

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier XX.XXXX/ACCESS.2022.DOI

# Multilingual Sentiment Analysis for Under-Resourced Languages: A Systematic Review of the Landscape

KOENA RONNY MABOKELA<sup>1,2</sup>, TURGAY CELIK<sup>2</sup> AND MPHO RABORIFE<sup>3</sup>

<sup>1</sup>University of Johannesburg, Applied Information Systems, Johannesburg, South Africa

<sup>2</sup>School of Electrical and Information Engineering, University of the Witwatersrand, Johannesburg, South Africa.

<sup>3</sup>University of Johannesburg, Institute for Intelligence Systems, Johannesburg, South Africa

Corresponding author: Koena Ronny Mabokela (e-mail: krmabokela@gmail.com)

“This work was supported by National Research Foundation (NRF) for the Black Academics Advancement Programme (BAAP) Grant: BAAP200225506825.”

**ABSTRACT** Sentiment analysis automatically evaluates people's opinions of products or services. It is an emerging research area with promising advancements in high-resource languages such as Indo-European languages (e.g. English). However, the same cannot be said for languages with limited resources. In this study, we evaluate multilingual sentiment analysis (MSA) techniques for under-resourced languages and the use of high-resourced languages to develop resources for MSA in low-resource languages, with the ultimate goal of identifying appropriate strategies for future MSA investigations. We report over 35 studies with different languages demonstrating an interest in developing MSA models for under-resourced languages in a multilingual context. Furthermore, we illustrate the drawbacks of each strategy used for the MSA task. Our focus is critically comparing MSA methods and employed datasets and identifying research gaps. Our comparative analysis study contributes to theoretical literature reviews with complete coverage of MSA studies from 2008 to date. Furthermore, we demonstrate how MSA studies have grown tremendously. Finally, because most studies propose MSA methods based on deep learning approaches, we offer a deep learning framework for MSA that does not rely on machine translation systems. According to the meta-analysis (PRISMA) protocol of this literature review, we found that, in general, just over 60% of the studies have used deep learning frameworks, which significantly improved the MSA performance. Therefore, deep learning methods are recommended for the development of MSA for under-resourced languages.

**INDEX TERMS** Multilingual, Sentiment Analysis, Code-switching, Deep Learning, Cross-lingual, Under-resourced languages, Systematic review

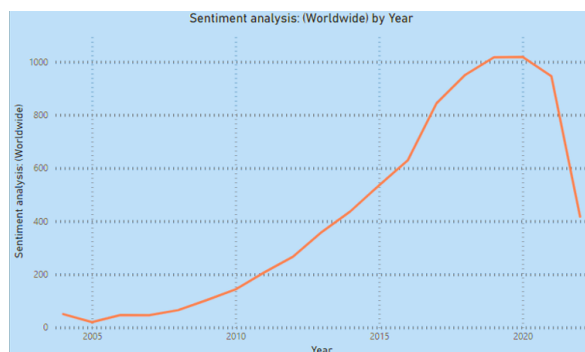
## I. INTRODUCTION

**S**ENTIMENT analysis is an intensive research activity in natural language processing (NLP). It uses NLP technologies to analyse textual messages and determine deeper contexts as they apply to a topic, brand, or theme [1], [2]. It is used to determine whether comments are subjective or objective and then classify such texts as *positive*, *negative* or *neutral* sentiment. sentiment analysis can be tackled at different classification levels such as document, sentence, aspect or feature [3], [4]. It has garnered considerable research attention, which can be attributed to its numerous essential NLP applications []. In recent years, sentiment analysis has gained even more interest owing to the rapid use of social media platforms. Its primary usage has been in businesses

and consumer care services [5], [6].

Social media users have made it a modern-day culture to use social media platforms to share their feelings or thoughts on various subjects [7]. As a result, social media platforms such as Facebook, Twitter, Instagram, WhatsApp and, most recently, TikTok, generate a large amount of potentially "rich" data [8]. Therefore, most sentiment analysis research studies have used social media data, particularly Twitter [9]–[11]. Notably, it is common for social media data to be multilingual and multicultural or have many linguistic variations, including mixed languages [12], [13].

Sentiment analysis is a growing research area with promising progress in high-resourced languages [14]–[16]. However, in another context, the same cannot be said for under-



**FIGURE 1.** Worldwide interest in sentiment analysis over time on Google Trends, 2004 - present.

resourced languages due to a lack of resources to develop NLP technologies. By under-resourced languages, we refer mainly to languages with little or no resources available to create digital language technologies [17]. Our study reviews multilingual sentiment analysis (MSA) methods for under-resourced languages. This paper uses the terms ‘*under-resourced languages*’ and ‘*low-resource languages*’.

Research on sentiment analysis has focused predominantly on single-language texts, mainly for high-resourced languages such as English [7], [18], [19]. Sentiment analysis research for high-resourced languages was actively studied due to the massive availability of resources such as benchmark datasets, annotated corpora and sentiment lexicons [13], [20]. In addition, sentiment analysis technologies developed for single-language tasks increase the risk of overlooking information in texts written in multiple languages [5], [21]. Deriu et al. [22] report that sentiment analysis methods developed for single-language texts could not be replicated for new or multilingual texts. Therefore, a concerted effort is necessary to create sentiment analysis models that cater for multiple languages. For example, some researchers proposed a cross-lingual sentiment analysis method with the help of a machine translation (MT) application and then applied machine learning (ML) techniques. The ML techniques such as support vector machine (SVM), Naive Bayes (NB), Maximum entropy (ME) are used for sentiments classification [23]–[25]. This cross-lingual sentiment analysis method has been successful in languages such as French, Spanish, Italian, Portuguese, Arabic, and Chinese [4], [12], [14].

MSA aims to recognise the sentiment of textual content written in multiple languages. It attempts to address issues presented by various languages, including code-switched comments. The success of monolingual sentiment analysis and MSA technologies mainly depends on the availability of labelled datasets to train computational language models [18], [26]. Recently, some resources and methods, including code-switched datasets, are available on SemEval, the largest workshop on computational semantic evaluations for multiple NLP research [27]. Although some of these MSA methods are already performing well for high-resourced

language datasets, they underperform for under-resourced languages, with English as the primary language and other contributing languages [8]. These sentiment analysis methods can only perform very well if labelled datasets are available or if methodologies that address the issues of under-resourced languages can be customised.

Most MSA approaches still rely on MT-based methods [23] or merging of monolingual datasets from different languages to build large-scale multilingual datasets [28]. Then apply ML techniques for sentiment classification. To some degree, MSA approaches that employ training of monolingual datasets from various languages cannot perform well for mixed-language texts. Some of the MSA methods are language-specific and may not be applied across distinct languages [21], [29]. In addition, supervised ML relies on a labelled dataset to produce accurate results [21]. Therefore, previous MSA research used manual data labelling methods, which is, to date the most labour-intensive and expensive process [21], [22], [29].

Code-switched texts originate from the most populous, multicultural societies and culturally diverse countries where more than one official language is spoken. Given this reality, social media users are more comfortable expressing their views in multiple languages [19], [22], [26]. Under-resourced languages are most commonly mixed with English [30], [31]. These mixed-language phenomena pose a significant challenge to existing MSA systems. In this context, we refer to multilingual data as sentences containing monolingual texts or code-switched data —texts written in more than one language. MSA for under-resourced languages is advancing gradually with progressive application of deep learning (DL) methods [15], [21], [32]. However, very few generic MSA methods have been developed for under-resourced languages [13]. In this study, we also examine MSA methods that considered mixed-language or code-switched texts intending to address under-resourced language challenges.

Prior studies on MSA explored the use of MT-based systems to transfer knowledge from resource-rich languages to under-resourced languages [17], [23]. These approaches translate text from an under-resourced language to English or vice-versa and then apply ML-based techniques to perform sentiment classification [14], [29]. Moreover, this method generally presents limitations like loss of meaning, and poor translation quality [12], [33]. In addition, [24] say that MT-based systems should be an obvious baseline system for any new MSA method [15]. In reality, the most recent development in the field of NLP has demonstrated that the effectiveness of MSA is significantly impacted by DL techniques [12], [20], [34]. Thus far, researchers have explored approaches such as convolutional neural networks (CNN), recurrent neural networks (RNN), adversarial neural networks (ANN), and generative adversarial networks (GAN) [20], [35].

The literature has made several attempts to address sentiment analysis in multilingual environments [12], [32] and others have addressed the problem from the perspective of creating methodologies that can actually operate on small

datasets [11], [36]. However, despite the fact that there is a multitude of studies on high-resource languages from the perspective of MSA [19], [37] but none of them focuses specifically on under-resourced languages in a multilingual setting. In this article, we concentrate on analysing the MSA literature survey from the perspective of under-resourced languages. Although there have been several literature surveys for MSA [13], [38]–[40], our study is by far the first to cover a mixture of high-resourced languages and under-resourced languages in a multilingual setting. We provide a detailed literature survey, along with the methods, models, mechanisms and performances, with a special focus on rule-based, cross-lingual, machine learning and deep learning techniques. The purpose is to the extent the research field provides scope for future research on under-resourced languages. We used the methods presented in Section III as a guideline to review the relevant studies. We categorised these studies as multilingual, cross-lingual, or code-switched approaches for under-resourced languages. Most significantly, we found that more than 40% of the research presented in this analysis had not been looked at in earlier literature reviews. The following contributions are made by our study:

- 1) To the best of our knowledge, we provide a detailed systematic review and an overview of MSA techniques for languages with limited resources in multilingual environments.
- 2) We provide the most recent comprehensive review of the MSA methods and an overview of MSA techniques for languages with limited resources in a multilingual environment.
- 3) We describe the outcomes of using cross-lingual sentiment analysis approaches to develop MSA methods for resource-constrained languages.
- 4) We further address the research questions raised in our systematic literature study.
- 5) Finally, we highlight the areas of research that need more investigation and offer suggestions for using MSA techniques in languages with limited resources in the future.

This paper is organised as follows: Section II presents the research questions. Section III presents the methodology used for this literature study. Section IV offers a summary of previous state-of-the-art studies. Section V describes MSA techniques and shortcomings. The evaluation metrics for MSA will be highlighted in section VI and the results will be discussed in section VII. Section VII will present the limitation of the study. In section IX, we will present the emerging MSA areas. Section X deals with emerging MSA areas. Lastly, we offer a conclusion and future suggestions.

## II. RESEARCH QUESTIONS

Our research aims to discover the most recent trends in MSA approaches for under-resourced languages. As a result, the following research questions guided this systematic literature review:

- 1) What are the existing MSA methods used to generate sentiment classification models, sentiment datasets and sentiment lexicons in a multilingual context?
- 2) What are MT system applications suitable for developing MSA methods and sentiment resources in multilingual environments?
- 3) What MSA techniques have been applied to sentiment classification for under-resourced languages in code-switched texts?
- 4) What are the DL and pre-trained techniques used to perform MSA for under-resourced languages?

Next, we will describe our research methods and limitations.

## III. RESEARCH METHODS

This section outlines the methodology utilised to achieve the purposes of this systematic review and provides a detailed description of the approaches and datasets used in MSA research. To prevent bias in our conclusion, we choose to undertake a thorough systematic literature study using high-quality peer-reviewed articles from 2008 to 2022 mainly because of the following reasons: (i) to bridge the gap between the research methods, which can be relevant and help address under-resourced language challenges, (ii) to provide a detailed overview of the most recent trends in MSA methods and offer an understanding of the research shift from 2008 to 2022, (iii) to recommend suitable research methods for future studies on under-resourced languages. Moreover, we clearly describe the differences across MSA methodologies and datasets. Furthermore, we examine how research has shifted from lexicon-based, cross-lingual methods and statistical ML techniques to more contemporary DL models, concentrating on low-resource languages and incorporating aspect-based sentiment analysis research. Lastly, a meta-analysis of the results from the selected articles is used to produce different summary tables.

### A. SEARCH STRATEGIES FOR RELEVANT STUDIES

The literature search was conducted by following the preferred reporting items for systematic reviews and meta-analyses (PRISMA) framework (i.e., Fig. 2). This framework includes the identification phase, screening phase, exclusion and inclusion phase [41]. Research communities widely use PRISMA for conducting a systematic literature review. The PRISMA framework was adopted in this study because it emphasizes the reporting of studies assessing the intervention's effect. It can be used as a foundation for publishing systematic reviews with goals other than considering interventions [41]. We began this literature review study by identifying the data sources and formulating the search keywords that eventually led to selecting the most relevant studies since the start of MSA research. Several peer-reviewed and published articles relating to sentiment analysis in multiple languages were used from 2008 to 2022. The methodology used in this study is depicted in the schematic diagram presented in Fig. 2.

