# Sentiment and Hate Speech datasets for more than 14 African languages

HAUSANLP, MASAKHANENLP, ICT4D, UHH LT GROUP, GOOGLE RESEARCH, UNIVERSITY OF PRETORIA, DATA SCIENCE FOR SOCIAL IMPACT     **Contact us!**
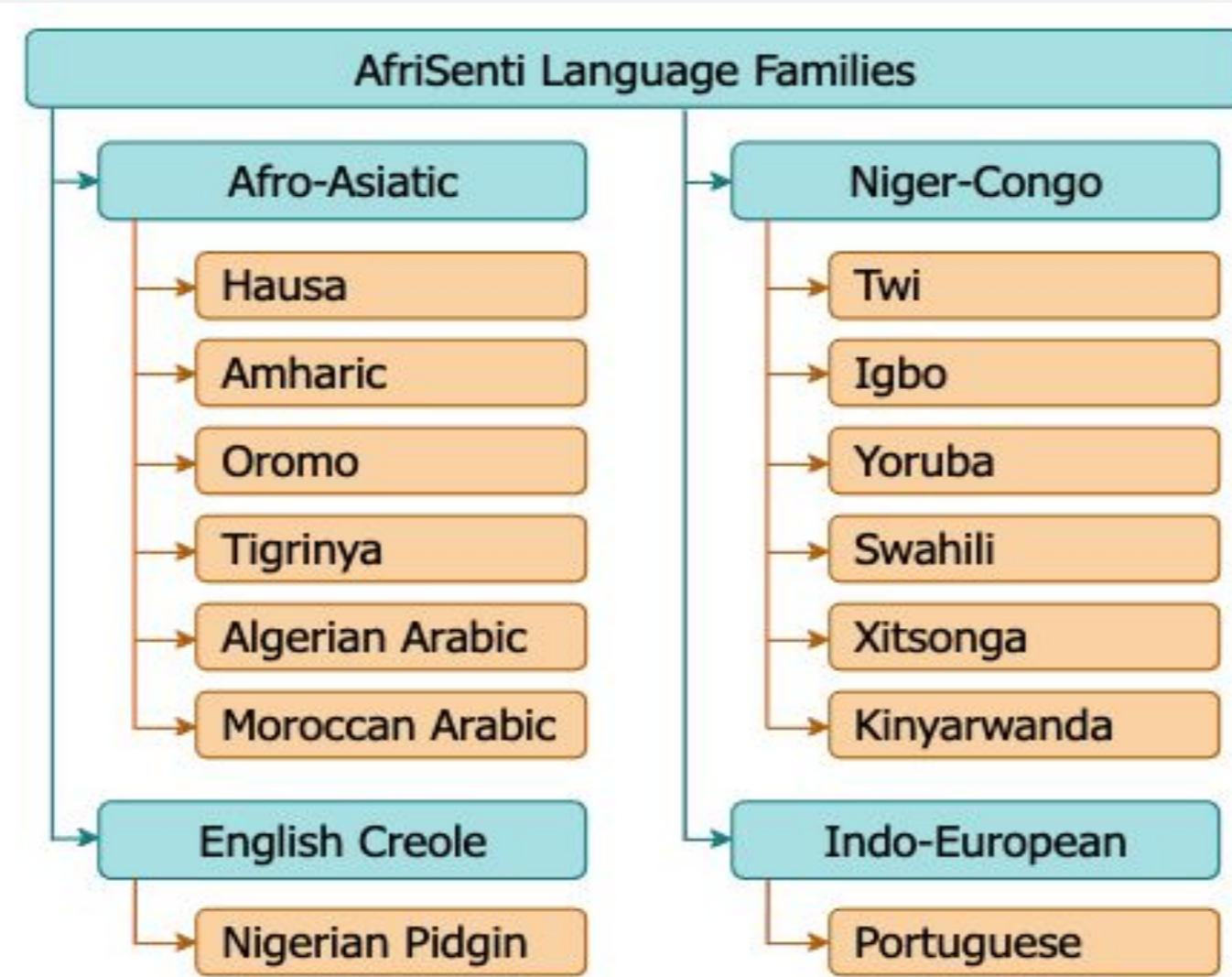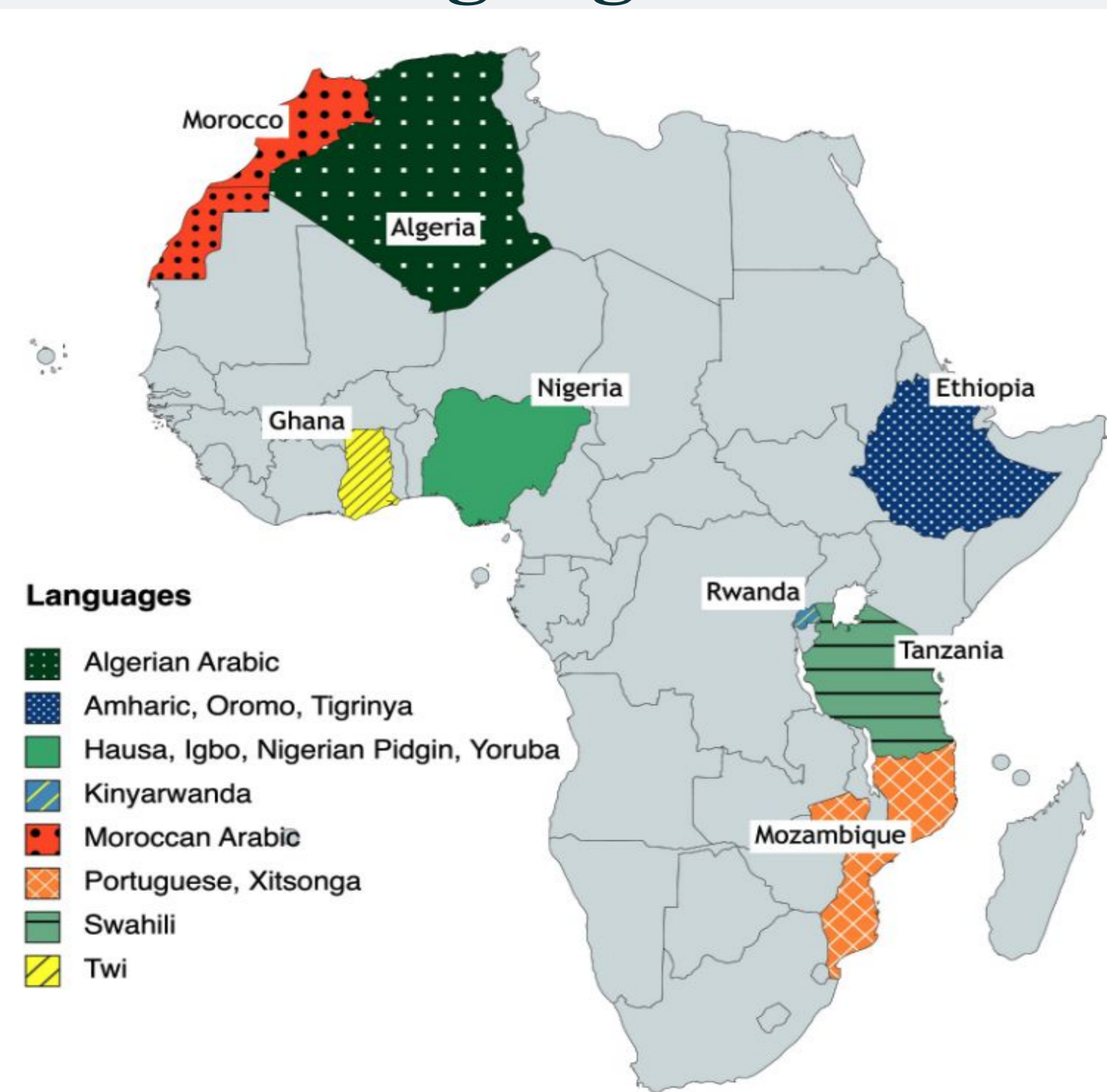
1. Dr. Seid Muhie Yimam
seid.muhie.yimam@uni-hamburg.de

