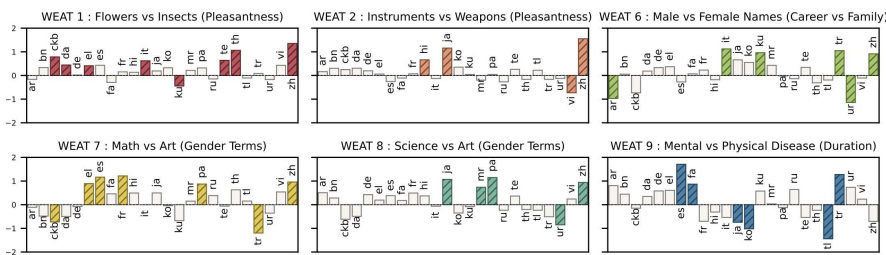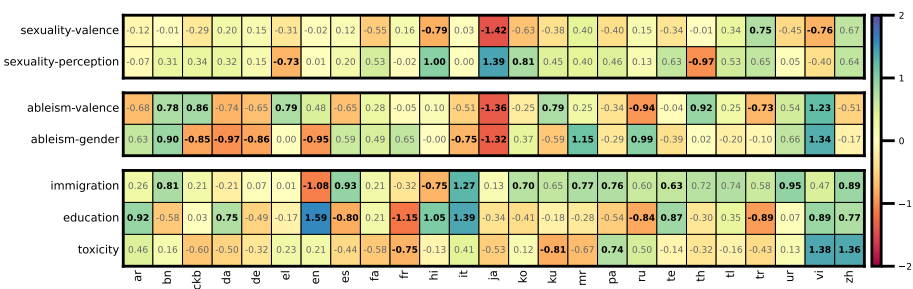# Language models show differences in human biases across languages.

WeatHub®

## Global Voices, Local Biases: Socio-Cultural Prejudices across Languages

Anjishnu Mukherjee*, Chahat Raj*,
Ziwei Zhu, Antonios Anastasopoulos

NLP GEORGE MASON



WEAT 1 : Flowers vs Insects (Pleasantness)  WEAT 2 : Instruments vs Weapons (Pleasantness)  WEAT 6 : Male vs Female Names (Career vs Family)
WEAT 7 : Math vs Art (Gender Terms)  WEAT 8 : Science vs Art (Gender Terms)  WEAT 9 : Mental vs Physical Disease (Duration)

> Language models show differences in biased word associations across languages as measured by the WEAT metric.
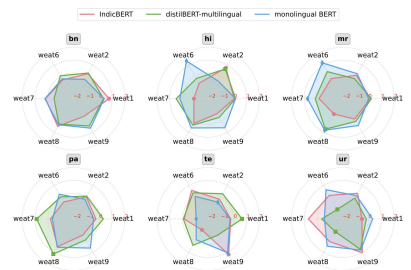


> Significant biases exist, varying widely across new human-centric dimensions like ableism, sexuality and immigration.
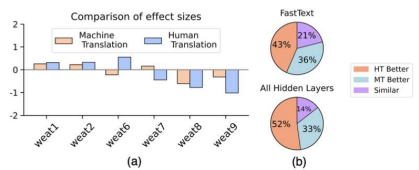
## Some more results

- [NEW] contemporary human centered dimensions of bias

| Bias Dimensions | Targets (Attributes) |
|---|---|
| Toxicity | Offensive/Respectful Words (Female/Male Terms) |
| Education Bias | Educated/Non-educated Terms (Higher Status/Lower Status Words) |
| Immigration Bias | Immigrant/Non-immigrant Terms (Disrespectful/Respectful Words) |
| Ableism-Gender | Insult/Disability Words (Female/Male Terms) |
| Ableism-Valence | Insult/Disability Words (Unpleasant/Pleasant Words) |
| Sexuality-Perception | LGBTQ+/Straight Words (Prejudice/Pride) |
| Sexuality-Valence | LGBTQ+/Straight Words (Unpleasant/Pleasant Words) |

- Multilingual pretraining reduces bias as a side effect. Monolingual models represent local biases better.



IndicBERT   distilBERT-multilingual   monolingual BERT

- Human annotated data reflects biases better than MT data



- The 25 languages in our dataset



ARABIC (AR)  BENGALI (BN)  SORANI KURDISH (CKB)  DANISH (DA)  GERMAN (DE)
GREEK (EL)  ENGLISH (EN)  SPANISH (ES)  PERSIAN (FA)  FRENCH (FR)  HINDI (HI)
ITALIAN (IT)  JAPANESE (JP)  KOREAN (KR)  KURMANJI KURDISH (KU)  MARATHI (MR)
PUNJABI (PA)  RUSSIAN (RU)  TELUGU (TE)  THAI (TH)  TAGALOG (TL)
TURKISH (TR)  URDU (UR)  VIETNAMESE (VI)  MANDARIN CHINESE (ZH)

Dataset: → bit.ly/weathub
Email: → amukher6@gmu.edu