

ACII 2023
The 11th International Conference on Affective Computing and Intelligent Interaction

BARLETTE MAKUWE, KOENA RONNY MABOKELA, TIM SCHLIPPE

SENTIMENT ANALYSIS FOR SHONA

Cambridge, USA

September 13, 2023

AGENDA

Introduction

1

Related Work

2

ShonaSenti – The Shona Sentiment Analysis Corpus

3

Sentiment Analysis for Shona

4

Conclusion and Future Work

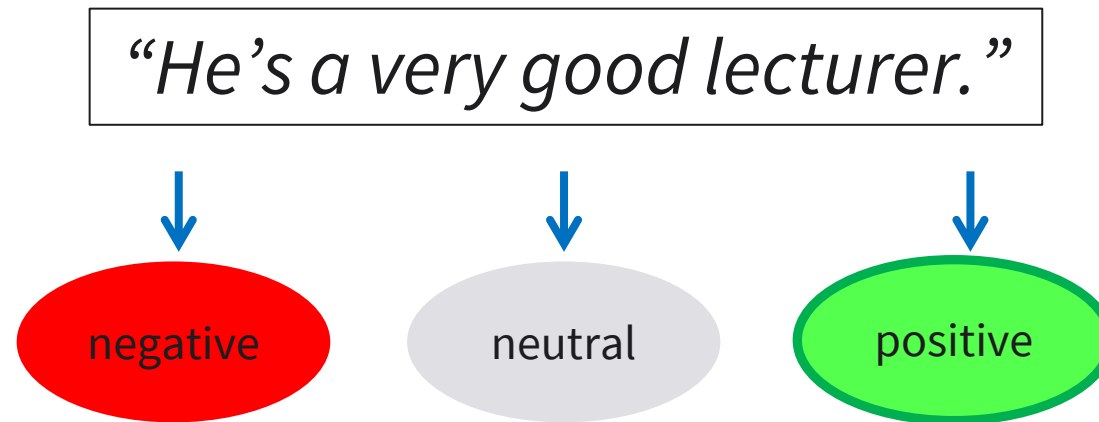
5

1

INTRODUCTION

SENTIMENT ANALYSIS

- Detecting & classifying sentiments, e.g.,



- Applications in business, medicine, education, AI for social good
(Rokade & Kumari, 2019; Zucco et al., 2018; Lalata et al., 2019; Rakhmanov & Schlippe, 2022; Mabokela & Schlippe, 2022)

SHONA

- African language, spoken by 16.6 million people (Ethnologue, 2023)
- Mother tongue of 75% of Zimbabwe's population (African Studies Center, 2023)
- Also spoken in Zambia, Botswana, Mozambique, South Africa, UK (Wikipedia, 2023)

➔ But: No sentiment analysis system

2

RELATED WORK

RELATED WORK

- Corpora usually retrieved from microblogging services like Twitter
(Go et al., 2009; Indriani et al., 2020)

RELATED WORK

- Corpora usually retrieved from microblogging services like Twitter
(Go et al., 2009; Indriani et al., 2020)
- Research on data collection/annotation methods
(Go et al., 2009; Pak & Paroubek, 2010; Agarwal & Sabharwal, 2012; Nakov et al., 2016; Indriani et al., 2020)

RELATED WORK

- Corpora usually retrieved from microblogging services like Twitter
(Go et al., 2009; Indriani et al., 2020)
- Research on data collection/annotation methods
(Go et al., 2009; Pak & Paroubek, 2010; Agarwal & Sabharwal, 2012; Nakov et al., 2016; Indriani et al., 2020)
- Distant supervised methods with emoticons and keywords
(Mabokela & Schlippe, 2022)

- Corpora usually retrieved from microblogging services like Twitter
(Go et al., 2009; Indriani et al., 2020)
- Research on data collection/annotation methods
(Go et al., 2009; Pak & Paroubek, 2010; Agarwal & Sabharwal, 2012; Nakov et al., 2016; Indriani et al., 2020)
- Distant supervised methods with emoticons and keywords
(Mabokela & Schlippe, 2022)
- For automatic sentiment analysis, Transformers usually outperform traditional and lexicon-based approaches
(Rakhmanov & Schlippe, 2022)