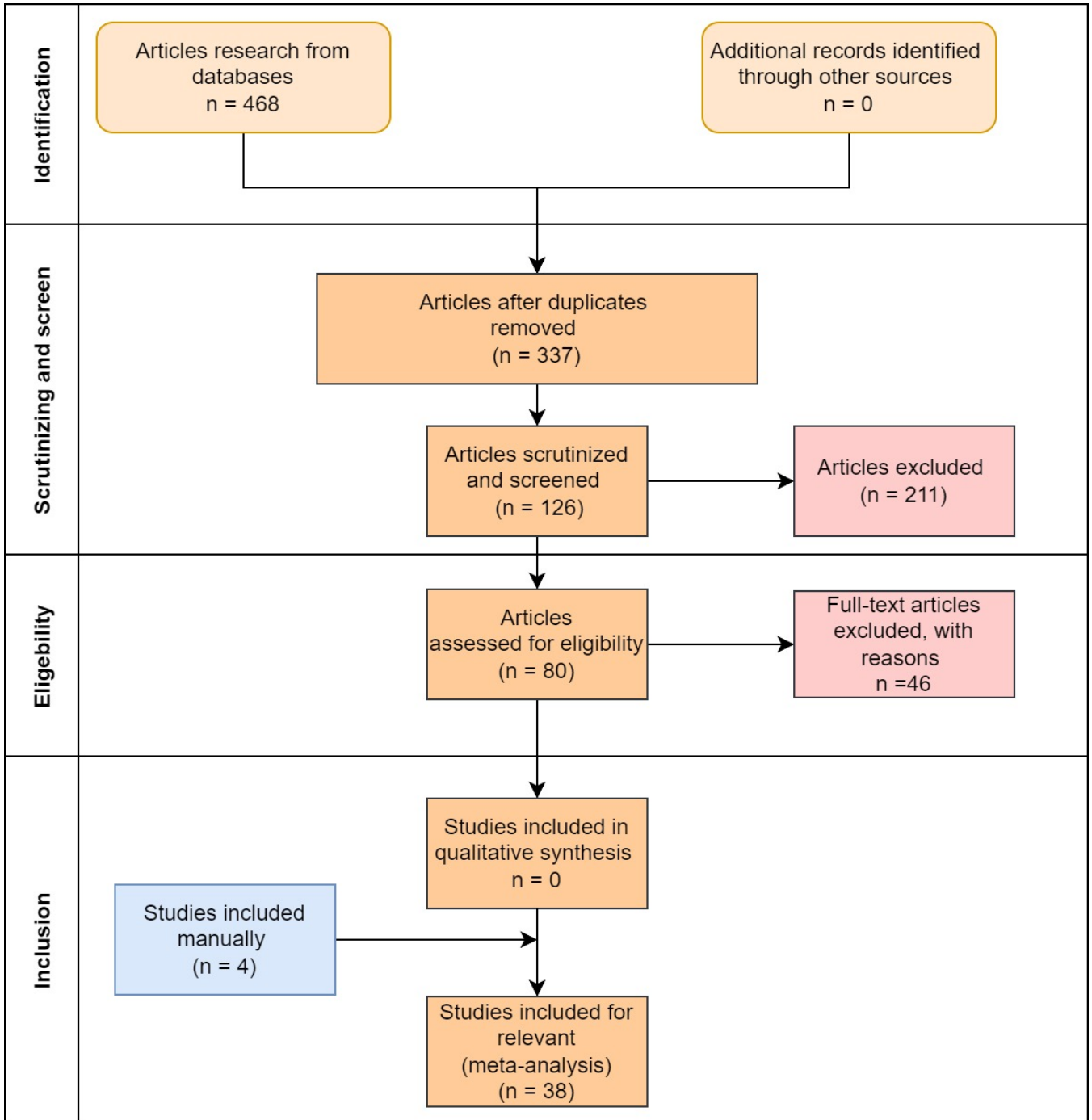


FIGURE 2. PRISMA process flowchart of the systematic literature survey methodology [42].

Note that the methodology used in this study is the same one used in the newly released comprehensive literature review on MSA for deep learning methods [13]. We used research papers produced from the search keywords as: "Multilingual Sentiment Analysis" OR "Cross-lingual AND sentiment analysis" OR "Multi-language AND NLP" OR "Multi-language AND sentiment" OR "Multi-lingual sentiment analysis" OR "Code mixed AND sentiment analysis"

OR "Code-switching AND sentiment analysis".

## B. INCLUSION AND EXCLUSION STRATEGIES

We provide a detailed explanation of our inclusion and exclusion criteria in this section.

**Inclusion Criteria.** In the next step, we selected articles based on their abstracts, methods, conclusions, and future directions. We included articles based on MSA studies for



under-resourced languages and studies on multilingual sentiment classification where the low-resource language is reported. We had relevant articles that investigated MSA from its infancy to date in the aspect of low-resourced languages. Peer-reviewed journals, book series and conference papers are selected because they are of high quality, and citation counts are also considered. We needed to read the entire article for some articles to determine whether they are to be included. We had a more detailed study version for studies published more than once. Finally, we selected conference proceedings papers that reported complete research during the same period.

**Exclusion Criteria.** There were few articles that are written in languages other than English. However, this systematic literature is based on peer-reviewed papers written exclusively in English. This approach neglects a requirement in a systematic literature review that discourages language limitations. Excluding these two articles from our analysis did not constitute a bias. Articles not published in computer science, decision science, mathematics and engineering journals and not using the techniques mentioned above were excluded, even if they were related to MSA studies for low-resource languages. Journal articles that did not present a complete or significant portion of their methodology are excluded. We also decided not to consider research studies that reported conceptual papers, work-in-progress, preliminary studies, or unfinished work. Review articles, survey papers and dissertations are excluded from this study.

### C. DIGITAL SOURCES AND DATA EXTRACTION

In the identification phase, research papers are searched from the following databases: Elsevier, Scopus, Google Scholar, Science Direct, IEEEExplore, ACM and Springer Link; the databases include studies done across the globe, and therefore geographical bias was not an issue. Furthermore, published books and book chapters are examined. We applied the filtering criteria to 468 papers. Our focus was on MSA for low-resource languages and the use of high-resourced languages to develop low-resource languages from a multilingual perspective. In the end, 38 primary studies were reviewed, and 4 of the 38 articles were later included manually in our research. The primary studies included models used to evaluate MSA systems, bilingual SA, cross-lingual sentiment analysis and code-switched SA, where knowledge/lexicon-based methods, ML-based and DL are reported to execute MSA tasks. Lastly, we derived the results and summary tables in the discussion section from the data extracted from the articles.

We screened and scrutinised the titles, abstracts, keywords, methods, conclusions, and citations and decided on potential eligibility. Studies were eligible if they reported on methods or models related to sentiment analysis in multiple languages. Studies of bilingual sentiment analysis are also considered. Studies with the same techniques used by other researchers were excluded, as well as the sentiment analysis methods developed for a single language. All data sources gathered

from social media platforms, and other related data sources are included in the datasets.

This research provides the most relevant and recent systematic literature survey for MSA in under-resourced languages. We aim to set the trend and suggest new approaches for under-resourced languages that can benefit sentiment analysis in a multilingual framework. While there are prior reported literature review studies, they paid more attention to sentiment analysis in a single language [5], [43]. Furthermore, several literature reviews that have been reported on the aspect of MSA [5], [13], [38]–[40], [44]. However, very few studies focused on MSA for under-resourced languages. For this, we briefly highlight the objectives of each literature survey and show the significant difference between the available literature review and our literature study. For example, the literature survey study presented by [13] only covered MSA methods for deep learning methods using social media data from 2017 to 2020. They only highlighted a shift of research from cross-lingual to code-switching MSA methods. Abdullah et al. [39] investigated a systematic literature review from 2010 to 2019 that covered the pre-processing methods, methods for sentiment analysis, the evaluation models utilised for MSA and the aspects of common languages supported in sentiment analysis.

Furthermore, Lo et al. [5] reviewed English-based sentiment analysis on social media as well as a few works on MSA for social media from 2010 to 2013. Santwana et al. [38] focused only on machine learning techniques for MSA in non-English languages from 2010 to 2018. However, our literature study covers even the most recent MSA methods employed from 2008 to 2022 whilst [40] only focused on the cross-lingual sentiment analysis methods for Chinese languages from 2004 to 2022. Lastly, Xu et al. [44] investigated a systematic literature review for sentiment analysis on social media in single languages from 2018 to 2021. Comparing [5] [13] and [39] with our literature survey, there is an overlap from 2010 to 2018 but [5], [39] provides very little information about recent methods and how the MSA methods work. However, our literature survey includes prior work and the most recent year's work on African languages. According to the best of our knowledge, there is currently a lack of systematic literature survey for MSA that is published, which covers rule-based/knowledge-based, cross-lingual with machine translation methods, traditional machine learning and deep learning models for under-resourced languages from 2008 to 2022.

### D. QUALITY ASSESSMENT

The quality assessment process shown in Fig. 3 was based on the following predefined questions:

- Are the aims of the study clearly stated with objectives and answers our research questions?
- Does the study provide new or unique techniques or contribution in MSA for low-resourced languages?
- Are there any major challenges identified in the study?

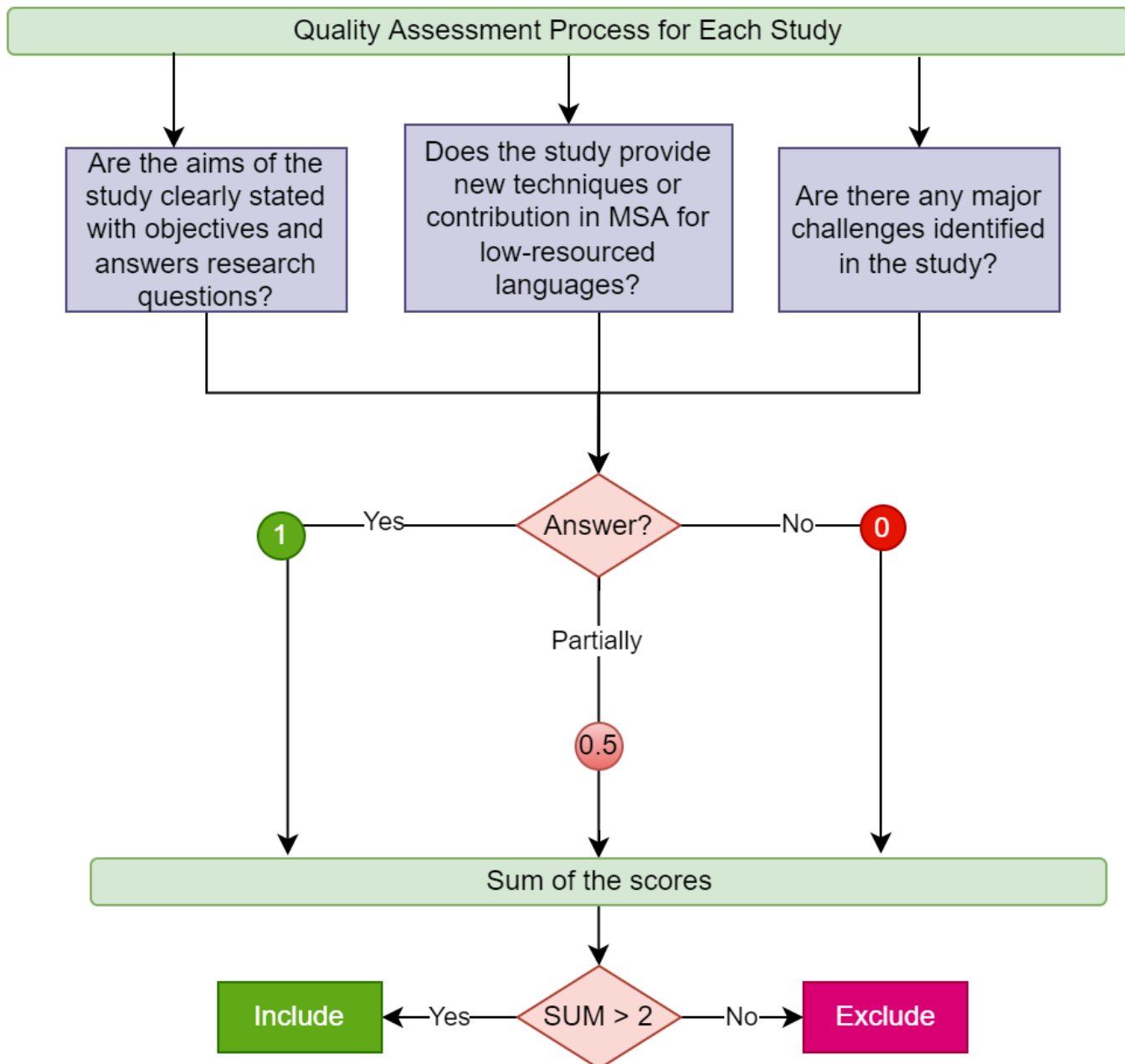


FIGURE 3. Quality Assessment process.

We selected several studies after excluding 84 articles. We assessed the quality of the research they presented. We used three quality assessment questions defined to evaluate the quality of the study and provide a quantitative comparison. To determine the quality assessment, we used the scoring procedure: Yes = 1, Partially = 0.5 or No = 0. Each study was given a score between 0 and 3. After the scoring, the points are summed for all the quality assessment questions. If the article received a non-integer total score, it was rounded to the nearest digit. A study is eliminated from our literature review if it receives a score of zero. Thirty-eight (38) articles with a score greater than two (2) are kept because they are

considered to meet this study criterion.

#### IV. STATE-OF-THE-ART STUDIES

Sentiment analysis is an active research field. Previous research has significantly improved the various tasks that makeup sentiment analysis systems. Even commercially produced technologies like Semantria<sup>1</sup> application programming interface (API), IBM Watson<sup>2</sup> and the iFeel tool for SA are accessible to the general public [15], [24]. Many studies have approached sentiment analysis research as a binary classifi-

<sup>1</sup><https://www.lexalytics.com/semantria>

<sup>2</sup><https://www.ibm.com/watson>

cation (i.e., *positive* or *negative*) or ternary classification (i.e., *positive*, *neutral*, and *negative*), and some have even gone as far as to investigate the fine-grained emotions [24]. For fine-grained emotions, researchers have investigated whether a text expresses emotions such as joy, happiness, love, or sadness. Several research studies have examined whether an objective or subjective sentence is positive or negative, as well as the subjective detection of that sentence [17], [45]. Another study for the Persian language focused on developing lexicon-based sentiment analysis [46] to evaluate Persian texts using online Persian language resources [47]. They used available texts from online and used native speakers to manually annotate the texts into positive, negative and neutral [46].