2. Dr. Shamsuddeen Hassan Muhammad
shamsuddeen2004@gmail.com

**Deep Learning Indaba 2024
September 1-7 2024, Senegal, Dakar**

## AfriSenti: A Twitter Sentiment Analysis Benchmark for African Languages



The language Family (in green) of each language (in yellow) included in AfriSenti.

### Abstract:

Africa is home to over 2,000 languages from more than six language families and has the highest linguistic diversity among all continents. These include 75 languages with at least one million speakers each. Yet, there is little NLP research conducted on African languages. Crucial to enabling such research is the availability of high-quality annotated datasets. In this paper, we introduce AfriSenti, a sentiment analysis benchmark that contains a total of $>\$110,000$ tweets in 14 African languages (**Amharic, Algerian Arabic, Hausa, Igbo, Kin-yarwanda, Moroccan Arabic, Mozambican Por-tuguese, Nigerian Pidgin, Oromo, Swahili, Tigrinya, Twi, Xitsonga, and Yorùbá**) from four language families. The tweets were annotated by native speakers and used in the AfriSenti-SemEval shared task.

| Lang. | Tweet | Sentiment |
|-------|-------|-----------|
| amh | ያ ጨካኝ አረመኔ ታስሮ ይኸው ካቴና ጉብቶለታል ይሉናል። ቆይ አስረው የጀበና ቡና አየጋበቱት ነው እንደ? | negative |
| arq | الشروق هذه من خرجت وهي تباع تهديل، مستوى منحط وشعبري @user .... | negative |
| ary | واش بغيتهم يداو يكرفسو على العادي والبادي عاد تبقاو أنتا على خاطركم | negative |
| ary | rabi ykhali alhbiba makayn ghir nachat o chi machat | positive |
| hau | @USER Aunt rahma i luv u wallah irin totally dinnan | positive |
| ibo | akowaro ya ofuma nne kai daalu nwanne mmadu | positive |
| kin | @user Ariko akokanu ngo inyebebe unyujijemo sisawa wangu | negative |
| orm | @user Jawaar Kenya OMN haala akkamiin argachuu dandeenya | neutral |
| por | Honestidade é algo que não se compra. Infelizmente a humanidade esqueceu disso por causa das suas ambições. | positive |
| pcm | E don tay wey I don dey crush on this fine woman ... | positive |
| swa | Asante sana watu wa Sirari jimbo la Tarime vijijini Huu ni Upendo usio na Mashaka kwa Mbunge wenu John Heche | positive |
| tir | @user ከመኸረኩም እንተተኾይኑ፡ንሕውሓት ነዝም ውሕዶ ቁጽሮም አባ ምጥፋአ ይሕሰ ኩም! | negative |
| tso | @user @user Yu , tindzava ? Tsika mbangui mpfana e nita ku despro-gramara | negative |
| twi | messi saf den check en bp na wo kwame danso wo di twe da kor aaa na wawu | negative |
| yor | onírèégbè aláàdúgbò ati olójúkòkòrò | negative |



Label distributions for the different AfriSenti datasets (i.e., number of **positive**, **negative**, and **neutral** tweets).

## AfriHate: Hate Speech dataset for Africa language – Upcoming

| Language | Country |
|----------|---------|
| Hausa, Igbo, Yoruba and Nigerian Pidgin | Nigeria |
| Amharic, Tigrinya, Oromo, Somali | Ethiopia |
| Swahilli | Kenya |
| Algerian Arabic | Algeria |
| Moroccan Arabic | Morocco |
| Sudanese Arabic | Sudan |
| Twi | Ghana |
| Mozambican Portuguese | Mozambique |
| Kinyarwanda | Rwanda |
| IsiZulu, Afrikaan, IsiXhosa | South Africa |