RELATED WORK

- AfroLM: Only Transformer-based pre-trained model to deal with Shona
- (But pre-trained on only 32.8 MB of Shona text) (Dossou et al., 2022)

RELATED WORK

- AfroLM: Only Transformer-based pre-trained model to deal with Shona
- (But pre-trained on only 32.8 MB of Shona text) (Dossou et al., 2022)

- Cross-lingual approaches to solve the problems of low resources
(Balahur & Turchi, 2014; Lin et al., 2014; Vilares et al., 2017; Can et al., 2018)

RELATED WORK

- AfroLM: Only Transformer-based pre-trained model to deal with Shona
 - (But pre-trained on only 32.8 MB of Shona text) (Dossou et al., 2022)

 - Cross-lingual approaches to solve the problems of low resources
(Balahur & Turchi, 2014; Lin et al., 2014; Vilares et al., 2017; Can et al., 2018)
-
- ➔ **Twitter** for corpus collection
 - ➔ **Distant supervised methods** to accelerate the labeling process
 - ➔ **AfroLM** for monolingual and **English BERT** for cross-lingual systems

3

SHONASENTI

COLLECTION OF THE TWEETS

- Essential Twitter API package

COLLECTION OF THE TWEETS

- Essential Twitter API package
- Specific collection with:
 - Topic-specific search terms to cover different topics
 - Geocode of Zimbabwe and radius to reduce other languages

COLLECTION OF THE TWEETS

- Essential Twitter API package
- Specific collection with:
 - Topic-specific search terms to cover different topics
 - Geocode of Zimbabwe and radius to reduce other languages
- Then: Frequency-based language identification

COLLECTION OF THE TWEETS

- Essential Twitter API package
- Specific collection with:
 - Topic-specific search terms to cover different topics
 - Geocode of Zimbabwe and radius to reduce other languages
- Then: Frequency-based language identification



Language	Frequency	Percentage
Shona	15,960	75.9%
English	2,776	13.2%
code-switch	2,184	10.4%
others	15,315	0.5%

COLLECTION OF THE TWEETS

- Essential Twitter API package
- Specific collection with:
 - Topic-specific search terms to cover different topics
 - Geocode of Zimbabwe and radius to reduce other languages
- Then: Frequency-based language identification



Language	Frequency	Percentage
Shona	15,960	75.9%
English	2,776	13.2%
code-switch	2,184	10.4%
others	15,315	0.5%

- Goal: Accelerate labeling process in a semi-automated manner
- Analyzed 3 distant supervision approaches:
 - *Distant Supervision with Sentiment Lexicon*
 - *Distant Supervision with Emoticons*
 - *Distant Supervision with Sentiment Lexicon + Emoticons*

LABELING STRATEGIES

- Goal: Accelerate labeling process in a semi-automated manner

- Analyzed 3 distant supervision approaches:

- *Distant Supervision with Sentiment Lexicon* → 21.0%

- *Distant Supervision with Emoticons*

- *Distant Supervision with Sentiment Lexicon + Emoticons*

correctly pre-labeled

LABELING STRATEGIES

- Goal: Accelerate labeling process in a semi-automated manner

- Analyzed 3 distant supervision approaches:

- *Distant Supervision with Sentiment Lexicon* → 21.0%

- *Distant Supervision with Emoticons* → 23.1%

- *Distant Supervision with Sentiment Lexicon + Emoticons*

correctly pre-labeled

LABELING STRATEGIES

- Goal: Accelerate labeling process in a semi-automated manner
- Analyzed 3 distant supervision approaches:
 - *Distant Supervision with Sentiment Lexicon* → 21.0%
 - *Distant Supervision with Emoticons* → 23.1%
 - *Distant Supervision with Sentiment Lexicon + Emoticons* → **23.6%**
correctly pre-labeled

LABELING STRATEGIES

- Goal: Accelerate labeling process in a semi-automated manner

- Analyzed 3 distant supervision approaches:

