



Hausa-NLP  
Research Group

U.PORTO



# NaijaSenti: A Nigerian Twitter Sentiment Corpus for Multilingual Sentiment Analysis

**Shamsuddeen Muhammad(@shmuhammad)**  
University of Porto, Bayero University, Kano - Nigeria, MasakhaneNLP, HausaNLP

# NaijaSenti: A Nigerian Twitter Sentiment Corpus for Multilingual Sentiment Analysis

**Shamsuddeen Hassan Muhammad<sup>1,2\*+</sup>, David Ifeoluwa Adelani<sup>3\*</sup>, Sebastian Ruder<sup>4</sup>  
Ibrahim Sa'id Ahmad<sup>5+</sup>, Idris Abdulmumin<sup>6\*+</sup>, Bello Shehu Bello<sup>5+</sup>, Monojit Choudhury<sup>7</sup>  
Chris Chinenyenye Emezue<sup>8\*</sup>, Saheed Salahudeen Abdullahi<sup>10+</sup>  
, Anuoluwapo Aremu<sup>11\*</sup>, Alipio Jeorge<sup>1,2</sup> Pavel Brazdil<sup>1</sup>**

<sup>1</sup> LIAAD - INESC TEC, <sup>2</sup>Faculty of Sciences-University of Porto, Portugal,

<sup>3</sup>Spoken Language Systems Group (LSV), Saarland University, Germany, <sup>4</sup>Google Research

<sup>5</sup>Faculty of Computer Science and Information Technology, Bayero University, Kano, Nigeria

<sup>6</sup>Department of Computer Science, Ahmadu Bello University, Zaria, Nigeria, <sup>7</sup>Microsoft Research India,

<sup>8</sup>Technical University of Munich, Germany, <sup>9</sup>Clear Global, <sup>10</sup>Kaduna state University

\*MasakhaneNLP, +HausaNLP

# Table of contents

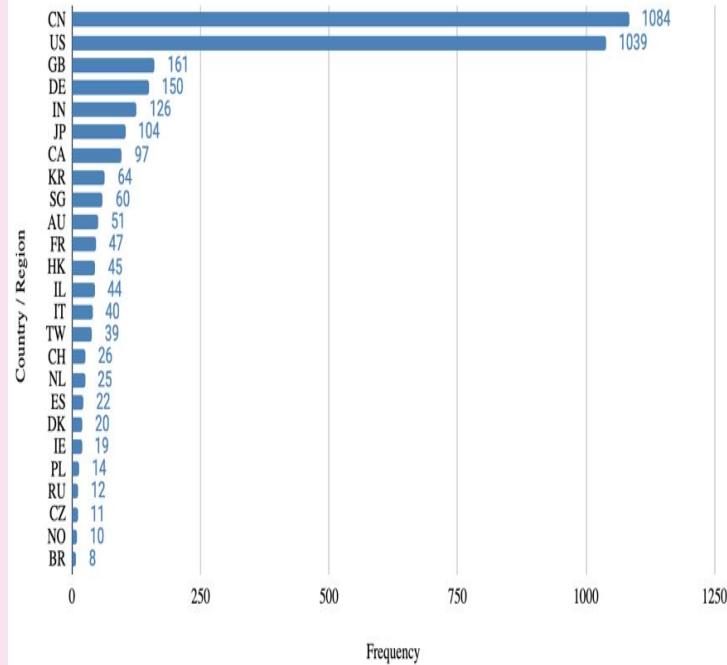
- |    |  |
|----|--|
| 01 | <b>Motivation</b>                          |
| 02 | <b>Sentiment Analysis for low-resource</b> |
| 03 | <b>Data collection and Annotation</b>      |
| 04 | <b>Benchmark Experiments</b>               |

# Motivation

30% of all living  
languages today are  
African languages

(Ethnologue)

Number of Submissions per Country/Region (Contact Author)