### Abstract:

Online hate is an escalating issue that negatively affects users and disrupts online communities, leading to psychological harm and potential offline violence. In Africa, efforts to combat this issue are often focused on **high-profile individuals** and rely heavily on **manual moderation**, making them inefficient and impractical for the wider population. African languages lack adequate natural language processing (NLP) tools, resulting in **blanket interventions** like **content removal** based on **keywords** without context consideration. To address these challenges, we introduce **AfriHate**, a comprehensive labeled dataset for identifying hate and abusive language in **18 African languages**. This dataset will facilitate the development of **classification models** that enhance **automatic moderation capabilities**. It can be utilized by **platform owners**, **peacebuilders**, **community service platforms**, and others to foster **innovation in NLP for African languages**.
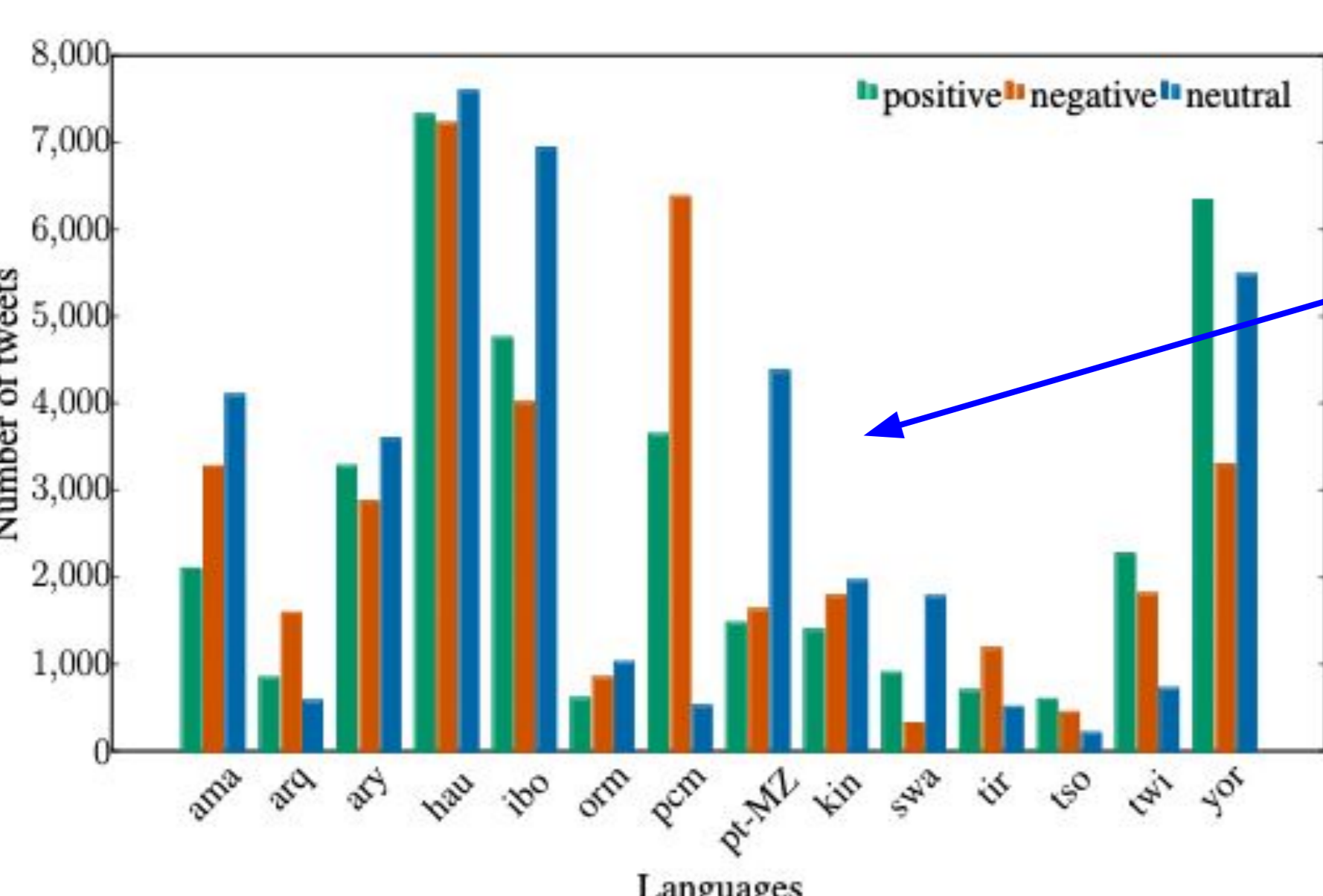
### Annotation Strategy

- Multi-layer annotation by five annotators per tweet
- Covers 18 African languages
- Uses Fleiss kappa for reliability
- Ensures cultural and context sensitivity

### Potential Applications

- **Develop models for automatic moderation**
  - Social media platforms - Peacebuilding organizations
  - Community service platforms
- **Supports research in:**
  - Sentiment analysis  - Socio-linguistic studies

### Impact

- Enhances NLP resources for African languages
- Promotes inclusivity and safer online spaces
- Supports linguistic diversity and understanding
- Lays foundation for future NLP advancements