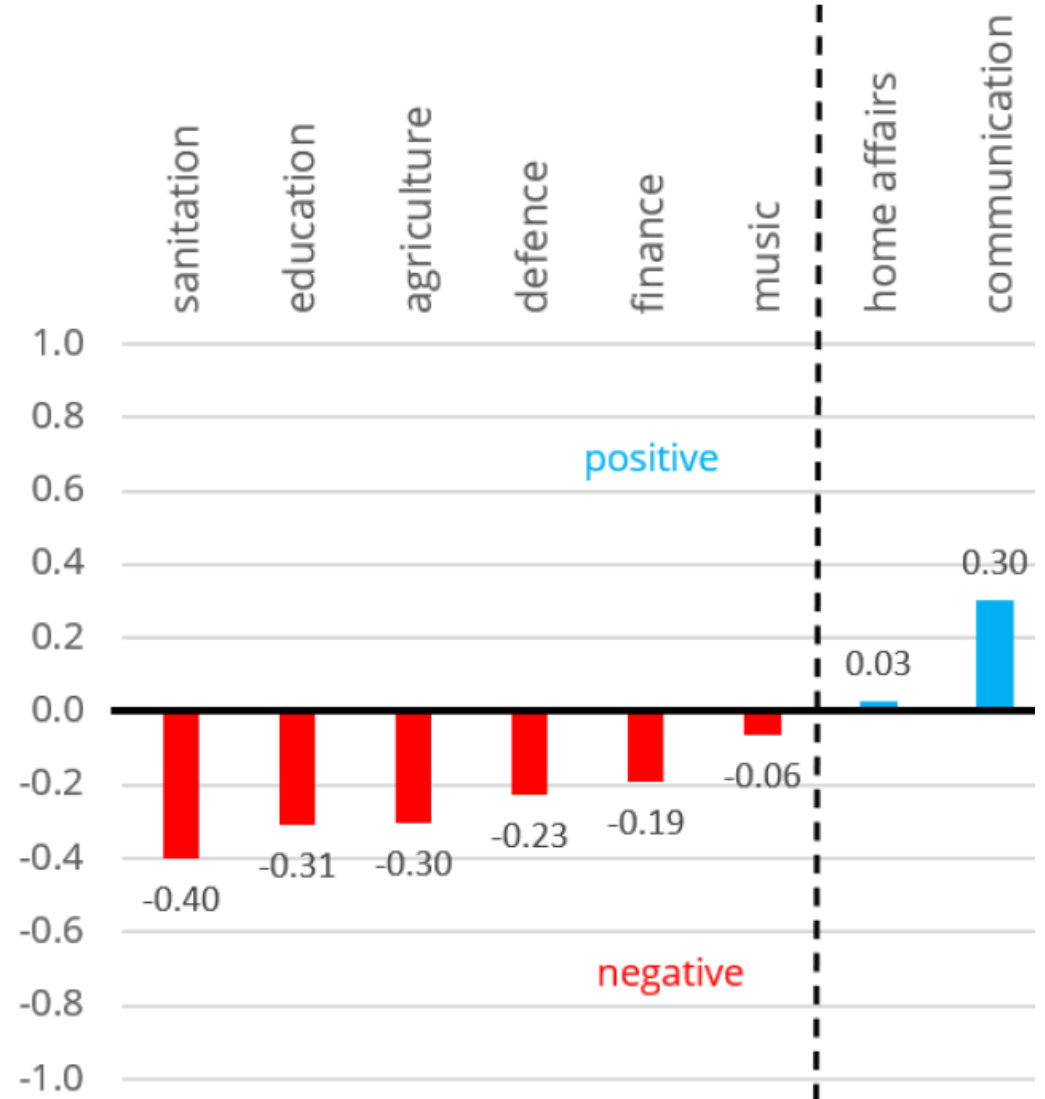
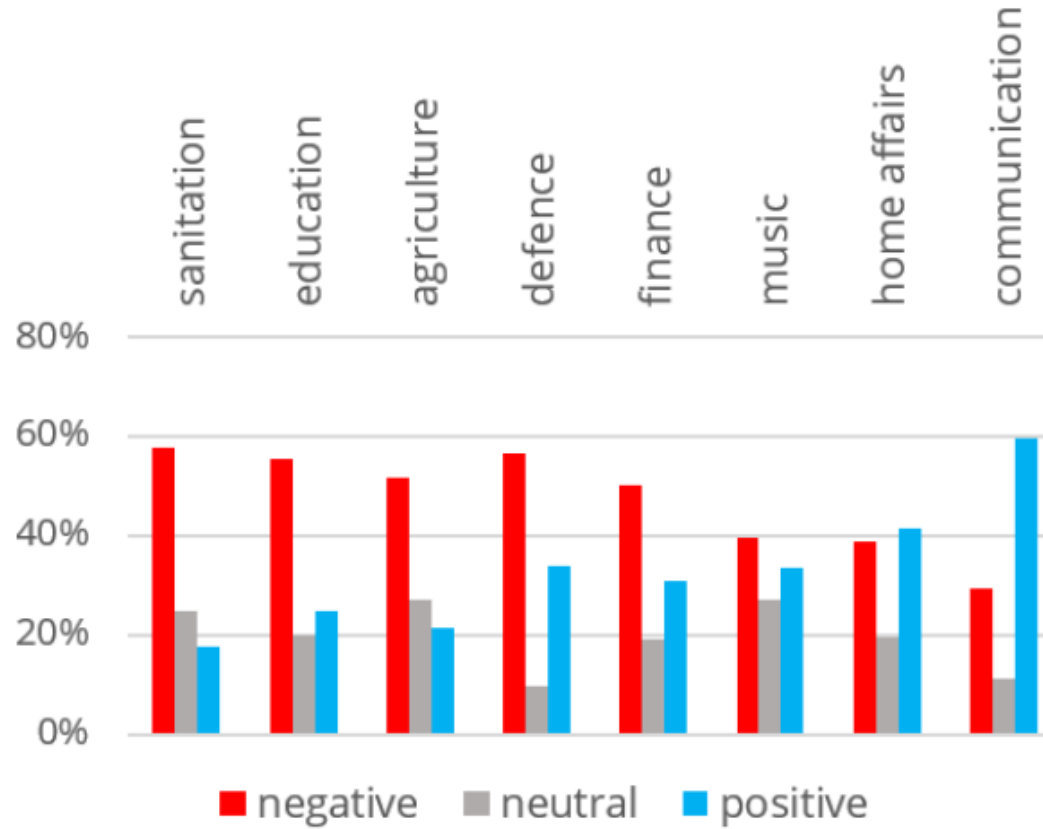
- *Distant Supervision with Sentiment Lexicon* → 21.0%
- *Distant Supervision with Emoticons* → 23.1%
- *Distant Supervision with Sentiment Lexicon + Emoticons* → **23.6%**