As social media platforms have grown in popularity, interest in studying several powerful sentiment analysis methods has increased. Progress in the field has moved from lexicon-based methods such as AFINN lexicon [48], Valence Aware Dictionary for Sentiment Reasoning (VADER) [49], SentiWordNet, SentiStrength [50], and statistical ML [15], [23] to DL methods like long short-term memory (LSTM), Bidirectional-LSTM and bidirectional encoder representations from transformers (BERT) [22], [51]. Twitter data is used mostly for NLP research, especially for sentiment analysis tasks [3], [9]. There are also platforms such as Amazon for sales reviews, music reviews, movie reviews and Twitter, which is the largest source of text datasets so far. The evolution of social media texts or microblogs (e.g. Twitter) has presented new opportunities for language technologies, but it has also posed many new challenges, making it one of the current prime research areas. Some interesting research has emerged using the Twitter dataset for multilingual sentiment classification in SemEval competitions [27], [36], and the introduction of code-switching texts has been studied for SA [52].

## V. MULTILINGUAL SENTIMENT ANALYSIS STRATEGIES

Various sentiment analysis methods for multilingual datasets have been explored [12], [15], [16], [21], [33]. There have also been proposed ways for classifying sentiment polarities on a multilingual dataset utilizing lexicon-based techniques along with MT systems and ML approaches [12], [21]–[23]. Most SA studies paid more attention to highly resourced languages than those with insufficient resources. However, because English language resources, such as sentiment lexicon, annotated corpora, and benchmark datasets, are easily accessible, most MSA approaches preferred strongly leveraging English language resources [22], [23], [53]. The details of MSA methods will be discussed in the following sections.

### A. MULTILINGUAL SUBJECTIVITY DETECTION

Previous sentiment analysis studies introduced the concept of subjectivity detection in multilingual sentiment. Subjectivity detection and sentiment analysis focus on identifying emotional states, such as opinions, emotions, feelings, evaluations, beliefs and speculations [17], [54]. Furthermore,

sentiment classification further refines the level of granularity by classifying subjective information as either positive, negative, or neutral.

Although there has been a lot of research on multilingual subjectivity detection, there is still a lot of room for future study in other languages [5], [45]. A lot of the research on the subjectivity detection task was done in English [55], [56]. As a result, most of the gold standard dataset is primarily written in English. Therefore, to create methods for detecting multilingual subjectivity, most studies attempt to use English-language resources [5], [57]. The lexicon and corpus methods dominated early research of multilingual subjectivity analysis [57]. They translated OpinionFinder (i.e., the English subjectivity analysis lexicon) to Romanian using a lexicon-based method and a lemmatized version of the English terminology. This research investigated the effects of corpus-based approaches on Romanian subjectivity-annotated corpora produced by translating English lexicons into Romanian.

Using linguistic resources in English, Banea et al. [58] investigated an MT-based method to conduct a subjectivity analysis of Romanian and Spanish. They used the MPQA corpus employed by Balahur and Turchi [14], which contains English-language news articles annotated for subjectivity from various sources. The authors showed that even though the translation system was employed, the results obtained were promising and comparable to those obtained by manually translating the corpora. Furthermore, Banea et al. [54] showed that using multilingual information, subjectivity classification (i.e., objective or subjective) in English could achieve 83% accuracy.

Banea et al. [59] explored the alignment of sense levels in different languages to reflect coherent subjectivity. The researchers claim that it is impossible to map one sense to another across languages because a particular purpose may have additional meanings or uses for a specific language. Additionally, they demonstrated that dual co-occurrence metrics could be used to model multilingual feature spaces, offering a more reliable model when compared to using individual languages as input. These metrics learn from comparable sense definitions. As a result of using a simple SVM classifier trained on multilingual space, the accuracy increased to 73% and 76% for English and Romanian, respectively. With an overall accuracy of > 73% across all iterations, they demonstrated that the multilingual model consistently outperformed its cross-lingual counterpart.

Another approach by [55] used a pre-annotated English corpus (i.e. 10,000 movie reviews) collected and annotated by Pang and Lee [56]. These reviews were obtained from the Rotten Tomatoes website<sup>3</sup> and IMDB<sup>4</sup> for subjective and objective reviews, respectively. They built a model to handle multilingual corpora annotated with opinion labels. Their models used Naïve Bayes (NB) techniques to classify

<sup>3</sup>[www.rottentomatoes.com](http://www.rottentomatoes.com)

<sup>4</sup><https://www.imdb.com/interfaces/>

reviews. In addition, they developed a method that can be used between topics and languages with high reliability using novel annotation methods. Parallel corpora of English and Arabic reviews are used for model evaluation. The findings show that the same annotations applied to English sentences in parallel corpora can also be applied to sentences in other languages [55].

## B. CROSS-LINGUAL METHODS FOR MSA

Several studies have employed MT to build sentiment analysis corpora for under-resourced languages [23]–[25], [28]. They utilised well-known MT applications such as Google Translate to translate a dataset existing in a high-resource language into an under-resourced language. However, translation quality is often affected by missing context information, cultural differences and lack of parallel corpora [21], [25], [60]. Some researchers proposed cross-lingual NLP approaches to solve the problem of low-resource languages by benefiting from high-resource languages like English [15], [25], [33], [61]. Previous sentiment analysis methods usually translate the comments from the original under-resourced language to English. This method allows the sentiment classification task to be performed on well-performing models. However, even though this approach was successful for high-resource languages like Russian, German and Spanish [62], it was reported in [60] that translation from English to German, Urdu, and Hindi had a harmful impact on the sentiment analysis performance. Ghafoor et al. [60] used Arabic social media comments to investigate the impact of MT on sentiment analysis performance. They reported that translation from English into German, Urdu and Hindi revealed poor sentiment analysis performance. According to studies on under-resourced languages, with the help of MT systems, cross-lingual sentiment analysis systems suffer performance degradation [23], [60]. Cross-lingual sentiment classification relies on MT approaches in which a source language is translated into the target language [16]. However, another challenge with approaches that rely on MT is that most APIs are not free of charge. Therefore, the task at hand may be costly when dealing with large text corpora [12]. Several authors have used MT systems to translate information directly from one language into another. However, due to differences in linguistic terms and writing styles, the translated data cannot cover the vital information found in the original data [63]. Some cross-lingual MSA approaches have been developed by training a sentiment polarity classifier in English and then employing MT, translating text written in another language into English and then applying a sentiment classifier. Fig. 4, shows the overview of a cross-lingual MSA method using MT-based techniques.

One approach uses SentiWordNet, which leverages English lexical resources to perform sentiment analysis [21], [29]. This approach focuses on extracting sentiments in languages other than English and then translating words into English using a standard MT system. Thus, translated documents are classified according to their sentiment, which is

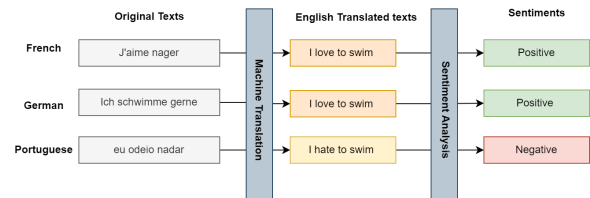


FIGURE 4. Overview of the MT-based approach with a simple example.

either positive or negative. The classification was performed by searching for sentiment-bearing words, such as adjectives, using SentiWordNet. The calculated score determines whether the words are positive or negative. This approach was investigated in German and English languages. The problem with this approach is that MT systems are not always accurate. Moreover, they have many issues, including data sparseness. The drawbacks of automatic MT systems have also been reported and further highlighted in studies [14], [23], [24].

Similarly, [64] proposed a unique way of leveraging reliable English resources to improve Chinese SA. Using MT systems, Chinese reviews were converted into English reviews, and then sentiment polarity in English reviews was identified by directly using English resources.

Balahur and Turchi [14] proposed an MSA approach that uses three distinct MT systems: Google Translate<sup>5</sup>, Bing<sup>6</sup> and Moses<sup>7</sup> translators. The approach used the English dataset from NTCIR 8 Multilingual Opinion Analysis (MOAT<sup>8</sup>). Three MT systems were employed to translate the gold standard dataset into French and German and build the training dataset and some testing datasets. This approach was further extended by using Yahoo systems to translate the dataset for testing into English, French, and German languages [23]. The test dataset translated using Yahoo systems was later corrected and verified by an expert [33]. Sentences containing no sentiments were omitted from the dataset to retain sentences with positive and negative sentiments.

SVM sequential minimal optimisation (SMO) was employed to classify sentiments for all languages to build a classification model for these three languages. First, the SVM SMO classification models for the three languages were trained separately for each language [23]. Then, the second experiment was executed by combining all separately trained models for each of the three languages and employing the unigram and bigram features extracted from the dataset [5]. The average performance accuracy reported for this approach is < 60% for all three MT systems. However, Balahur and Turchi [14] observed that translated data increased features, sparseness and more issues separating positive and negative sentiments during the training phase. This happened due to the low quality of the MT's data, which led to a decrease in

<sup>5</sup><https://translate.google.com/>

<sup>6</sup><https://www.bing.com/translator>

<sup>7</sup><https://www.statmt.org/moses/>

<sup>8</sup><http://research.nii.ac.jp/ntcir/data/data-en.html>



performance accuracy. Moreover, the extracted features had insufficient information to allow the classifier to learn.

Balahur and Turchi [14] suggest that the quality of the MT process has implications for the set of features used to build models. According to Becker et al. [12], MT systems are costly, and the results are limited because of the quality of the data translation. In addition, monolingual datasets that are combined to train MSA classification models have shown no impact in improving performance accuracy.

### C. IMPROVING MT-BASED METHODS

Balahur and Turchi [23] attempted to improve the performance of an MT system for MSA to obtain the best possible results. They employed MT systems and a supervised ML technique to perform sentiment analysis on a multilingual dataset (i.e. English, Spanish and French). The MT system translated training and test data into a single language, and then a monolingual sentiment classifier was applied. They concluded that the MT system reached a level of maturity and obtained good performance for languages other than English. Although the translated data produced reliable training data, the approach did not address the drawbacks reported in [12], [22]. They also concluded that the gap in classification performance between the models trained in English and translated data was somewhat in favour of source language data. Nonetheless, with MT systems, there is room for translation errors [14]. Although this technique helps to disambiguate the use of specific words, it does not eliminate translation errors. In addition, adopting this approach requires a more reliable MT system for the accurate performance of MSA models. However, several attempts have been made to improve MT for MSA. Becker et al. [12] argue that even if a perfect MT is readily available, there is always a potential cultural difference between the source and target languages, which may have implications for final classification results. Consequently, the approach mentioned above may not be reliable and will therefore not address the task of MSA, particularly in under-resourced languages.

Similarly, [65] presented a standalone MSA method for English using a gold standard dataset and Google Translate system to translate the dataset from English into four other languages (Italian, Spanish, French and German) to redesign their sentiment analysis system, which caters for data in multilingual settings. The approach employs a supervised learning method (i.e. SVM SMO) which was used previously [14], [23] with a linear kernel on unigrams and bigrams as features. In addition, they used tweet normalisation and MT to obtain high-quality training data for sentiment analysis in the four languages.

In their study, it was further shown that the joint use of training data from different languages, especially a closely related family of languages, can significantly improve the results of the sentiment analysis system. The authors claim that their proposed sentiment analysis approach can perform multilingual sentiment classification with up to 70% performance accuracy. The dataset used in their study was

sufficiently small for training and testing. A small dataset allows for easy manual correction of translation errors and eliminates incorrect translations [65]. Furthermore, Balahur and Turchi [65] claim that this approach can be extended to other languages using similar dictionaries created in this work. However, their study focused only on four different languages in which the dataset was presented in a monolingual setting to build a multilingual system. Therefore, we can assume that this method may not be easily adopted for MSA in a mixed-language context.

### D. TRANSLATION WITH ML-BASED METHODS

Prior studies indicate that numerous sentiment analysis strategies have explored different methods, but these methods usually rely on lexical resources or ML techniques [21], [24]. Many of these existing methods involve adapting lexical resources without proper comparisons, and validations [15], [24]. Araujo et al. [24] took a different step in the field by evaluating 21 methods for English multi-linear sentence-level SA. These methods are compared with two language-specific strategies based on nine language-specific datasets consisting of Arabic, Dutch, French, German, Italian, Portuguese, Russian, Spanish and Turkish. The two language-specific strategies were a multi-language version of SentiStrength (i.e. ML-SentiStrength) and a commercial API for Semantria.

They investigated these methods to address the problem of multiple languages in SA. First, they used MT (i.e. Google translate) to translate texts from a specific language into English and then employed the existing English-based sentiment classification methods [15], [24] to translate texts from a particular language into English and then employed the current English-based sentiment classification methods in Table 1.

TABLE 1. Overview of MSA methods with classes

Methods	Classes	L	ML
AFINN	3	✓	
Emolex	2	✓	
Emoticons	2	✓	
Happiness Index	9	✓	
Opinion Finder (MPQA)	3	✓	✓
Opinion Lexicon	3	✓	
PANAS-t	11	✓	
Pattern.en	2	✓	
SO-CAL	3	✓	
NRC Hashtag	8	✓	
Stanford Recursive Deep Model	6		✓
Sentiment140Lexicon	4	✓	
SA sentiment analysis	4		✓
SentiStrength	3		✓
Umigon	3	✓	
Vader	3	✓	

These methods were evaluated and compared across all nine languages [24], [66]. Further details of the classification classes of these methods are presented by [15], [24]. Although the datasets were small, the researchers concluded

that the existing English methods performed better than the two language-specific approaches. In this regard, SentiStrength was shown to be the most accurate method for SA, but language-specific techniques significantly impacted MT-based approaches.