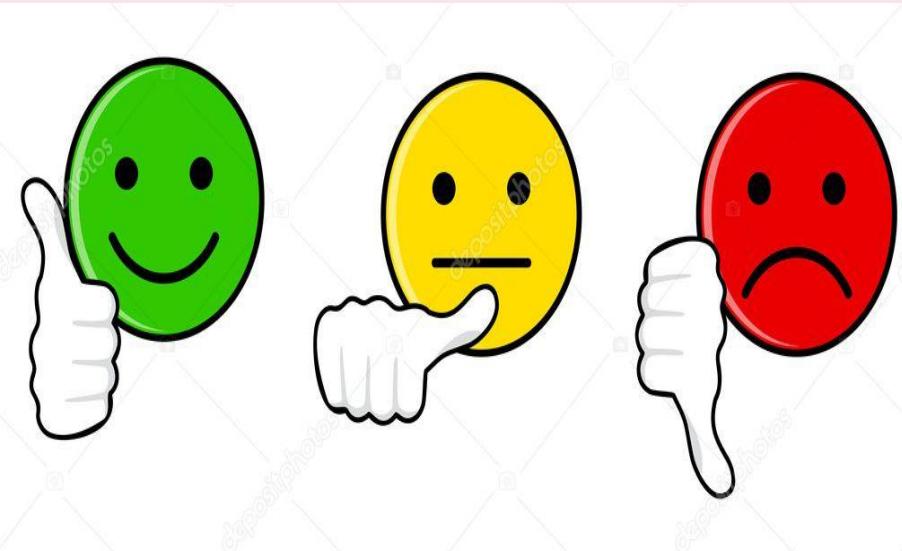
<https://acl2020.org/blog/general-conference-statistics/>

02/04



## Sentiment Analysis for low-resource

# Sentiment Analysis for Low-resource



SA involves finding opinion or sentiment in a text !

- Considerable increase of interest
- English centric
- Low-resource:
  - Lack of available annotated data
  - Translation
    - Focus > English
    - English > Focus

# Translation Alter Sentiment

## Development of a General Purpose Sentiment Lexicon for Igbo Language

**Emeka Ogbuju**

Department of Computer Science, Federal University Lokoja, Nigeria  
emeka.ogbuju@fulokoja.edu.ng

**Moses Onyesolu**

Department of Computer Science, Nnamdi Azikiwe University, Awka, Nigeria  
mo.onyesolu@unizik.edu.ng

### Abstract

## How Translation Alters Sentiment

Saif M. Mohammad

National Research Council Canada

SAIF.MOHAMMAD@NRC-CNRC.GC.CA

Mohammad Salameh

University of Alberta

M.SALAMEH@UALBERTA.CA

Svetlana Kiritchenko

National Research Council Canada

SVETLANA.KIRITCHENKO@NRC-CNRC.GC.CA

## SEMANTIC ENRICHMENT OF NIGERIAN PIDGIN ENGLISH FOR CONTEXTUAL SENTIMENT CLASSIFICATION.

**Wuraola Fisayo Oyewusi**

Data Science Nigeria  
Lagos, Nigeria  
wuraola@datasciencenigeria.ai

**Olubayo Adekanmbi**

Data Science Nigeria  
Lagos Nigeria  
olubayo@datasciencenigeria.ai

**Olaelekan Akinsande**

Data Science Nigeria Lagos,Nigeria  
olaelekan@datasciencenigeria.ai



**Jabeer Muhammad**

3 March · 2 ·

Masha Allah,Allah ya Azur tani da samin karuwa ta da Namiji  
Masha Allah

Uwa Iafiya Danma Shima Iafiya

Allah ya rayamanashi Amin akan tafarkin { SHUGABA  
S.A.W. }

Masha Allah, may Allah bless me with getting my prostitute  
and my husband Masha Allah  
Mother is fine, son is fine too  
May God grant him long life on the path of {leader S.A.W.}

· Hide Translation · Rate this translation

...