correctly pre-labeled

→ cross-checked, corrected, and added missing labels

TOPICS

$$\frac{\#negative * (-1) + 0 * \#neutral * (+1) * \#positive}{\#allsentiments}$$





“I have nothing much to say but he’s a very good lecturer.”



- (1) corrected
- (2) cross-checked
- (3) extended to 5 labels



“I have nothing much to say but he’s a very good lecturer.”

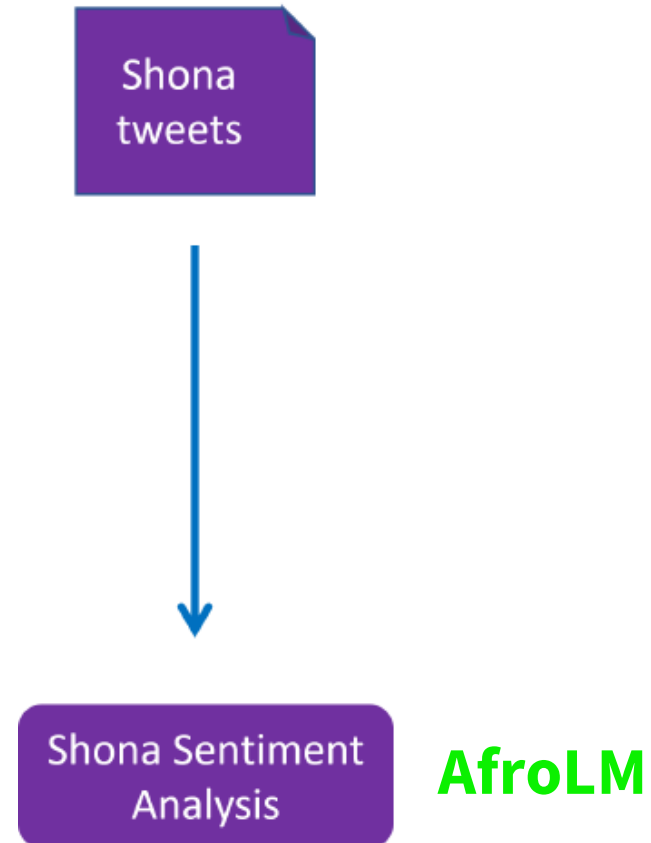
	#sentences	Percentage
Train	12,767	80%
Test	3,193	20%