Similarly, Araújo et al. [15] evaluated the performance of 16 English-based methods (including Google prediction API) for multilingual sentence-level sentiment analysis across 14 languages. They added five more languages: Chinese, Greek, Hindi, Czech and Haitian Creole to expand their previous work. They compared these English methods with three language-specific methods, including the IBM Watson API commercial sentiment analysis system developed by IBM. The sentiment analysis approaches employed previously [24] were explored. They investigated how the methods of using MT systems addressed multiple languages and found that the MT strategy should be used as a baseline system for new MSA systems [15]. The goal was to evaluate how effectively non-English texts could be analysed using English sentiment analysis methods, and MT systems [67]. As a final contribution, they developed iFeel 3.0, a web-based framework tool for multilingual sentence-level SA [15], [68]. Their methods, datasets and codes of the research work are freely accessible online to the research community.

#### E. CO-TRAINING MSA METHODS

Pan et al. [69] investigated a cross-lingual approach that employed an annotated sentiment corpus in English to predict the sentiment polarities in the Chinese language. They used an MT system to create the training dataset. This approach used the co-training of the two models simultaneously and added lexical knowledge to improve model accuracy. Co-training is training two or more monolingual models of the languages involved to build an MSA model [63]. This approach showed that adding lexical knowledge could improve the accuracy of the sentiment classification model.

In another study by [45], [70], they used the co-training method to overcome the problem of cross-lingual SA. He used the cross-lingual method, which employed a readily available English corpus for Chinese sentiment classification, using the English corpus as training data [45]. He then exploited a bilingual co-training approach to leverage the annotated English resources to perform sentiment classification in Chinese reviews. An MT system translated English labelled documents into Chinese, and a similar approach was used to translate unlabelled Chinese documents into English. Fig. 5 shows a co-training model used for cross-lingual sentiment analysis where English labelled data is transferred to another language, then apply sentiment classification on English data. In this work, an SVM-based classifier was adopted for sentiment classification. The co-training models are designed to select the high-confidence samples suited for training data. However, the classifiers in each language view will increase the probability of adding incorrect labels to the training set. Furthermore, adding such samples increased the accuracy of the learning model but gradually decreased the

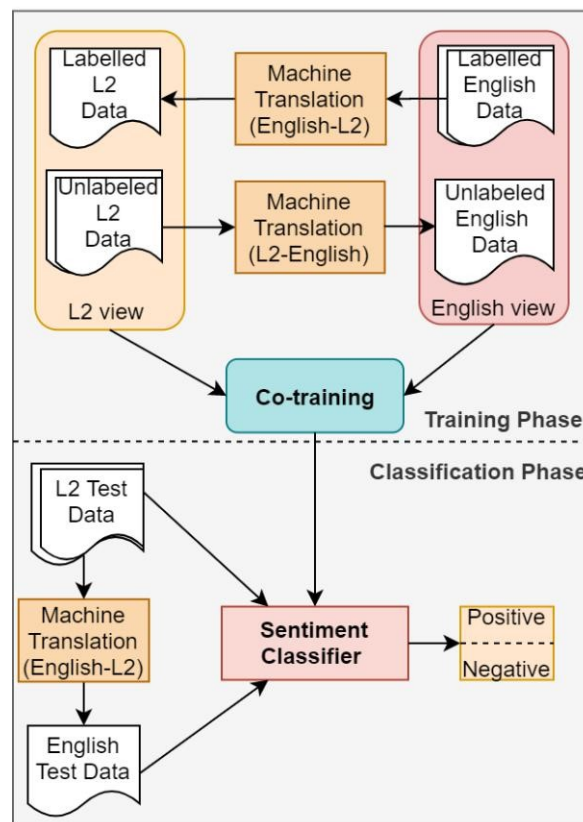


FIGURE 5. A co-training model for cross-lingual sentiment analysis.

performance of the initial classifiers. Nevertheless, the co-training approach can evaluate the different languages when the datasets are readily available.

Another approach incorporated an ensemble of English and Chinese sentiment classifications. Peng et al. [4] used an MT system, conducted an analysis of English and Chinese reviews, and the results were combined to improve the overall performance of the sentiment classification system. However, this approach is not reliable due to the poor results of the translation system when the domain knowledge is different from the target language. Furthermore, the technique used for the Chinese language cannot easily be adopted without modifying the language models. Thus, a sentiment lexicon in the target language is required for this method to work optimally on other languages and cannot be applied to other languages with no lexical resources. In addition, the MT system led to the accumulation of errors and reduced the accuracy of the translation. Despite the use of the MT system in this approach, a structural correspondence learning technique was applied to find a low-dimensional representation shared by the two languages at the feature level. This technique was done to reduce translation errors [71].

Although the most apparent solution to multilingual sentiment classification is by employing MT and using existing English methods to deal with multiple languages, which may not be an easy and reliable solution to MSA referred

to in our study (i.e., mixed-language texts in a sentence). Secondly, the earlier approaches have studied how sentiment analysis can be done for languages other than English using MT. In our view, these methods are "*short-cut*" solutions to address issues presented by multiple languages. Some of these methods have explored the MSA task as a 2-class polarity classification task, while others tackle MSA methods as a 3-class polarity detection problem (i.e., positive, neutral, and negative) for mixed-language comments [66]. However, an extra class is added in some methods to transform the MSA task into a 3-class polarity detection task.

## F. MULTILINGUAL CORPUS-BASED SENTIMENT ANALYSIS

The following section deals with studies focusing on MSA methods for low-resource languages rather than multilingual subjectivity analysis in multiple languages. Texts written in different languages pose a considerable challenge to sentiment polarity classification. However, [16], [28] proposed a multilingual approach that addresses the problem of sentiment polarity classification using Twitter data from different languages. They employed and compared three techniques in English and Spanish and used three ML models to address the issues presented in other languages. As a result, the following models have been developed [28], [61]:

- **Multilingual approach model:** This approach is achieved by training a multilingual dataset that does not require prior language identification or recognition phases. To accomplish a multilingual model, they merged two or more monolingual datasets to train and develop a single-pass multilingual sentiment classifier.
- **Dual monolingual approach models:** These two or more monolingual models know the origin of the text's language. Each model is trained and adjusted by using a monolingual corpus. In this case, the correct monolingual model is executed for sentiment classification once the language of the text is known.
- **Monolingual pipeline with language detection model (pipe model):** This model acts based on the decision provided by the language identification tool. This approach identifies a language given unknown text using a language identification tool. The training is similar to the monolingual system, as the language of the texts is known before using the correct sentiment classifier.

Vilares et al. [16] evaluated sentiment analysis approaches using Spanish and English datasets. They created a Spanish-English multilingual and monolingual English and Spanish models with language detection tools [28], [61]. They used monolingual corpora from the SemEval 2014 task-B corpus and the TASS 2014 corpus for English and Spanish, respectively. These monolingual corpora were combined to create a multilingual corpus for training and testing the classification model. L2-regularised logistic regression was employed for sentiment classification, which was then compared with a super-supervised model based on the bag-of-words.

Four features were considered: words, lemmas, psychometric properties and part of speech tags. The word features are obtained using a simple statistical model for counting word frequencies in texts; psychometric properties refer to emotions such as anger or topics (e.g. job) that commonly appear in messages [16].

Three approaches are evaluated using monolingual, synthetic multilingual and code-switching corpora of English, and Spanish tweets [28]. First, code-switching for the testing set was obtained by filtering tweets containing Spanish and English words. Then, three annotators labelled these tweets manually using the SentiStrength strategy, which uses a dual score to indicate positive or negative sentiments [24]. The conclusion was that the multilingual model approach was the best option when Spanish was the majority language. It was due to the high number of English words in Spanish tweets. Furthermore, monolingual models with language detection performed well only when English was the dominant language. Again, this was because of the lower number of Spanish words in the English corpus. Therefore, the monolingual approach cannot be used for multilingual settings. They also reported that the monolingual (pipeline) model with a language identification tool performed worse on the code-switching test set for most of the features used. Finally, the multilingual model approach obtained the best performance of 59.34%, using features such as lemmas and psychometrics. The lemmas are simply terms labelled using set rules to reduce data sparsity. However, in general, the atomic set of features, such as words, psychometric properties, or lemmatisation and their combinations, performed better under the proposed multilingual model approach [16].

The proposed multilingual model approach appears more robust in environments containing code-switched tweets and tweets written in multiple languages. However, again, it was concluded that neither dual monolingual nor multilingual strategies based on language detection are optimal for addressing code-switching texts. Notably, the performance accuracy of these systems on the experimented features was still < 70%, even after the improvements reported by [61]. In addition, they felt it would be interesting to explore the performance of MSA using neural deep network methods. Finally, the authors suggested that using DL architectures can help deal with code-switching texts [18], [61].

Tho et al. [72] investigated a code-mixed sentiment analysis of the Indonesian language and Javanese language using a lexicon-based approach. The authors compared two translated lexicon models such as SentiNetWord and VADER. They collected 3,963 tweets from two accounts that provide code-mixed tweets. The results of the manual labelling with the lexicons mentioned above showed that SentiNetWord outperformed the VADER lexicon. However, the overall performance showed that the VADER lexicon performed better than SentiNetWord.



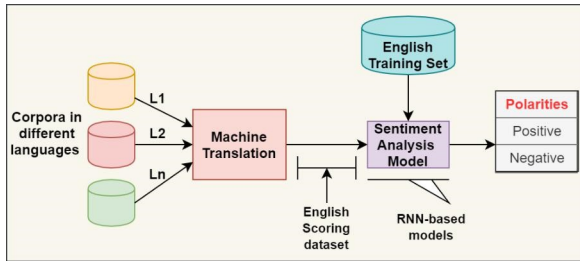


FIGURE 6. Multilingual sentiment analysis approach using RNN models [25].

### G. DEEP LEARNING METHODS FOR MSA

Recently, [25] presented an approach for MSA based on an RNN framework which aimed to answer the question: "Can a sentiment analysis model trained on a language be reused for sentiment analysis in other languages, Russian, Spanish, Turkish, and Dutch, where the data is more limited?" A single multilingual sentiment model utilising English data was developed for four different languages. The approach was built by training sentiment models using RNN methods with English reviews. An MT strategy was employed, translating Russian, Turkish, Dutch and Spanish studies into English and then reusing the English-based RNN model to classify sentiments. Fig. 6 shows an RNN-based structure for the MSA task where an MT-based system is utilised to translate the text in non-English language to English and apply deep learning sentiment classification on English data.

The MSA approach in this study was developed to eliminate the need to train language-dependent models and sentiment word embeddings in four languages. Can et al. [25] claim that other languages with low resources can utilise this multilingual model. In their study, the method was compared with a lexicon-based method that uses SentiWordNet to obtain positive and negative sentiment scores considerably better, with an accuracy of  $> 80\%$  for the three languages and  $74\%$  for Turkish. The RNN-based model eliminates the feature extraction process. Can et al. [25] concluded that the RNN-based model performed significantly better than language-specific models in all four languages, despite the misclassification encountered during translation. Furthermore, their study offers a solution that employs a single multilingual model, but it does not consider mixed-language texts.

Previous MSA methods have utilised English methods for which robust classifiers are readily available [15], [61]. Another work by [12] proposed an approach for the multilingual sentiment classification of Twitter data (i.e. English, German, Portuguese and Spanish), namely an efficient translation-free DL architecture to perform MSA on tweets in multiple languages. This approach was implemented by employing cost-effective character-level embeddings and Adhoc convolutions to learn a different language. The MSA model could learn hidden features from the four languages used during the training phase in their study. The authors compared their work with three different neural network architectures [73].

Each neural network was trained using different embedding strategies.

Their study concluded that the proposed multilingual approach achieved the best performance accuracy with the LSTM-based model. The LSTM models also performed optimally on all four languages when the  $F_1$ -score was evaluated. However, this classification model did not undergo the pre-processing phase, affecting performance accuracy [12]. The authors claim that these models can be extended to other languages and handle texts written in multiple languages. Furthermore, they suggest that employing a fully supervised CNN model will increase performance accuracy. However, it is not worth the time manually labelling thousands of tweets; therefore, a different approach can be investigated. Finally, the authors argued that a multilingual strategy offers several advantages over a language-specific sentiment model.

Deriu et al. [22] proposed a novel approach for multilingual sentiment classification of short texts in four languages (i.e. English, Italian, German and French) to enhance the system's ability to deal with mixed languages. They used weakly supervised data trained only on a CNN method for up to three layers. This approach trains multi-layer CNN where word embeddings (i.e. word2vec) are created on a large corpus of unlabelled tweets. Word embeddings are generally numerical representations of words input into DL-based methods. These are used for language modelling and feature learning (i.e., Word2Vec, GloVe and FastText) [13], [51].

The CNN model was trained in an unsupervised phase, where word embeddings are created on a large corpus, a distantly supervised step trained on the weakly labelled dataset, and a supervised stage, where the network was fully trained manually annotated tweets [73]. They evaluated the performance of the sentiment model with different datasets, including the benchmark sentiment prediction dataset from SemEval-2016 Task 4. They demonstrated that a single-CNN model could be trained successfully for MSA tasks rather than separate classification models for each language. However, the performance of the model can be improved by training a large number of convolutional layers. This method can be easily extended to new languages, and multilingual texts [22]. Deriu et al. [22] concluded that CNN models require a large amount of training data for the model to perform well, as well as the labelled dataset.