# Previously Annotated Datasets

Dataset	Language	Open-source	Annotated/translated	Code-mixed	Source
Abubakar et al. (2021)	Hausa	✗	annotated	✓	Twitter
Ogbuju and Onyesolu (2019)	Igbo	✗	translated	✗	General
Umoh et al. (2020)	Igbo	✗	annotated	✗	General
Oyewusi et al. (2020)*	Pidgin	✗	annotated/translated	✓	Twitter
Orimaye et al. (2012)	Yorùbá	✓	annotated	✓	Youtube
Iyanda and Abegunde (2019)	Yorùbá	✗	annotated	✗	General
<b>Ours</b>	Hausa, Igbo, Yorùbá, Pidgin	✓	annotated	✓	Twitter

Table 1: Summary of datasets used in six existing datasets on sentiment analysis in four major Nigerian languages in comparison to ours. \*: The provided URL is no longer accessible.

# NaijaSenti Contributions

**1st**



## Labelled Dataset

Labelled tweets in four Nigerian languages (Hausa, Yoruba, Igbo, Pidgin).

**2nd**



## Unlabelled Dataset

Released the largest Twitter corpus in four languages

**3rd**



## Sentiment Lexicon

Released annotated sentiment lexicon and stopword in Hausa, Yoruba and Igbo

**4th**



## Trained Models

Train various PLM and Released the code

03/04



## Data collection and Annotation

# Dataset collection

## Twitter Academic API



- Support 70 languages ( only one African lang: Amharic )
- None among the Nigerian languages is supported

## Stopwords, emoji, sentiment lexicon



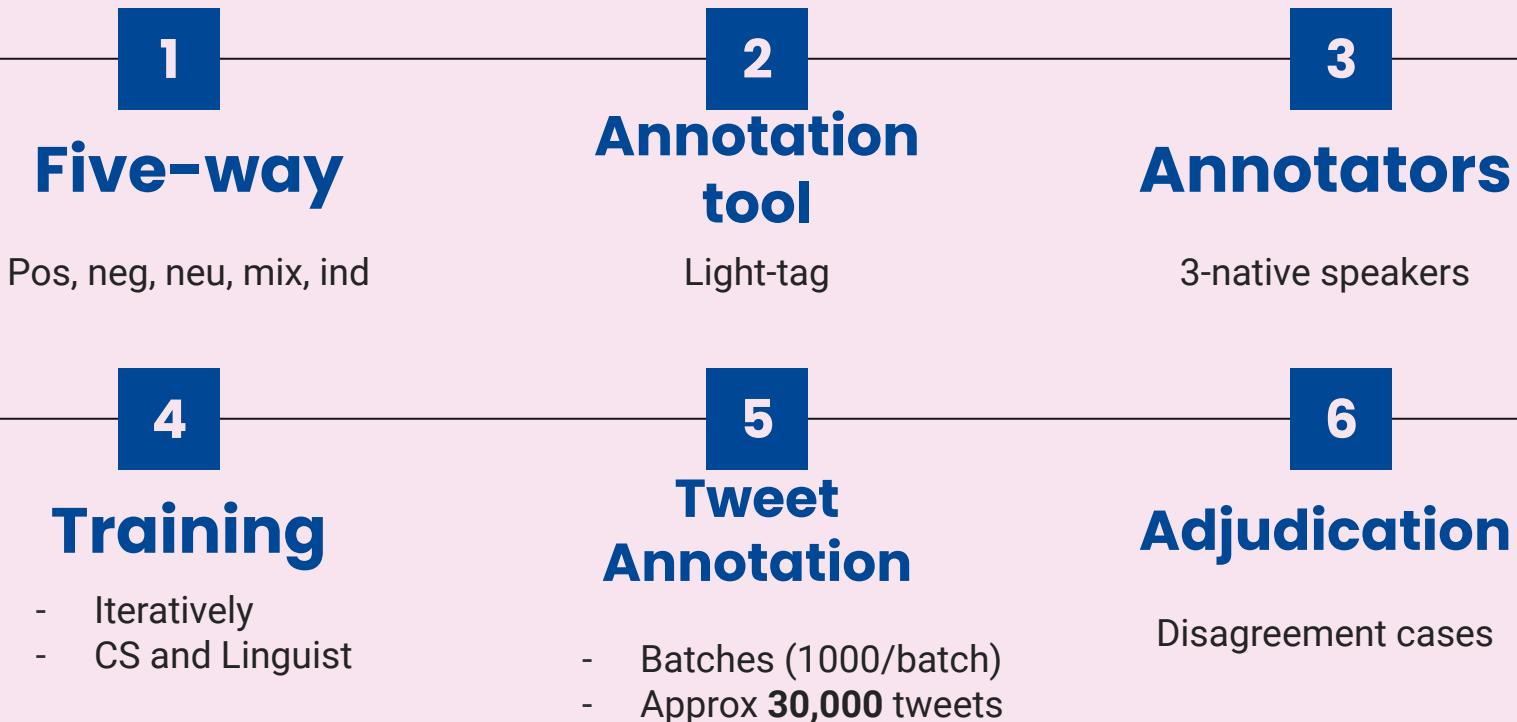
- Stopword + Emoji
- Translated sent.lexicon
- Hashtag e.g #Yorubaday
- Handles e.g @bbchausa