4

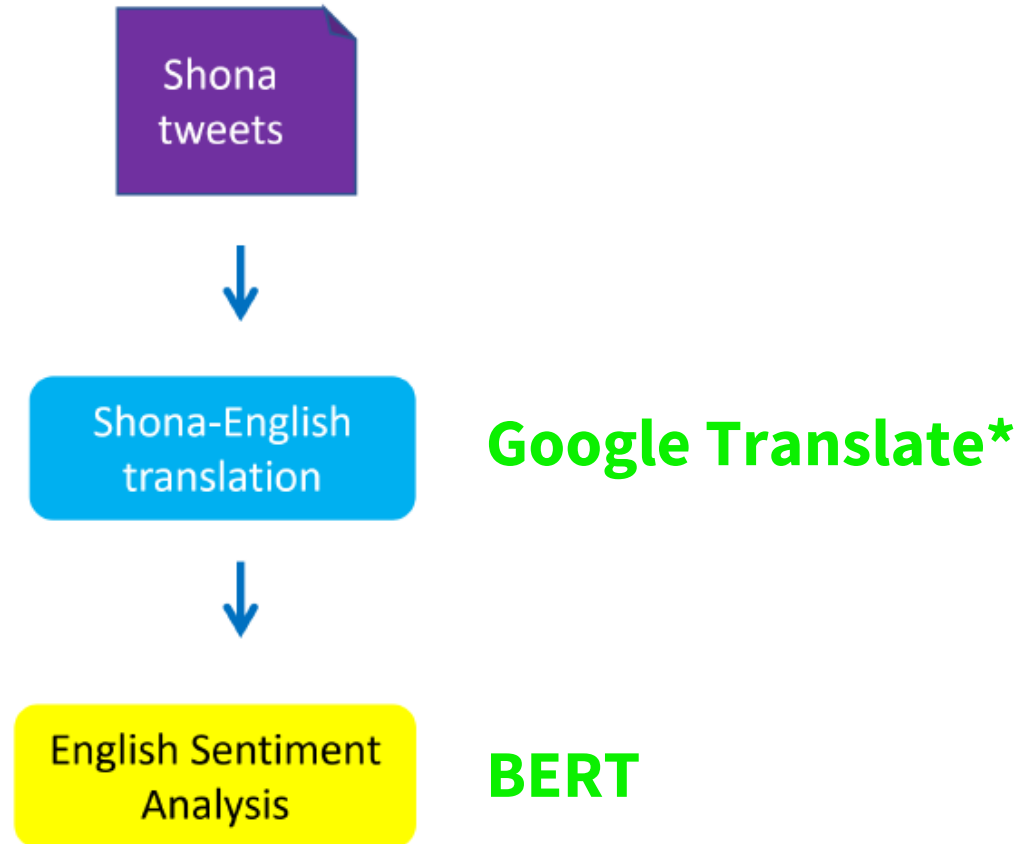
SENTIMENT ANALYSIS FOR SHONA

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

MONO-LINGUAL SENTIMENT ANALYSIS



CROSS-LINGUAL SENTIMENT ANALYSIS



*BLEU = 0.57 on ShonaSenti subset

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

System	Emoticons	Accuracy	F-score
<i>monoLingualNoEmoticons3classes_{AfroLM}</i>	no	78.60%	76.39%
<i>monoLingualEmoticons3classes_{AfroLM}</i>	yes	80.56%	79.44%
<i>crossLingualNoEmoticons3classes_{BERT}</i>	no	77.63%	76.56%
<i>crossLingualEmoticons3classes_{BERT}</i>	yes	83.88%	82.34%
<i>monoLingualNoEmoticons5classes_{AfroLM}</i>	no	68.00%	67.00%
<i>monoLingualEmoticons5classes_{AfroLM}</i>	yes	70.38%	69.56%
<i>crossLingualNoEmoticons5classes_{BERT}</i>	no	60.02%	59.23%
<i>crossLingualEmoticons5classes_{BERT}</i>	yes	61.56%	60.03%