Similarly, [74] followed an approach that achieved the best results in the SemEval 2017 task. Even the system by Nguyen and Nguyen [75] that employs deep CNN and Bi-LSTM has shown that word2vec strategies can significantly improve classification accuracy. Additionally, recent improvements in DL techniques, especially the combination of CNN and LSTM techniques, have produced greater accuracy than per-language models [20], [74].

Medrouk et al. [76] proposed an approach which employs a deep neural network for sentiment analysis in a multilingual corpus. The deep neural networks used in their study employed CNNs (i.e. feed-forward).



**TABLE 2.** Summary of the methods, corpus, and techniques for MSA.

Sources	Methods	Languages	Sentiment Corpus
[29]	MT, SentiWordNet	English, German	2,000 reviews
[54], [58]	MT, NB, SVM	English, Romanian, Arabic and Spanish , French, German	MPQA MPQA corpus - 9,700 labelled sentences
[55]	NB	English and Arabic	Movie reviews, news corpora TED talks corpus - 183,000 manually annotated sentences
[77]	SVM, NB, ME	English, Dutch, French	2,500 blogs reviews and news sites.
[45], [70]	SVM	Chinese, English	Amazon product reviews (8,000)
[17]	MT, NB, DT, SVM	Chinese, English, Japanese	MPQA corpus - 535 for English, newspaper headlines - 1,200 movie reviews-12K (objective , subjective)
[14], [23]	MT, SVM	English, Spanish, French	Annotated NTCIR 8 Multilingual Opinion Analysis Task (MOAT - 6,200)
[33], [65]	SVM SMO	Italian, Spanish, French, German	Annotated NTCIR 8 Multilingual Opinion Analysis Task (MOAT - 6,200)
[67]	AFINN, SentiStrength, Vader MPQA, SO-CAL, Opinion Lexicon SASA, Umigon, Sent140lex Emoticons, NRC, Pattern.en PANAS-t, Opinion Finder, EmoLex Stanford Recursive Deep Model	Arabic, Dutch, French, German, Italian, Portuguese Russian, Spanish, Turkish	Human labelled products, food reviews and tweets - 23,000
[15], [24]	AFINN, ML-SentiStrength, Vader MPQA, SO-CAL, Opinion Lexicon SASA, Umigon, Sent140lex Emoticons, NRC, Pattern.en PANAS-t, Opinion Finder, EmoLex Stanford, Semantria, IBM Watson	Arabic, Dutch, French, German, Italian, Portuguese, Russian, Spanish, Turkish English, Croatian, Hindi, Haitian Creole, Chinese	Human labelled products, food reviews and tweets - 23,000
[16], [28], [61]	L2-RLR	Spanish and English (mixed)	SemEval 2015 Task B ( 8,200 tweets) TASS 2014 - (7,200 tweets)
[78]	LSTM and MT	Chinese and English	NLP& CC 2013 (91,600 unlabelled reviews)
[12]	CNN, LSTM	English, German Portuguese, Spanish	Annotated tweets (128,200, subset of 1.6M)
[76]	CNN	English, French, Greek	Labelled restaurant reviews (62,600)
[22]	CNN + word2vec	English, Italian, German, French	Unlabeled tweets (300M), weakly-labelled data (40-60M), & annotated tweets (71,000)

*continued on the next page*

Sources	Methods	Languages	Sentiment Corpus
[36]	CNN	English, Spanish, French, Russian, Arabic, Dutch, Turkish, Chinese	SemEval 2016 Task 5 corpus
[79]	CNN	English, German, Portuguese, Spanish	Annotated tweets (128,000, subset of 1.6M)
[11]	BiLSTM	English, Hindi	English annotated tweets (114,000) Hindi-English labelled Facebook posts (3,800)
[25]	MT, RNN	Russian, Turkish, Dutch, Spanish	Yelp Dataset Challenge (8,000), Amazon reviews (8M) Kaggle competition (68,000) Restaurant reviews
[80]	CNN	English, Spanish, Dutch, German, Russian, Italian, Czech, Japanese, French	Movie reviews (12,200 ), labeled reviews (20,800, TripAdvisor, Amazon Fine Food & labeled tweets (4,800)
[81]	LSTM	English, French, Greek	Labeled restaurant and hotel reviews (91,800)
[35]	GAN + DAN	English, Chinese, Arabic	Annotated tweets (48,100) & Weibo posts (53,600)
[32]	BiLSTM	English, Bengali, Hindi, Kannada (English mixed)	Labelled Facebook comments (22,500)
[18]	CNN	Bambara & French (mixed)	Facebook comments (17,000, subset of 74,000)
[82]	SVM_Linear, SVM_RBF	English, Greek	(5,300 reviews)
[15]	AFINN, SO-CAL	Arabic, Dutch, French, German, Italian, Portuguese, Russian	Human labelled products, food reviews and tweets (23,000)
[83]	CNN, Bi-LSTM	English, Hindi	Facebook posts (3,800)
[19]	Double LSTM	English & Hindi (mixed)	Hindi-English labeled sentences of Facebook posts (3,800)
[26]	Bi-LSTM	English, Hindi, Bengali	Annotated tweets: English (9,200), Bengali-English (5,500) English-Hindi (18,400) labeled Facebook posts English-Hindi (3,800)
[84]	mBERT XLM-RoBERTa	Hindi-English, Spanish-English, Tamil-English, Malayalam-English. Malayalam-English	Hinglish Tweets (14,000), Tanglish - YouTube comments (9,600) Spanglish - tweets (12,000) YouTube comments (3,900)
[10]	XLM-RoBERTa	Malayalam-English, Tamil-English YouTube Tamil-English	6,740 YouTube comments comments (15,740)
[72]	SentiNetWord VADER	Indonesian, Javanese (mix)	3,963 tweets
[85]	BERT	Persian-English	3,640 labeled tweets
[86]	LR, NB, DT, RF, SVM, BERT, DistilBERT, ALBERT RoBERTa,XLM, XLM-R Character BERT	Tamil-English, Kannada-English, Malayalam-English	60,000 YouTube comments 44,000 - Tamil-English, 20,000 Malayalam-English
[87], [88]	NB, SVM TeluguSentiWordNet	Telugu-English	15,744 YouTube comments
[89]	CNN, LSTM, GRU, BiGRU BiLSTM, BiLSTM+CNN, LSTM+CNN BiGRU+CNN, GRU+CNN, XLM-R	Malayalam- English	7,000 comments (FIRE 2020) 20,000 comments (EACL 2021)

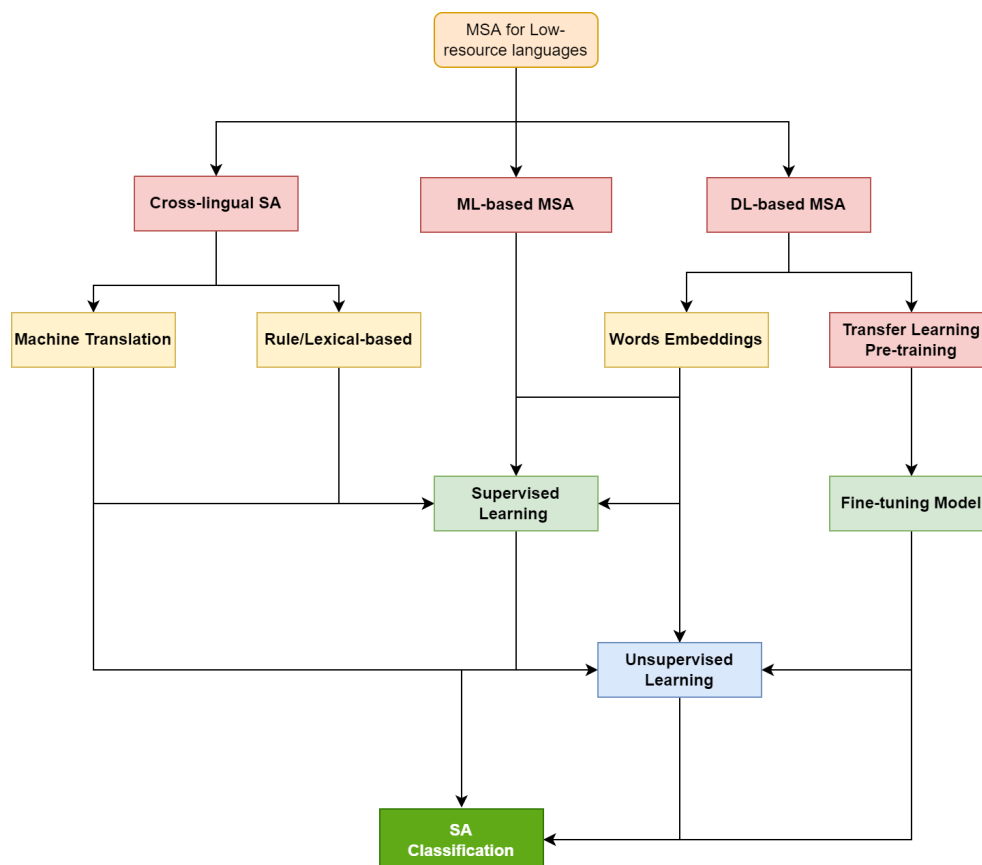


FIGURE 7. A general taxonomy of the MSA for low-resource languages.

The authors constructed their multilingual opinion corpus in three languages (English, French and Greek). The CNN exploited n-gram level information, and the system achieved high accuracy for sentiment polarity prediction. They concluded that the model used for feature extraction was language-independent. Another study by [83] presented a hybrid architecture for sentiment analysis of English-Hindi code-mixed data. They trained sub-word-level representations for sentences using the CNN model and employed a dual-encoder network consisting of two Bi-LSTMs. The model combined a network of orthographic features and word embeddings and achieved the best results with an accuracy of 83.54%.

Zhou et al. [78] proposed a cross-lingual sentiment classification approach using attention-based bilingual LSTM networks. The attention-based bilingual representation learning model was used to learn document distributed semantics for both the source and target languages. This approach was implemented for languages like Chinese and English. The authors used Google Translate MT to translate the training data into the desired target languages and then employed a bidirectional LSTM network to model the documents for both the source and target languages. A dataset from the cross-language sentiment evaluation of NLP&CC 2013 was used [78]. Reviews are divided into three categories: books,

DVDs and music. On average, the bilingual model yielded an accuracy of 82.4% across all domains. They concluded that LSTM could capture the compositional semantics of bilingual texts and that the proposed model achieved promising results on the dataset used [13], [78]. It also outperformed the best results in NLP&CC cross-language systems. An interesting part of their study is that the attention model could find the key sentences in a document, and the sentiment signals were captured with the help of word-level attention [22], [78].

Another work by [90] proposed a character-to-sentence CNN method that exploits characters to sentence-level information to perform sentiment analysis of short texts (i.e. Twitter data). They used the CharSCNN model with two convolutional layers to extract semantic information from word and sentence features to improve the performance of the sentiment analysis system. Irsoy and Cardie [91] improved sentiment classification accuracy using an RNN model on time-series information to obtain sentence representations. Socher et al. [34] improved sentiment analysis by using a recursive neural tensor network model, which synthesises the semantics of the syntactic tree of binary sentiment polarity. A good sentiment classification accuracy was also obtained using a tree-structured LSTM model with semantic association [92]. Furthermore, Baziotis et al. [93] presented an attention strategy for the LSTM model to achieve good sentiment

analysis results on the SemEval-2017 Task-4 dataset.

#### H. CODE-SWITCHED MSA METHODS

To some extent, most of the research on the code-mixed text has focused on the English-Hindi setting [11], [19], [83]. However, Code-mixed challenges can be addressed by learning the sentiment feature space and preserving the similarity of the sentences in which the sentiment is portrayed [11]. It allows a straightforward measure of the relatedness between code-switched content and labelled data from a resource-rich corpus. Choudhary et al. [11] demonstrated this using Siamese Bi-LSTM with tri-gram embeddings and a fully connected layer. Additionally, they compared the model, which was trained with a pair of Hindi-English texts, one with pairs of English sentences and another with code-mixed texts, yielding a lower *F-score* by 8.7% to the other with 75.9%. Their results suggest that adding more resource-rich data (i.e. the English dataset) is beneficial, as it increases model performance.

In other cases, most researchers have used DL techniques to model sentiment analysis for code-switched datasets. For example, Konate and Du [18] used Facebook comments of Bambara-French with different DL architectures such as LSTM, Bi-LSTM, CNN and Bi-LSTM-CNN, together with embeddings as input at the word or character level. Their proposed model can learn multilingual embeddings from input characters and words to mitigate the embeddings from the language model (Elmo) lack of pre-trained embeddings in the code-mixed corpus. They obtained the best results using a single-layer CNN model with an accuracy of > 80%. In addition, they compared LSTM and CNN, where the latter showed the best results in such a domain.

Furthermore, Kusampudi et al. [88], proposed a sentiment analysis in code-mixed Telugu-English text with Unsupervised Data Normalization. They reported accuracy of 80.22% on this dataset using novel unsupervised data normalisation with an MLP model, which is an increase of 2.53% accuracy due to this data normalisation. According to Saikrishna et al. [87], combining Telugu and English or Tamil and English in the same sentence is commonly observed. Saikrishna et al. [87] developed a sentiment analysis system for Telugu-English code-mixed sentences. They classified the polarity of the code-mixed sentences collected from Youtube comments into positive and negative sentiments using lexicon-based approaches and ML approaches such as NB and SVM classifiers. They achieved an accuracy of 82% and 85%, respectively, outperforming the lexicon-based method.