## Language Detection



- Stopword overlap
- "nke", produce tweet in Hausa *amin ya rabbi godiya nke*"
- longitude, latitude, and radius parameters (25 miles)

# Annotation Process



# Inter-annotators progress over Thirty batches

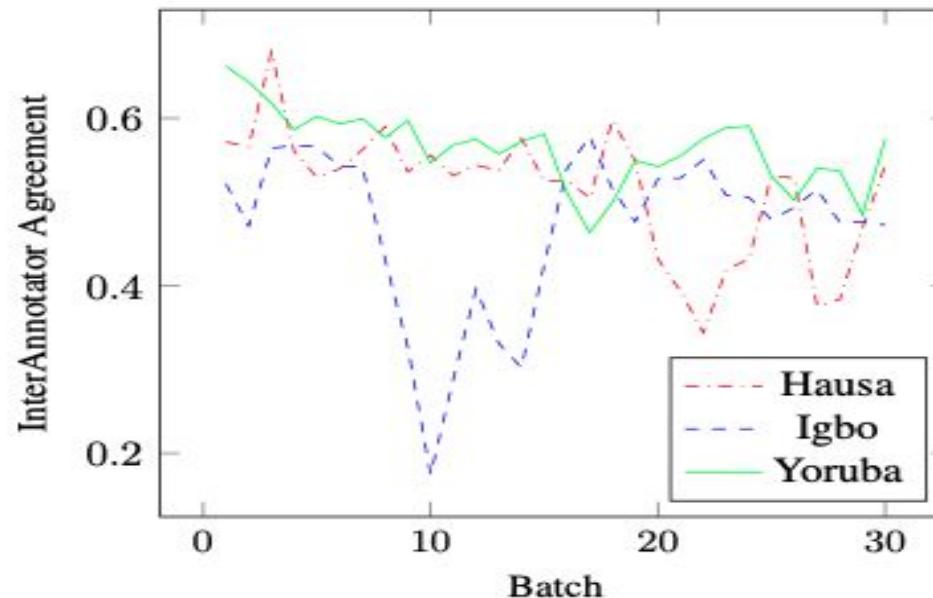
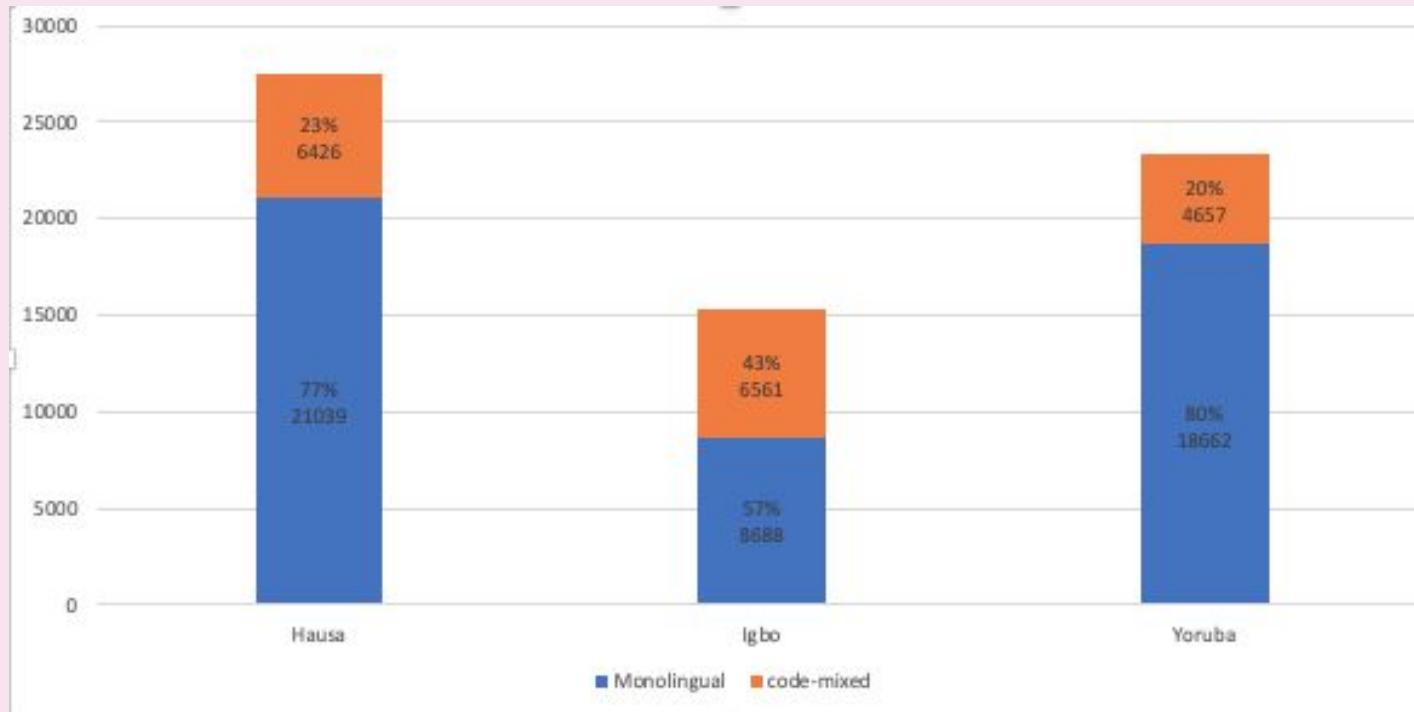