➔ Impact of emoticons in the tweets on system performance

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

System	Emoticons	Accuracy	F-score
<i>monoLingualNoEmoticons3classesAfroLM</i>	no	78.60%	76.39%
<i>monoLingualEmoticons3classesAfroLM</i>	yes	80.56%	79.44%
<i>crossLingualNoEmoticons3classesBERT</i>	no	77.63%	76.56%
<i>crossLingualEmoticons3classesBERT</i>	yes	83.88%	82.34%
<i>monoLingualNoEmoticons5classesAfroLM</i>	no	68.00%	67.00%
<i>monoLingualEmoticons5classesAfroLM</i>	yes	70.38%	69.56%
<i>crossLingualNoEmoticons5classesBERT</i>	no	60.02%	59.23%
<i>crossLingualEmoticons5classesBERT</i>	yes	61.56%	60.03%

➔ Impact of emoticons in the tweets on system performance

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

System	Emoticons	Accuracy	F-score
<i>monoLingualNoEmoticons3classes_{AfroLM}</i>	no	78.60%	76.39%
<i>monoLingualEmoticons3classes_{AfroLM}</i>	yes	80.56%	79.44%
<i>crossLingualNoEmoticons3classes_{BERT}</i>	no	77.63%	76.56%
<i>crossLingualEmoticons3classes_{BERT}</i>	yes	83.88%	82.34%
<i>monoLingualNoEmoticons5classes_{AfroLM}</i>	no	68.00%	67.00%
<i>monoLingualEmoticons5classes_{AfroLM}</i>	yes	70.38%	69.56%
<i>crossLingualNoEmoticons5classes_{BERT}</i>	no	60.02%	59.23%
<i>crossLingualEmoticons5classes_{BERT}</i>	yes	61.56%	60.03%

➔ Impact of emoticons in the tweets on system performance

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

System	Emoticons	Accuracy	F-score
<i>monoLingualNoEmoticons3classesAfroLM</i>	no	78.60%	76.39%
<i>monoLingualEmoticons3classesAfroLM</i>	yes	80.56%	79.44%
<i>crossLingualNoEmoticons3classesBERT</i>	no	77.63%	76.56%
<i>crossLingualEmoticons3classesBERT</i>	yes	83.88%	82.34%
<i>monoLingualNoEmoticons5classesAfroLM</i>	no	68.00%	67.00%
<i>monoLingualEmoticons5classesAfroLM</i>	yes	70.38%	69.56%
<i>crossLingualNoEmoticons5classesBERT</i>	no	60.02%	59.23%
<i>crossLingualEmoticons5classesBERT</i>	yes	61.56%	60.03%

➔ Impact of emoticons in the tweets on system performance

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

System	Emoticons	Accuracy	F-score
<i>monoLingualNoEmoticons3classesAfroLM</i>	no	78.60%	76.39%
<i>monoLingualEmoticons3classesAfroLM</i>	yes	80.56%	79.44%
<i>crossLingualNoEmoticons3classesBERT</i>	no	77.63%	76.56%
<i>crossLingualEmoticons3classesBERT</i>	yes	83.88%	82.34%
<i>monoLingualNoEmoticons5classesAfroLM</i>	no	68.00%	67.00%
<i>monoLingualEmoticons5classesAfroLM</i>	yes	70.38%	69.56%
<i>crossLingualNoEmoticons5classesBERT</i>	no	60.02%	59.23%
<i>crossLingualEmoticons5classesBERT</i>	yes	61.56%	60.03%

- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

System	Emoticons	Accuracy	F-score
<i>monoLingualNoEmoticons3classesAfroLM</i>	no	78.60%	76.39%
<i>monoLingualEmoticons3classesAfroLM</i>	yes	80.56%	79.44%
<i>crossLingualNoEmoticons3classesBERT</i>	no	77.63%	76.56%
<i>crossLingualEmoticons3classesBERT</i>	yes	83.88%	82.34%
<i>monoLingualNoEmoticons5classesAfroLM</i>	no	68.00%	67.00%
<i>monoLingualEmoticons5classesAfroLM</i>	yes	70.38%	69.56%
<i>crossLingualNoEmoticons5classesBERT</i>	no	60.02%	59.23%
<i>crossLingualEmoticons5classesBERT</i>	yes	61.56%	60.03%

- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes
- Mono-lingual system better for 5 classes

MONO-LINGUAL + CROSS-LINGUAL SENTIMENT ANALYSIS

System	Emoticons	Accuracy	F-score
<i>monoLingualNoEmoticons3classesAfroLM</i>	no	78.60%	76.39%
<i>monoLingualEmoticons3classesAfroLM</i>	yes	80.56%	79.44%
<i>crossLingualNoEmoticons3classesBERT</i>	no	77.63%	76.56%
<i>crossLingualEmoticons3classesBERT</i>	yes	83.88%	82.34%
<i>monoLingualNoEmoticons5classesAfroLM</i>	no	68.00%	67.00%
<i>monoLingualEmoticons5classesAfroLM</i>	yes	70.38%	69.56%
<i>crossLingualNoEmoticons5classesBERT</i>	no	60.02%	59.23%
<i>crossLingualEmoticons5classesBERT</i>	yes	61.56%	60.03%

- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes
- Mono-lingual system better for 5 classes
- Performances with 3 classes are higher than with 5 classes

5

CONCLUSION AND FUTURE WORK

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (***sentiment lexicon + emoticons***)

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (***sentiment lexicon + emoticons***)
- Investigated **mono-lingual and cross-lingual sentiment analysis approaches** for Shona

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (**sentiment lexicon + emoticons**)
- Investigated **mono-lingual and cross-lingual sentiment analysis approaches** for Shona
- Impact of emoticons in the tweets on system performance

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (**sentiment lexicon + emoticons**)
- Investigated **mono-lingual and cross-lingual sentiment analysis approaches** for Shona
- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes: 83.88% accuracy
- Mono-lingual system better for 5 classes: 70.38% accuracy

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (**sentiment lexicon + emoticons**)
- Investigated **mono-lingual and cross-lingual sentiment analysis approaches** for Shona
- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes: 83.88% accuracy
- Mono-lingual system better for 5 classes: 70.38% accuracy

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (**sentiment lexicon + emoticons**)
- Investigated **mono-lingual and cross-lingual sentiment analysis approaches** for Shona
- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes: 83.88% accuracy
- Mono-lingual system better for 5 classes: 70.38% accuracy

Future Work

- **Optimize data annotation process** with the help of machine learning

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (**sentiment lexicon + emoticons**)
- Investigated **mono-lingual and cross-lingual sentiment analysis approaches** for Shona
- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes: 83.88% accuracy
- Mono-lingual system better for 5 classes: 70.38% accuracy

Future Work

- **Optimize data annotation process** with the help of machine learning
- **System combination**

CONCLUSION AND FUTURE WORK

Conclusion

- Shona is spoken by 16.6 million people
- Corpus collection of 16k Tweets covering 8 topics
 - the Shona Sentiment Analysis Corpus (**ShonaSenti**)
- 23.6% less labeling effort with distant supervision approach (**sentiment lexicon + emoticons**)
- Investigated **mono-lingual and cross-lingual sentiment analysis approaches** for Shona
- Impact of emoticons in the tweets on system performance
- Cross-lingual system better for 3 classes: 83.88% accuracy
- Mono-lingual system better for 5 classes: 70.38% accuracy

Future Work

- **Optimize data annotation process** with the help of machine learning
- **System combination**
- expand **SAfriSenti** with more African low-resourced languages

