপনি (<i>apni,</i> you)	আপনারা (<i>apnara,</i> you)
		F	ଏ (<i>e,</i> he/she)	এরা (<i>era,</i> they)
	н 3 Т	Р	ইনি (<i>ini,</i> he/she)	এঁরা (<i>ẽra,</i> they)
		I	এটি/এটা (<i>eţi/eţa,</i> it)	এগুলো (<i>egulo</i> , these)
		F	ર (<i>o,</i> he/she)	ওরা (<i>ora,</i> they)
3		Р	উনি (<i>uni,</i> he/she)	ওঁরা (<i>õra,</i> they)
		I	ওটি/ওটা (<i>oţi/oţa,</i> it)	ওগুলো (<i>ogulo,</i> those)
	F	সে (<i>she,</i> he/she)	তারা (<i>tara,</i> they)	
	E	Р	তিনি (<i>tini,</i> he/she)	তাঁরা (<i>tãra,</i> they)
		I	সেটি/সেটা (<i>sheţi/sheţa</i> , it)	সেগুলো (<i>shegulo,</i> those)

However, the exact opposite pattern is found in speakers of Bengali (বাংলা): they use less, no more, pronouns

English personal pronouns (NOM)

			Nominative
First-person	Sir	I	
	Plural		we
	Singular	Standard	you
Second-person		Poetic/dialectal	thou
	Ρ	you	
Third-person	Singular	Masculine	he
		Feminine	she
		Neuter	it
		Epicene	they
	Plural		they

Bengali personal pronouns (NOM)

Subject	Proximity	Honor	Singular	Plural
	1		আমি (<i>ami</i> , I)	আমরা (<i>amra</i> , we)
VF		<u>ত</u> ুই (<i>tui</i> , you)	তোরা (<i>tora,</i> you)	
	2	F	তুমি (<i>tumi</i> , you)	তোমরা (<i>tomra,</i> you)
		Р	আপনি (<i>apni,</i> you)	আপনারা (<i>apnara,</i> you)
	H	F	ଏ (<i>e,</i> he/she)	এরা (<i>era,</i> they)
		Р	ইনি (<i>ini,</i> he/she)	এঁরা (<i>ẽra,</i> they)
		I	এটি/এটা (<i>eţi/eţa,</i> it)	এগুলো (<i>egulo</i> , these)
		F	૩ (<i>o,</i> he/she)	ওরা (<i>ora,</i> they)
3	т	Р	উনি (<i>uni,</i> he/she)	उँत्रा (<i>õra,</i> they)
		I	ওটি/ওটা (<i>oţi/oţa,</i> it)	ওগুলো (<i>ogulo,</i> those)
	F	সে (<i>she,</i> he/she)	তারা (<i>tara,</i> they)	
	E	Р	তিনি (<i>tini,</i> he/she)	তাঁরা (<i>tãra,</i> they)
			সেটি/সেটা (<i>sheţi/sheţa</i> , it)	সেগুলো (<i>shegulo,</i> those)

This suggests it might not be a problem of perspective taking but just AD patients using "easy" pieces of language for communication

Even with the limited number of studies on non-behemoth languages, this issue is **widespread** in the diagnose of neurodegenerative diseases

Linguistic biomarkers rarely generalise across languages

Examples of cross-linguistic differences

DISORDER	LANGUAGES	STRUCTURAL CONTRAST	DISTINCT MARKER
	English	Greater phonetic and lesser morphosyntactic complexity	Phonetic distortions as most salient symptom
nfvPPA	Italian	Lesser phonetic and greater morphosyntactic complexity	Distinct syntactic alterations
	English	Alphabetic script (letters represent phonemes)	High prevalence of surface dysgraphia
svPPA	Chinese	Logographic script (logograms convey semantic or phonological information)	Low prevalence of surface dysgraphia
WORDS	English	Less diverse morphosyntactic patterns	Frequent sentence repetition deficits
IvPPA	German	More diverse morphosyntactic patterns	Infrequent sentence repetition deficits
	English	Simpler pronominal system	Overuse of pronouns
	Bengali	More complex pronominal system	Underuse of pronouns
	Spanish	Verb-framed language with rich verb vocabulary	Selective action-verb deficits
PD R	Dutch	Satellite-framed language with fewer verbs	Non-selective action-verb deficits

As language diversity researchers, what can we do to tackle the Behemoth Language Bias?





I. Get involved with the relevant parties



"And now, a bit of feedback from our silent partner." Credit: Adobe Stock

I. Get involved with the relevant parties



"And now, a bit of feedback from our silent partner." Credit: Adobe Stock



Shameless staged picture of me at the UNESCO headquarters in Paris. Disclosure: I do not represent my country in any official capacity.

II. Be demographically aware of the future



Credit: r/MapPorn user Sppoderman89 based on US Census Bureau and the World Bank

III. Enhance academic credit

Systematic Inequalities in Language Technology Performance across the World's Languages

Damián Blasi Harvard University dblasi@fas.harvard.edu Antonios Anastasopoulos George Mason University antonis@gmu.edu Graham Neubig Carnegie Mellon University gneubig@cs.cmu.edu

Circle size reflects number of publications in the language, proxying magnitude of R&D in language technologies



Contemporary language technologies, are (1) heavily centered on English and other Behemoth Languages, and (2) explained by economic centrality rather than demographic need





III. Enhance academic credit

The number of languages in a publication is not associated with its relative citation rate over time



→ No marginal incentive for researchers to further the language diversity of their data and models











The **Branco Weiss** Fellowship Society in Science