#### I. PRE-TRAINED METHODS FOR MSA

The BERT model has been applied to several NLP tasks. BERT has demonstrated outstanding performance in state-of-the-art text classification, including MSA tasks [51], [94]. BERT is a Google model created and reported in [51]. It was created using a significant amount of plain text data freely available on the Internet, and the model was trained

unsupervised. Before BERT, a few other pre-trained language models used bidirectional unsupervised learning.

The ELMo is one such model that focuses on contextualised word representations [95]. It constructs word embeddings by utilising LSTM, which separately trains left-to-right and right-to-left word representations and then concatenates these embeddings [95]. On the other hand, BERT does not use LSTM to obtain word context characteristics; instead, it employs attention-based transformers [51]. These models are beneficial for low-resource languages when there is a large amount of unlabelled data, but not much task-specific labelled data. Tomohiro et al. [96] presented a novel model that learns from sentences by labelling them with emojis, utilising English and Japanese tweets to create the corpus. The authors validated and evaluated many models based on attention LSTM, CNN and BERT. In addition, they compared the BERT model with the standard models CNN, FastText and attention Bi-LSTM, all of which received good results in prior studies. Compared to the traditional models, the authors performed better using the BERT model.

Gupta et al. [84] maintain that their unsupervised model understood code-switched languages or learnt only their representations. They introduced an unsupervised self-training method as a generic framework and demonstrated its applicability to the specific use of code-switched data. They exploited the power of pre-trained BERT models to initialise and fine-tune them using only pseudo-labels generated via zero-shot transfer. Their study was conducted in four code-mixed languages: Hinglish (Hindi-English), Spanglish (Spanish-English), Tanglish (Tamil-English) and Malayalam-English. They concluded that their unsupervised models outperformed their supervised counterparts, with performance ranging from 1% to 7%. Another study by [85] used a pre-trained multilingual BERT model to learn the polarity scores of these tweets for code-mixed Persian-English sentiment analysis. They collected tweets and employed two annotators to label the code-mixed tweets. Their multilingual BERT (mBERT) model outperformed the baseline models that use NB and random forest (RF).

Ou and Li [10] proposed a system to identify the sentiment polarity of the code-mixed dataset of the Dravidian dataset. They built on a pre-trained multi-language model such as the cross-lingual language model Roberta (XLM-Roberta) [97], and their system employed a k-folding approach to the ensemble and addressed the sentiment analysis problem of multilingual code mixed across language models. They took part in two code-mixed language challenges (Malay-English and Tamil-English). Their system had the highest *F-Score* of > 0.7 in Malayalam-English and ranked third in Tamil-English with an *F-score* > 0.6.

Chakravarthi et al. [86] introduced a code-mixed dataset of the under-resourced Dravidian languages. They manually annotated a dataset from social media comments for three under-resourced Dravidian languages. For over 60,000 YouTube comments, the dataset was annotated for sentiment analysis and identifying offensive language. The collection



includes roughly 20,000 comments in Malayalam and English, 7,000 comments in Kannada, and 44,000 comments in Tamil [98]. Unpaid volunteers manually annotated the data, and Krippendorff's alpha indicates a high level of inter-annotator agreement. Utilising machine learning and deep learning techniques, they provided baseline studies to create benchmarks on the dataset with the highest accuracy of 71% with the XLM technique. Notably, traditional machine learning methods have suffered low performance.

Thara et al. [89] investigated two major aspects of the code-mixed text: offensive language identification and sentiment analysis for Malayalam–English code-mixed data set. Their framework utilises different word embedding methods, such as Word2vec and FastText. They evaluated different deep learning methods (CNN, LSTM, Gated Recurrent Unit (GRU), BiLSTM, and Bidirectional GRU (BiGRU) ) on Forum for Information Retrieval Evaluation (FIRE<sup>9</sup>) 2020 and European Chapter of the Association for Computational Linguistics (EACL<sup>10</sup>) 2021 dataset. Among the hybrid models, GRU+CNN and BiLSTM+CNN turned in the highest F1-score of 0.9969. The challenge with this study is that the training dataset for sentiment analysis was minimal. They obtained the best performance accuracy of 99% using the transformer-based model XLM-R. Next, we will outline the evaluation metrics for MSA.

## VI. EVALUATION METRICS FOR MSA

In addition, we looked at evaluation measures for sentiment classification model performance. Several evaluation metrics have been identified from the systematic literature review [18], [22], [99]. These are reported in the SemEval 2016 challenge [36], averaging the macro F1-score of the positive and negative classes. Confusion matrix is also used as evaluation parameter to measure sentiment analysis performance [12], [13], [22]. Researchers use four metrics such as: *True positive*, *True negative*, *False negative* and *False positive*. These metrics are described as in [12], [22]. Metrics such as *accuracy*, *precision*, *recall* and *F-score* are generally used to evaluate the performance of the sentiment analysis classifiers [12], [24].

## VII. RESULTS AND DISCUSSION

In this section, we describe the results and then further discuss them.

### A. RESULTS

We explore how we answer our research questions as we present the findings. Research question 1 aims to identify MSA datasets and resources for under-resourced languages. Research question 2 aims to determine if MT methods are suitable for building MSA systems. To achieve this, we used the information from the literature review. A summary of the methods from the selected studies is presented in Table 2.

Additionally, Table 3 shows a quantitative summary of the results of our research questions 1 to 4. Our results show that DL methods at 61% were the most common MSA techniques for multiples languages, followed by ML methods at 40% and lexicon methods at 37%. Lastly, 29% of these studies used MT systems to help build their MSA resources. From the results in Table 3, we can conclude that DL methods are the leading techniques for MSA, including those where the English language is mixed with other languages, followed by ML and lexicon methods. In this case, lexicon-based, ML and MT methods were almost equally adopted by some of the studies presented in Table 2. Additionally, our research reveals that 63% of the publications studied ternary classification, while 37% and 13% looked into binary and five-category classification. Furthermore, only 31% of the DL methods have explored the binary classification, while 45% of the DL methods have concentrated on ternary classification.

Furthermore, we identified languages studied for MSA and the methods focused on mixed languages. Table 5 presents a summary of results for the languages involved. 42% of the selected studies used two languages in their proposed methods, 18% of the selected studies used three languages, and about 21% used four languages. Four studies with 3% used five, six, eight and fourteen languages in their studies, 34% of these studies have focused on code-switched datasets and only 8% used nine languages. Moreover, 63% of these studies have tackled MSA as a 3-class problem followed by a 2-class problem at 37%. Furthermore, out of the 34% of the studies which focused on a mixed-language dataset, English is the most commonly mixed language except for studies where the French language was mixed with the Bambara language, and Indonesian was combined with the Javanese language. Furthermore, we have noticed that the Persian language is mixed with English. In table 4, we presented our top thirty-five articles, which are highly cited. We also indicate their FWCI score which shows how well cited the article is compared to other similar articles. It is suggested that a value greater than 1.00 means the document is more cited than expected according to the average as in Table 4. Although some studies used a mixture of English, French, Hindi, and Bambara, their classification techniques are mainly based on DL methods.

### B. DISCUSSION

Table 2 and Table 3 are the summarised versions of the different methods and techniques used for sentiment analysis in multiple languages. Although most of the research studies utilised English methods, it is still a challenge for languages that do not have sufficient resources. Even methods that use MT systems are unreliable in tackling the task of MSA, owing to the limitations of MT systems [23], [33], [65]. Furthermore, many state-of-the-art sentiment analysis classification methods are based primarily on supervised learning algorithms. This means that a large amount of manually labelled data is required. Therefore, there is currently an immense

<sup>9</sup><https://dravidian-codemix.github.io/2020/>

<sup>10</sup><https://competitions.codalab.org/competitions/27654>

**TABLE 3.** Summary of the methods and techniques for MSA with their classification categories where L = lexicon, B = binary, T = ternary, F = five categories

Sources	L	MT	ML	DL	B	T	F
[29]	✓	✓	✓		✓		
[58]	✓	✓	✓		✓		
[55]	✓		✓		✓		
[77]		✓	✓			✓	
[45], [70]		✓	✓		✓		
[17]	✓	✓	✓		✓		✓
[14], [23]	✓	✓	✓			✓	
[33], [65]	✓	✓	✓			✓	
[16]	✓		✓		✓		
[28], [61]	✓		✓		✓		
[24], [67], [68]	✓	✓	✓			✓	✓
[12]				✓	✓		
[76]	✓			✓	✓		
[82]	✓		✓			✓	
[76]				✓	✓		
[22]				✓		✓	
[78]		✓		✓	✓		
[36]				✓		✓	
[79]				✓	✓		
[25]		✓		✓		✓	
[11]				✓		✓	
[80]				✓	✓		
[81]				✓	✓		
[35]				✓		✓	✓
[32]				✓		✓	
[18]				✓		✓	
[83]				✓		✓	
[15]	✓	✓	✓			✓	✓
[26]				✓		✓	
[37]				✓		✓	
[96]				✓		✓	
[84]				✓		✓	
[10]				✓		✓	
[72]	✓					✓	
[85]				✓		✓	
[86]			✓	✓		✓	
[87]	✓		✓			✓	
[89]				✓		✓	
<b>Total</b>	14	11	15	23	14	24	5
<b>%</b>	37%	29%	40%	61%	37%	63%	13.2%

need for techniques that require less human intervention [12], [22], [103] and even for data annotation of mixed-language texts. In this study, we found that research on MSA has shifted from lexicon or corpus-based and MT-based methods to a multilingual approach using DL techniques, which currently show incredibly encouraging results. Research on low-resource languages is gradually gaining momentum, and studies are paying more attention to code-mixed, and code-switched texts [26], [84]. It can also be noted from Table 2 that over 30% of the MSA studies preferred to use data collected from the Twitter platform. Moreover, most of these datasets were hand labelled by hand annotators despite efforts to build auto-labelling methods [50].

We are guided by the methods in the systematic literature

survey to draw a general taxonomy of MSA for low-resource languages, as shown in Fig 7. In Table 5, we can deduce that the number of languages studied increases yearly. The selected studies show that an increasing number of languages are gaining traction in the context of MSA. Several studies have used Twitter to develop sentiment analysis datasets. The fact that more languages are studied means that there are still more under-resourced languages to consider since many different languages are used in social media. Additionally, although researchers are using MT-based methods to generate language resources for under-resourced languages, the universal approach to cater for many languages is still far from being achieved. Therefore, there is an excellent

**TABLE 4.** Top 35 and leading articles for MSA studies with citations and field-weighted citation impact (FWCI)

No.	Sources	Year	No. of Languages	No. of Citations	FWCI
1	[29]	2008	2	233	7.82
2	[45]	2009	2	356	11.4
3	[77]	2009	3	291	5
4	[54]	2010	4	96	7.01
5	[17]	2011	3	175	14.42
6	[69]	2011	2	29	3.19
7	[33]	2014	4	123	7.94
8	[65]	2013	4	24	3.12
9	[16]	2015	2	36	1.17
10	[78]	2016	2	166	14.74
11	[100]	2016	2	100	3.44
12	[21]	2016	2	131	5.18
13	[36]	2016	8	54	7.45
14	[24]	2016	9	42	3.46
15	[28]	2016	2	21	2.13
16	[82]	2017	2	194	15.58
17	[5]	2017	3	106	12.67
18	[22]	2017	4	124	9.8
19	[61]	2017	2	48	4.14
20	[76]	2017	3	16	3.18
21	[32]	2017	5	17	2.53
22	[12]	2017	4	12	1.2
23	[11]	2018	2	33	4.75
24	[25]	2018	4	31	2.12
25	[18]	2018	2	15	0.69
26	[101]	2018	4	14	1.16
27	[83]	2019	2	28	4.6
28	[102]	2019	9	17	4.51
29	[19]	2019	2	5	1.88
30	[15]	2020	14	32	4.75
31	[85]	2021	2	4	-
34	[88]	2021	2	6	-
32	[86]	2022	4	24	3.94
33	[87]	2022	2	-	-
35	[89]	2022	2	-	-

**TABLE 5.** Number of languages used for MSA models including mixed languages

No.	No. of Languages	Studies	Total	%
1	2	[16], [21], [45], [78], [82], [83], [100]	16	<b>42%</b>
		[11], [18], [29], [55], [69], [82], [85], [87], [89]		
2	3	[10], [17], [26], [35], [76], [77], [81]	7	<b>18%</b>
3	4	[12], [22], [25], [32], [65], [81], [86], [101]	8	<b>21%</b>
4	5	[84]	1	<b>3%</b>
5	6	[54]	1	<b>3%</b>
6	8	[36]	1	<b>3%</b>
7	9	[24], [68], [102]	3	<b>8%</b>
8	14	[15]	1	<b>3%</b>
9	Mixed	[10], [12], [18], [22], [26], [28], [32], [61], [84]–[87], [89]	13	<b>34%</b>

opportunity for future sentiment analysis studies to focus on developing versatile techniques [13].

We noticed a shift in how MT systems were used for MSA research in languages with limited resources. Many studies used the monolingual dataset to address the issues of multiple languages [104]. In contrast, other studies looked at utilis-

ing a cross-lingual approach using MT systems [80]. The evidence is presented in Table 2, which demonstrates that DL methods account for 62% of all methods employed, followed by ML-based methods at 41%, lexicon-based methods (38%), and MT methods (35%). The MT-based techniques are used mainly because of their advantage of reproducing

the language resources from English to other languages where MT APIs are supported. It is, therefore, difficult for languages with limited resources and not supported by various MT APIs. It is also clear that only a few studies have utilised pre-trained models to fine-tune the MSA task in low-resource languages, although there is a significant increase in DL methods. This is mainly because the methods are still new in the NLP community. In addition, pre-trained models for low-resource languages still need to be explored. Annotated dataset remains a challenge for low-resource languages. For this, another approach used BERT methods to fine-tune the downstream tasks for low-resource languages [51], but these have been unable to outperform the existing DL methods [13].