Figure 1: Inter-annotators progress over thirty batches—one-thousand tweets per batch.

# Monolingual vs Code-mix



# Diacritics challenge

E.g., Yoruba

- “Awon omo fo abo” does not have a meaning without diacritics,
- annotators classify it as indeterminate

8

Hausa

1904

Igbo

1754

Yoruba

# Tonality

- Tone in Yorùbá and Igbo helps to give meaning to words in context,
- The same sentence but different tones may have opposite sentiment:
  - Àwon omó fo abó (The children washed the dishes) has a **positive** meaning,
  - Àwon omó fó abó ..(The children broke the dishes) has a **negative** meaning

# Tonality

- Similarly, tonality is heavily used in Igbo.
  - ò nwèkwàrà mgbe i naenwe sense ? – Will you ever be able to talk sensibly? – You're a fool.
  - • ò nwèkwàrà mgbe i naenwe sense – Sometimes you act with great maturity – I'm impressed.

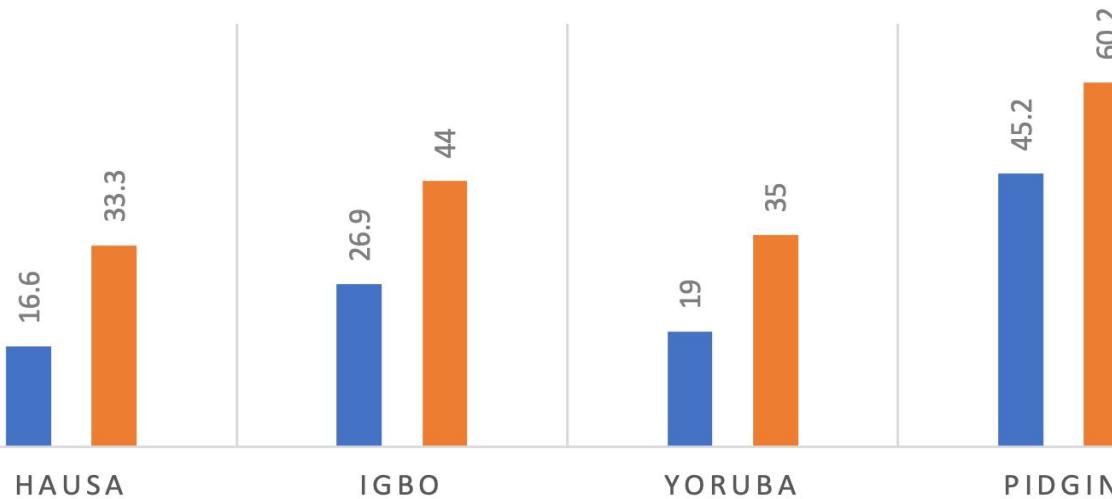
04



# Benchmark Experiments

# MAJORITY CLASSIFIER

■ Weighted-F1 ■ Micro-F1



# Multilingual PLMs

AfriBERTa  
winner

Model	NG lang. supported	PLM size	hau	ibo	pcm	yor	Avg
Multilingual PLMs							
AfriBERTa-large	hau, ibo, pcm, yor	126M	$81.0_{\pm 0.2}$	<b><math>81.2_{\pm 0.5}</math></b>	$75.0_{\pm 0.6}$	$80.2_{\pm 0.6}$	$79.3_{\pm 0.3}$