THANK YOU

Tim Schlippe

 tim.schlippe@iu.org

Literature

- **Agarwal, A. & Sabharwal, J. S. (2012):** *End-to-End Sentiment Analysis of Twitter Data*, in Workshop on Information Extraction and Entity Analytics on Social Media Data.
- **Balahur, A. & Turchi, M. (2014):** *Comparative Experiments using Supervised Learning and Machine Translation for Multilingual Sentiment Analysis*, *Comput. Speech Lang.*, vol. 28, pp. 56–75.
- **Can, E. F. , Ezen-Can, A. , & Can, F. (2018):** *Multilingual Sentiment Analysis: An RNN-Based Framework for Limited Data*, in ACM SIGIR 2018 Workshop on Learning from Limited or Noisy Data.
- **Dossou, B. F. P., Tonja, A. L. , Yousuf, O., Osei, S., Oppong, A., Shode, I., Awoyomi, O. O., & Emezue, C. (2022):** *AfroLM: A Self-Active Learning-Based Multilingual Pretrained Language Model for 23 African Languages*, in Proceedings of The Third Workshop on Simple and Efficient Natural Language Processing (SustainLP). Abu Dhabi, United Arab Emirates (Hybrid): Association for Computational Linguistics, pp. 52–64.
- **Ethnologue (2023):** *Shona*, <https://www.ethnologue.com/language/sna>, accessed: 01-2023.
- **A. Go, R. Bhayani, and L. Huang (2009):** *Twitter Sentiment Classification using Distant Supervision*, *Processing*, vol. 150.
- **D. Indriani, A. H. Nasution, W. Monika, and S. Nasution (2020):** *Towards a Sentiment Analyzer for Low-Resource Languages*, *CoRR*, vol. abs/2011.06382,

Literature

- **J. Lalata, B. Gerardo, and R. Medina (2019):** *A Sentiment Analysis Model for Faculty Comment Evaluation Using Ensemble Machine Learning Algorithms*, in The 2019 International Conference on Big Data Engineering, ser. BDE 2019. New York, NY, USA: ACM, p. 68–73.
- **Z. Lin, X. Jin, X. Xu, Y. Wang, S. Tan, and X. Cheng (2014):** *Make It Possible: Multilingual Sentiment Analysis Without Much Prior Knowledge*, in IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT), vol. 2, pp. 79–86.
- **K. R. Mabokela and T. Schlippe (2022):** *A Sentiment Corpus for South African Under-Resourced Languages in a Multilingual Context*, in The 1st Annual Meeting of the ELRA/ISCA Special Interest Group on Under-Resourced Languages (SIGUL 2022), p. 70–77.
- **K. R. Mabokela and T. Schlippe (2022):** *AI for Social Good: Sentiment Analysis to Detect Social Challenges in South Africa*, in The South African Conference for Artificial Intelligence Research (SACAIR 2022).
- **P. Nakov, A. Ritter, S. Rosenthal, F. Sebastiani, and V. Stoyanov (2016):** *SemEval-2016 Task 4: Sentiment Analysis in Twitter*, in International Workshop on Semantic Evaluation (SemEval).
- **A. Pak and P. Paroubek (2010):** *Twitter as a Corpus for Sentiment Analysis and Opinion Mining*, in The 7th Edition of the Language Resources and Evaluation Conference (LREC 2010), pp. 1320–1326.
- **O. Rakhmanov and T. Schlippe (2022):** *Sentiment Analysis for Hausa: Classifying Students' Comments*, in The 1st Annual Meeting of the ELRA/ISCA Special Interest Group on Under-Resourced Languages (SIGUL 2022), Marseille, France.

Literature

- **P. P. Rokade and A. Kumari D. (2019):** *Business Intelligence Analytics using Sentiment Analysis—A Survey*, International Journal of Electrical and Computer Engineering (IJECE).
- **University of Pennsylvania, African Studies Center (2023):** *Shona*, accessed: 01-2023. Available: <https://plc.sas.upenn.edu/shona>
- **D. Vilares, M. A. Alonso, and C. Gómez-Rodríguez (2015):** *Sentiment Analysis on Monolingual, Multilingual and Code-switching Twitter Corpora*, in The 6th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, pp. 2–8.
- **Wikipedia (2023):** *Shona People*, accessed: 01-2023. Available: https://en.wikipedia.org/wiki/Shona_people
- **C. Zucco, H. Liang, G. D. Fatta, and M. Cannataro (2018):** *Explainable Sentiment Analysis with Applications in Medicine*, in 2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pp. 1740–1747.