Furthermore, multilingual aspect-based sentiment analysis is still in its infancy [105]. Despite its first study in SemEval 2016 Task 5, it produced the highest accuracy of 88.13% for English and the lowest of 73.35% for Chinese [36]. Efforts are devoted to comparing DL-based methods such as CNN, LSTM, or Bi-LSTM performance and improving performance by adding attention mechanisms. However, few methods focus on self-learning sentiment analysis classification, with less attention paid to multilingual contexts. Recently, a study by [106] used Twitter to develop NaijaSenti corpus (i.e. languages such as Hausa, Igbo, Pidgin, and Yorùbá), and they evaluated their corpus using mBERT, XLM-R and Roberta. This study demonstrated that model fine-tuning on pre-trained models performs well.

On a different note, while conducting this research study, we were able to identify several challenges with some MSA studies: (i) some of the research methodologies applied are difficult to follow to build baseline studies, (ii) some research cannot be easily replicated to obtain the exact results reported in the published research papers [5], (iii) some of the MSA research studies have not yet been released or published their datasets, tools and other resources for easy rebuilding or reproduction, (iv) although some studies provided Internet links for their resources used, some were not available for use or were not updated, and some resources and datasets are available only on request [15], [24].

The practical implications of this research are to identify the gaps that can be filled in the future and to set the trend of research shift concerning sentiment analysis of under-resourced languages. For diverse MSA datasets, our systematic literature review offers a variety of tools and techniques. This systematic literature review aims to provide several contributions, covering different application methods used for MSA for under-resourced languages. We focused on illustrating the contributions of each research work and observing the type of language-specific methods, transfer learning methods with MT systems, machine learning and deep learning algorithms used. Our investigations also focus on identifying the type of dataset used, how it was gathered, and how these datasets were annotated. Additionally, they used the environment and the performance measures covered in each study, evaluating them and concluding with appar-

ent research gaps and obstacles, which aids in identifying the non-saturated applications for which the MSA is most required in future research. For instance, aspect-based MSA deep learning systems need more sophisticated learning algorithms to produce better results.

Our findings indicate that deep learning methods are used in more than 60% of the studies, which is where significant research innovation can occur. This comprehensive literature review results bring alternative study directions for languages with limited resources and shed insight into current MSA research trends. Last but not least, this research will assist in identifying current, significant challenges in MSA for low-resource languages. Using computational models created for English or other rich languages with plenty of resources by many NLP systems puts technology developed for languages with limited resources at risk. It is beneficial to develop strategies that support languages with limited resources. We also support the development of platforms from languages with few resources accessible to everyone in other societies and using new NLP technologies for under-resourced languages. This comprehensive literature review, which examines MSA studies from the past and recent years, is anticipated to be helpful to other researchers in the future. The most popular datasets for MSA study are provided in Table 2.

## VIII. STUDY LIMITATIONS

The systematic literature study may have a few limitations. There may have been published papers we missed due to our selection criteria or search keywords, even though those studies examined the MSA approaches per the period specified in the methodology section. The review was conducted by only a few researchers, meaning there may be bias in selecting studies for comparison. Only original and unique studies published from 2008 to 2022 were included, with MSA studies with knowledge-based/lexicon techniques, multilingual subjectivity and MT-based methods for optimal comparison. The intention is to provide a complete picture of the origin of the MSA methods and the direction of progress, including ML and DL methods. Comparing techniques that appear to be out-of-date may also be a drawback. However, we believe they could help develop baseline systems for other languages and dialects.

## IX. EMERGING MSA AREAS

The findings from this research are discussed in this part, along with some recent development in the field that may necessitate further investigation.

**Textual representations strategies.** Text representation methods have been explored for low-research languages [95]. Research is still needed, especially for mixed-language contexts, on the desire to progress the text representation problem for under-resourced languages. Word2vec has known limitations in handling similar words. The ELMo approach was used to lessen these limitations [95]. ELMo extracts context-based word representations. For mixed languages,



research on cross-lingual word embedding has drawn considerable attention. Cross-lingual word embeddings are vector representations of words from multiple languages in the same vector space, allowing words with the same meaning but different languages to have the same vector representation. However, this technique is determined by the nature of the data requirements rather than the structure of the model. Despite this efforts, the question of which type of embedding captures better text features in the MSA task remains unanswered.

Recent research to address word representation in multilingual texts will increase as researchers are trying to study the effect of different languages and closely related language families. An interest in using DL models and pre-trained language models in under-resourced languages has grown in the most recent MSA models. However, many under-resourced languages continue to struggle with annotated datasets. Despite the advances in NLP, many under-resourced languages are still not covered by pre-trained models like BERT, RoBERTa, and XML-RoBERTa. The use of fine-tuning language models is one strategy for addressing this problem [106].

**MSA for pre-trained models.** The use of transformer-based models like BERT, RoBERTa, and mBERT allows researchers to focus on fine-tuning models for downstream tasks rather than training models from scratch. The recent interest in using pre-trained models for MSA tasks has become more useful for under-resourced languages with promising results for high-resourced languages [10], [51]. The research work of [37] presents extra language-specific pre-training for multilingual contextual word representations in a low-resource situation before usage in a downstream task. They enhanced the current vocabulary with frequent tokens from the low-resourced language (i.e. Irish, Maltese, Vietnamese, and Singlish (Singapore-English) and mimicked better language-specific phrases [37]. Chau et al. [37] suggests that we can improve the performance of the multilingual models on low-resourced language variations similarly by applying additional pre-training on language-specific corpora. They examined dependency parsing in four topologically varied low-resource language varieties with varying degrees of similarity to the pre-training data of a multilingual model. According to their findings, these methods consistently improve performance for each target variety, especially in low-resource conditions.

**Multilingual aspect-based Sentiment analysis.** Although there is an interest in continuing research on multilingual aspect-based sentiment analysis, there is a need to standardise the approach to this issue [36]. It is quite clear that this problem has not been widely addressed, especially with using DL-based architectures and considering under-resourced languages. Research on this topic suggests that the use of attention mechanisms and aspect-based embedding may significantly help resolve this problem. Perhaps MSA in an aspect-based context should be adopted for the methods that handle code-switched setups. It is essential to evaluate

whether the coupling of processes will impact word embeddings and sentiment analysis classification.

**Under-resourced languages.** High-resource languages from different language families have been studied extensively [32], [80]. Despite steady interest in under-resourced languages, code-switching and code-mixing remain challenging within multilingual communities, preferably mixing languages with high resources. A general approach to address these challenges is lacking. Although there is promising progress in some Indian languages like Tamil, Urdu [107], [108] and Telugu [109] and Iranian languages like Persian [85], [110], much is still required to develop models that employ deep learning techniques [111]. Also, they attempted to address challenges in Persian language by applying DL methods which only achieved the *f-score* of 55.5%. Ghasemi et al. [112] developed a sentiment analysis task in Persian language by proposing a cross-lingual deep learning framework to benefit from available training data of the English language. Deep learning models such as CNN and LSTM and their combinations were experimented with to achieve the *f-score* of 91.8% on LSTM-CNN. Recent work by [113] explored a cross-lingual sentiment analysis approach in Bengali language. Cross-lingual sentiment classification is another process to handle low-resource language issues. Bengali is considered a low-resourced language due to the scarcity of annotated data and the lack of text processing tools. They created and annotated a comprehensive corpus of around 12,000 Bengali reviews using the MT system and prior sentiment information to generate accurate pseudo-labels from English-based lexicons. The best F1 score of 0.897 is achieved by integrating LR and SVM classifiers as a hybrid method. For the SVM, the best accuracy of 91.5% was achieved. The decision tree (DT) based methods, RF and Extra Trees Classifier (ET), achieved the lowest F1 scores. sentiment analysis for monolingual, code-switched and multilingual comments in under-resourced languages has been studied only for a few African languages, e.g. several Nigerian languages [99], [106], Swahili [114] and Bambara [18]. Annotated datasets for MSA are lacking. A concerted effort to build datasets for sentiment analysis is required, especially for under-resourced languages such as African languages [99], including the languages in South Africa [115].

A study by [106] investigated the development of NaijaSenti—an introduction of the first African large-scale human-annotated dataset for sentiment analysis in Nigerian languages (i.e. Igbo, Yoruba, Hausa and Nigerian-Pidgin). The authors evaluated their methods on several pre-trained models such as mBERT [51], XML-R [97], mDeBERTaV3, mDeBERTaV3 and AfriBERTa [116]. They further evaluated their dataset using language-adaptive fine-tuning methods. To address this under-representation, AfriBERTa was developed—an African version of BERT trained from scratch to accommodate some of the African languages [116]. AfriBERTa has been trained in 11 languages: Afaan Oromoo (also known as Oromo), Amharic, Gahuza (i.e. a hybrid language that includes Kinyarwanda and Kirundi), Hausa, Igbo, Nigerian

Pidgin, Somali, Swahili, Tigrinya, and Yoruba [106], [116]. AfriBERTa was tested for named entity recognition and text categorisation in ten languages. In numerous languages, it outperformed mBERT and the cross-lingual language model with RoBERTa (XLM-R) [10], [97], and it was reported to be a competitive model.

Some of the NLP models derived from BERT, such as mBERT [51], RoBERTa [117] or XML-RoBERTa [97], were trained with many languages and can classify comments straight-forward from those languages. Unfortunately, XLM-RoBERTa [97] and mBERT are not trained with any data containing South African languages [118]. Most of these PLMs models cover 50 to 110 languages with only a few African languages, which are represented due to a lack of large monolingual corpora [116]. AfriBERTa, RoBERTa and XML-R have not yet been evaluated with any South African languages from the sentiment analysis perspective. IsiZulu, Sesotho and IsiXhosa are only now represented and assessed using a multilingual adaptive fine-tuning model [99] trained on XML-RoBERTa, and AfriBERTa [99], [116] but for a different NLP task.

In the context of South Africa, a concerted effort is required to create resources for South African under-resourced languages. The research could start by curating the SAfriSenti corpus—a multilingual sentiment corpus for South African languages and then expand to other African languages such as Lingala, Shona and Swahili and other African languages. In South Africa, there are 11 official languages. Languages like Sepedi (i.e. *Northern Sotho*), isiZulu, isiXhosa, Setswana, siSwati, Tshivenda, Xitsonga, and Sesotho (i.e. *Southern Sotho*) and their dialects are spoken by large populations. In addition, the lack of NLP resources for under-resourced languages makes it difficult to develop digital language technologies. Therefore, it is for these reasons that a massive data collection for under-resourced languages, in general, is necessary to address the under-resourced language challenges. Another exciting area for future research is code-switching between English and under-resourced languages.

**Automatic data annotation.** Supervised algorithms rely on the labelled dataset. Recent studies have investigated models which can be used to reduce human effort during data annotation [119]. For example, Kranjc et al. [103] developed active learning methods using pre-trained BERT language models. These techniques appear to be effective for text labelling vast amounts of data. Another emerging area for further research is the investigation of multi-class sentiment classification using active learning methods [119].

## X. CONCLUSION

As many MSA strategies have been proposed and experimented with, many studies have evaluated and contrasted the performances of different statistical ML models and DL-based methods using MT strategies for MSA over the years. However, there is currently no method identified as the best-performing model. In this study, we reviewed the most used

methods for MSA that apply traditional ML and DL models, together with those that employ MT-based methods. The performances of the MT-based methods and statistical ML and DL models reported in the literature were compared. Although some studies suggest that MT-based methods should be a baseline system for newly proposed MSA systems, these methods are still not proven for under-resourced languages. The literature results show that most DL architectures have recently spiked research attention compared to traditional ML-based classifiers, including even methods that rely on MT-based systems. Furthermore, the literature has proven that a combination of DL models such as CNN, Bi-LSTM and LSTM can significantly improve the performance of the MSA system. This study also highlighted the limitations of the use of MT-based methods. Furthermore, we propose a DL learning framework for MSA without using an MT-based system. We have also stressed that the BERT or XML-RoBERTa model can play a pivotal role in the performance of MSA models if it is fine-tuned to handle the downstream task.

Future studies on sentiment analysis must focus more on developing gold-standard datasets suitable for sentiment analysis of multiple languages. Researchers should focus more on developing sentiment lexicons for low-resourced languages so that, in future, they can concentrate on developing advanced models. Investigation of how much sentiment is lost in translation, for instance, when moving between multilingual and a single language versus the sentiment of the original text, should be studied to report MT challenges. This study did not consider single-language sentiment analysis evaluation; future research should focus on multilingual SA, particularly in under-resourced languages. Furthermore, future research should focus more on developing DL MSA models for multiple languages or robust languages with independent MSA techniques that can be used to analyse monolingual, numerous and mixed languages. An exciting research direction is to focus on methods that address code-switching sentiment analysis and multilingual aspect-based sentiment analysis in a multilingual setting. In addition, the significant challenges and current research gaps in MSA were reviewed. Finally, future directions for research in MSA will be investigated.

## XI. ACKNOWLEDGEMENTS

We acknowledge the National Research Foundation (NRF) for the Black Academics Advancement Programme (BAAP) grant (REFERENCE NO: BAAP200225506825).