mBERT-base	yor	172M	$77.8_{\pm 0.5}$	$79.8_{\pm 0.5}$	$72.4_{\pm 1.5}$	$77.6_{\pm 0.9}$	$76.9_{\pm 0.3}$
XLM-R-base	hau	270M	$78.4_{\pm 1.0}$	$79.9_{\pm 0.7}$	$76.3_{\pm 0.6}$	$76.9_{\pm 0.4}$	$77.9_{\pm 0.2}$
mDeBERTaV3-base	None	276M	$79.3_{\pm 0}$	$80.7_{\pm 0.2}$	$77.6_{\pm 0.8}^*$	$78.4_{\pm 0.5}$	$79.0_{\pm 0.3}$
RemBERT	hau, ibo, yor	559M	$79.0_{\pm 0}$	$79.9_{\pm 0.4}$	<b><math>78.4_{\pm 1.4}^*</math></b>	$78.0_{\pm 0.6}$	$78.8_{\pm 0.6}$

Better than  
others

# Language adaptive fine-tuning (LAFT)

AfriBERTa  
Competitive with  
LAFT

Model	NG lang. supported	PLM size	hau	ibo	pcm	yor	Avg
Multilingual PLMs+LAFT							
mBERT+LAFT (General)	hau / ibo / pcm / yor	172M	$80.8 \pm 0.3$	$80.4 \pm 0.4$	$74.2 \pm 0.5$	$80.8 \pm 0.5$	$79.1 \pm 0.3$
mBERT+LAFT (Tweet)	hau / ibo / pcm / yor	172M	$79.3 \pm 0.6$	$77.7 \pm 0.6$	$74.0 \pm 0.7$	$76.8 \pm 0.3$	$77.0 \pm 0.3$
XLM-R-base+LAFT (General)	hau / ibo / pcm / yor	270M	$81.5^* \pm 0.7$	$80.8^* \pm 0.8$	$74.7 \pm 1.5$	$80.9^* \pm 0.4$	$79.5^* \pm 0.3$
XLM-R-base+LAFT (Tweet)	hau / ibo / pcm / yor	270M	$79.5 \pm 0.9$	$77.0 \pm 0.5$	$74.8 \pm 0.7$	$76.2 \pm 0.4$	$76.9 \pm 0.2$

- CC100 > 318MB general purpose
- Hausa tweets =32MB

# Sentiment classification model

1st step



**Multilingual PLMs**  
( mBert, XLM-R,  
AfriBERTa, mDeBER-  
TaV3, Remberta)

2nd step



**Language  
adaptive  
fine-tuning (LAFT)**

3rd step



**Multi-task  
sentiment  
classification**

4th step



**Zero-shot  
Cross-Lingual  
Transfer**

# Multi-task Sentiment Classification

Model	NG lang. supported	PLM size	hau	ibo	pcm	yor	Avg
Multi-task Multilingual PLMs							
AfriBERTa-large	hau, ibo, pcm, yor	126M	81.2 $\pm$ 0.2	80.8 $\pm$ 0.2	74.5 $\pm$ 0.6	80.4 $\pm$ 0.7	79.3 $\pm$ 0.3
mDeBERTaV3-base	None	276M	79.0 $\pm$ 0.2	79.3 $\pm$ 0.5	76.8 $\pm$ 0.6	78.7 $\pm$ 0.4	78.4 $\pm$ 0.3

**AfriBERTa  
winner**

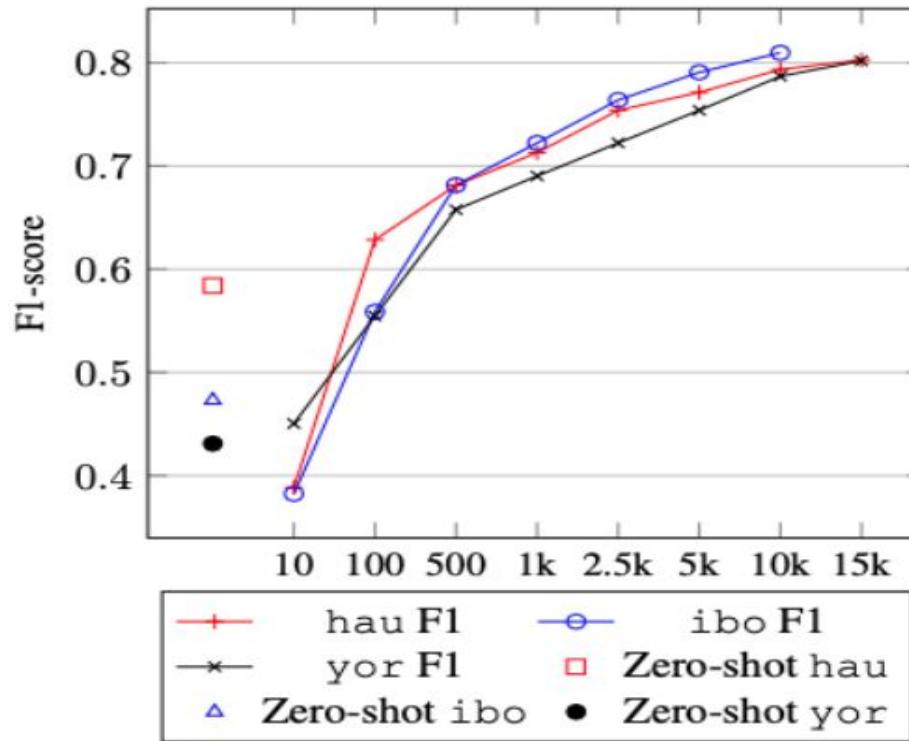
**mDeBERTaV3-base**  
sight drop (-0.6%)

# Zero-shot Cross-Lingual Transfer

Model	hau	ibo	pcm	yor	Avg
AfriBERTa-large	<b>58.4</b>	<b>47.7</b>	62.0	<b>43.1</b>	<b>52.8</b>
mBERT-base	31.0	37.0	57.4	39.5	41.2
XLM-R-base	38.4	37.8	62.4	26.7	41.3
mDeBERTaV3-base	50.1	47.2	64.9	36.4	49.7
RemBERT	54.0	45.4	<b>66.2</b>	30.2	49.0

Table 8: Transfer Learning experiments. PLMs are trained on English SemEval 2017 and evaluated on NG languages in a zero-shot setting

# Sample Efficiency in Transfer



# Conclusions and Future Work

- We present NaijaSenti
- Sentiment Lexicon
- Benchmark experiments
  - AfriBERTa
- NaijaSenti has the potential to spark interest in sentiment analysis and other downstream NLP tasks in the languages involved.
- Extend the corpus : AfriSenti

# Thanks!

Paper: <https://arxiv.org/abs/2103.11811>

Code: <https://github.com/hausanlp/NaijaSenti>