## REFERENCES

- [1] W. Medhat, A. Hassan, and H. Korashy, "Sentiment analysis algorithms and applications: A survey," *Ain Shams Engineering Journal*, vol. 5, no. 4, pp. 1093–1113, 2014.
- [2] M. Wankhade, A. Rao, and C. Kulkarni, "A Survey on Sentiment Analysis Methods, Applications, and Challenges," *Artificial Intelligence Review*, pp. 1–50, 02 2022.
- [3] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up?: sentiment classification using machine learning techniques," in *Proceedings of the ACL-02*

- conference on Empirical methods in natural language processing-Volume 10, pp. 79–86, Association for Computational Linguistics, 2002.
- [4] H. Peng, E. Cambria, and A. Hussain, "A review of sentiment analysis research in chinese language," *Cognitive Computation*, vol. 9, no. 4, pp. 423–435, 2017.
- [5] S. L. Lo, E. Cambria, R. Chiong, and D. Cornforth, "Multilingual sentiment analysis: from formal to informal and scarce resource languages," *Artificial Intelligence Review*, vol. 48, no. 4, pp. 499–527, 2017.
- [6] P. D. Turney, "Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews," in *Proceedings of the 40th annual meeting on association for computational linguistics*, pp. 417–424, Association for Computational Linguistics, 2002.
- [7] R. Bhargava and Y. Sharma, "Msats: Multilingual sentiment analysis via text summarization," in *2017 7th International Conference on Cloud Computing, Data Science Engineering - Confluence*, pp. 71–76, 2017.
- [8] A. Ghafoor, A. S. Imran, S. M. Daudpota, Z. Kastrati, Abdullah, R. Batra, and M. A. Wani, "The impact of translating resource-rich datasets to low-resource languages through multi-lingual text processing," *IEEE Access*, vol. 9, p. 124478 – 124490, 2021. Cited by: 0; All Open Access, Gold Open Access, Green Open Access.
- [9] A. Pak and P. Paroubek, "Twitter as a corpus for sentiment analysis and opinion mining," in *LREC*, pp. 1320–1326, 2010.
- [10] X. Ou and H. Li, "Ynu@dravidian-codemix-fire2020: Xlm-roberta for multi-language sentiment analysis," in *FIRE*, 2020.
- [11] N. Choudhary, R. Singh, I. Bindlish, and M. Shrivastava, "Sentiment analysis of code-mixed languages leveraging resource rich languages," *CoRR*, vol. abs/1804.00806, 2018.
- [12] W. Becker, J. Wehrmann, H. E. Cagnini, and R. C. Barros, "An efficient deep neural architecture for multilingual sentiment analysis in twitter," in *The Thirtieth International Flairs Conference*, pp. 246–251, 2017.
- [13] M. M. Aguero-Torales, J. I. Abreu Salas, and A. G. Lopez-Herrera, "Deep learning and multilingual sentiment analysis on social media data: An overview," *Applied Soft Computing*, vol. 107, pp. 107–373, 2021.
- [14] A. Balahur and M. Turchi, "Multilingual sentiment analysis using machine translation?," in *Proceedings of the 3rd workshop in computational approaches to subjectivity and sentiment analysis*, pp. 52–60, Association for Computational Linguistics, 2012.
- [15] M. Araújo, A. Pereira, and F. Benevenuto, "A comparative study of machine translation for multilingual sentence-level sentiment analysis," *Information Sciences*, vol. 512, pp. 1078–1102, 2020.
- [16] D. Vilares, M. A. Alonso, and C. Gómez-Rodríguez, "Sentiment analysis on monolingual, multilingual and code-switching twitter corpora," in *Proceedings of the 6th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pp. 2–8, 2015.
- [17] C. Banea, R. Mihalcea, and J. Wiebe, "Multilingual sentiment and subjectivity analysis," *Multilingual natural language processing*, vol. 6, pp. 1–19, 2011.
- [18] A. Konate and R. Du, "Sentiment analysis of code-mixed bambara-french social media text using deep learning techniques," *Wuhan University Journal of Natural Sciences*, vol. 23, pp. 237–243, 06 2018.
- [19] S. Mukherjee, "Deep learning technique for sentiment analysis of hindi-english code-mixed text using late fusion of character and word features," in *2019 IEEE 16th India Council International Conference (INDICON)*, pp. 1–4, 2019.
- [20] S. Seo, C. Kim, H. Kim, K. Mo, and P. Kang, "Comparative study of deep learning-based sentiment classification," *IEEE Access*, vol. 8, pp. 6861–6875, 2020.
- [21] K. Dashtipour, S. Poria, A. Hussain, E. Cambria, A. Y. Hawalah, A. Gelbukh, and Q. Zhou, "Multilingual sentiment analysis: state of the art and independent comparison of techniques," *Cognitive computation*, vol. 8, no. 4, pp. 757–771, 2016.
- [22] J. Deriu, A. Lucchi, V. D. Luca, A. Severyn, S. Müller, M. Cieliebak, T. Hofmann, and M. Jaggi, "Leveraging large amounts of weakly supervised data for multi-language sentiment classification," *Proceedings of the 26th International Conference on World Wide Web*, pp. 1045–1052, 2017.
- [23] A. Balahur and M. Turchi, "Comparative experiments for multilingual sentiment analysis using machine translation," in *SDAD@ ECML/PKDD*, pp. 75–86, 2012.
- [24] M. Araújo, J. Reis, A. Pereira, and F. Benevenuto, "An evaluation of machine translation for multilingual sentence-level sentiment analysis," in *Proceedings of the 31st Annual ACM Symposium on Applied Computing*, pp. 1140–1145, 2016.
- [25] E. F. Can, A. Ezen-Can, and F. Can, "Multilingual sentiment analysis: An rnn-based framework for limited data," *arXiv preprint arXiv:1806.04511*, 2018.
- [26] A. Jamatia, S. D. Swamy, B. Gamback, A. Das, and S. Debbarma, "Deep learning based sentiment analysis in a code-mixed english-hindi and english-bengali social media corpus," *International Journal on Artificial Intelligence Tools*, vol. 29, pp. 399 – 408, 2020.
- [27] P. Nakov, A. Ritter, S. Rosenthal, F. Sebastiani, and V. Stoyanov, "Semeval-2016 task 4: Sentiment analysis in twitter," *arXiv preprint arXiv:1912.01973*, 2019.
- [28] D. Vilares, M. A. Alonso, and C. Gómez-Rodríguez, "En-es-es: An english-spanish code-switching twitter corpus for multilingual sentiment analysis," in *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pp. 4149–4153, 2016.
- [29] K. Denecke, "Using sentiwordnet for multilingual sentiment analysis," in *2008 IEEE 24th international conference on data engineering workshop*, pp. 507–512, IEEE, 2008.
- [30] K. R. Mabokela, M. J. Manamela, and M. Manaileng, "Modeling code-switching speech on under-resourced languages for language identification," in *Spoken Language Technologies for Under-Resourced Languages*, 2014.
- [31] R. Mabokela, "Phone clustering methods for multilingual language identification," in *In Computer Science & Information Technology (CS & IT) Conference*, pp. 285–298, 2020.
- [32] K. Shalini, H. B. Ganesh, M. A. Kumar, and K. P. Soman, "Sentiment analysis for code-mixed indian social media text with distributed representation," in *2018 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, pp. 1126–1131, 2018.
- [33] A. Balahur and M. Turchi, "Comparative experiments using supervised learning and machine translation for multilingual sentiment analysis," *Computer Speech & Language*, vol. 28, no. 1, pp. 56–75, 2014.
- [34] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1631–1642, 2013.
- [35] X. Chen, Y. Sun, B. Athiwaratkun, C. Cardie, and K. Weinberger, "Adversarial Deep Averaging Networks for Cross-Lingual Sentiment Classification," *Transactions of the Association for Computational Linguistics*, vol. 6, pp. 557–570, 12 2018.
- [36] S. Ruder, P. Ghaffari, and J. G. Breslin, "Insight-1 at semeval-2016 task 5: Deep learning for multilingual aspect-based sentiment analysis," *arXiv preprint arXiv:1609.02748*, 2016.
- [37] E. C. Chau, L. H. Lin, and N. A. Smith, "Parsing with multilingual bert, a small corpus, and a small treebank," *CoRR*, vol. abs/2009.14124, 2020.
- [38] S. Sagnika, A. Pattanaik, B. Mishra, and S. Meher, "A review on multilingual sentiment analysis by machine learning methods," *Journal of Engineering Science and Technology Review*, vol. 13, pp. 154–166, 04 2020.
- [39] N. A. S. Abdullah and N. I. A. Rusli, "Multilingual sentiment analysis: A systematic literature review," *pertanika journal of science and technology*, vol. 29, 2021.
- [40] Y. Xu, H. Cao, W. Du, and W. Wang, "A survey of cross-lingual sentiment analysis: Methodologies, models and evaluations," *Data Science and Engineering*, vol. 7, pp. 1–21, 06 2022.
- [41] D. Moher, A. Liberati, J. M. Tetzlaff, and D. G. Altman, "Preferred reporting items for systematic reviews and meta-analyses: the prisma statement," *International journal of surgery*, vol. 8 5, pp. 336–41, 2010.
- [42] A. Qazi, N. Hasan, O. Abayomi-Alli, G. Hardaker, R. Scherer, Y. Sarker, S. Paul, and J. Maitama, "Gender differences in information and communication technology use & skills: a systematic review and meta-analysis," *Education and Information Technologies*, vol. 27, p. 4225–4258, 04 2022.
- [43] A. Ghallab, A. Mohsen, Y. Ali, and C. W. Dawson, "Arabic sentiment analysis: A systematic literature review," *Appl. Comp. Intell. Soft Comput.*, vol. 2020, jan 2020.
- [44] Q. A. Xu, V. Chang, and C. Jayne, "A systematic review of social media-based sentiment analysis: Emerging trends and challenges," *Decision Analytics Journal*, vol. 3, p. 100073, 2022.
- [45] X. Wan, "Co-training for cross-lingual sentiment classification," in *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-volume 1*, pp. 235–243, Association for Computational Linguistics, 2009.



- [46] M. E. Basiri, A. R. N. Nilchi, and N. Ghassem-Aghaee, "A framework for sentiment analysis in persian," in *Open Transactions on Information Processing*, 2014.
- [47] F. Amiri, S. Scerri, and M. H. Khodashahi, "Lexicon-based sentiment analysis for persian text," in *Proceedings of Recent Advances in Natural Language Processing*, pp. 9–16, 2015.
- [48] F. Å. Nielsen, "A new ANEW: evaluation of a word list for sentiment analysis in microblogs," *CoRR*, vol. abs/1103.2903, 2011.
- [49] C. Hutto and E. Gilbert, "Vader: A parsimonious rule-based model for sentiment analysis of social media text," in *Proceedings of the 8th International Conference on Weblogs and Social Media, ICWSM 2014*, pp. 218–225, 01 2015.
- [50] M. A. Thelwall, K. Buckley, and G. Paltoglou, "Sentiment in twitter events," *Journal Association Information Sci. Technology*, vol. 62, pp. 406–418, 2011.
- [51] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [52] J. Angel, S. T. Aroyehun, A. Tamayo, and A. Gelbukh, "NLP-CIC at SemEval-2020 task 9: Analysing sentiment in code-switching language using a simple deep-learning classifier," in *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pp. 957–962, International Committee for Computational Linguistics, Dec. 2020.
- [53] O. Araque, I. Corcuera-Platas, J. F. Sánchez-Rada, and C. A. Iglesias, "Enhancing deep learning sentiment analysis with ensemble techniques in social applications," *Expert Systems with Applications*, vol. 77, pp. 236–246, 2017.
- [54] C. Banea, R. Mihalcea, and J. Wiebe, "Multilingual subjectivity: Are more languages better?," in *COLING*, 2010.
- [55] M. Saad, D. Langlois, and K. Smaïli, "Building and modelling multilingual subjective corpora," in *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, (Reykjavik, Iceland, Iceland), European Language Resources Association (ELRA), May 2014.
- [56] B. Pang and L. Lee, "A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts," *CoRR*, vol. cs.CL/0409058, 2004.
- [57] R. Mihalcea, C. Banea, and J. Wiebe, "Learning multilingual subjective language via cross-lingual projections," in *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, (Prague, Czech Republic), pp. 976–983, Association for Computational Linguistics, June 2007.
- [58] C. Banea, R. Mihalcea, J. Wiebe, and S. Hassan, "Multilingual subjectivity analysis using machine translation," in *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, (Honolulu, Hawaii), pp. 127–135, Association for Computational Linguistics, Oct. 2008.
- [59] C. Banea, R. Mihalcea, and J. Wiebe, "Sense-level subjectivity in a multilingual setting," *Computer Speech and Language*, vol. 28, pp. 7–19, 2014.
- [60] A. Ghafoor, I. A.S., D. S.M., K. Z., Abdullah, B. R., and W. M.A., "The impact of translating resource-rich datasets to low-resource languages through multi-lingual text processing," *IEEE Access*, vol. 9, p. 124478 – 124490, 2021. Cited by: 0; All Open Access, Gold Open Access, Green Open Access.
- [61] D. Vilares, M. A. Alonso, and C. Gómez-Rodríguez, "Supervised sentiment analysis in multilingual environments," *Information Processing & Management*, vol. 53, no. 3, pp. 595–607, 2017.
- [62] G. Shalunts, G. Backfried, and N. Commeignes, "The impact of machine translation on sentiment analysis," in *The Fifth International Conference on Data Analytics*, 2016.
- [63] M. S. Hajmohammadi, R. Ibrahim, and A. Selamat, "Bi-view semi-supervised active learning for cross-lingual sentiment classification," *Information processing & management*, vol. 50, no. 5, pp. 718–732, 2014.
- [64] X. Wan, "Using bilingual knowledge and ensemble techniques for unsupervised chinese sentiment analysis," in *2008 Conference on Empirical Methods in Natural Language Processing, EMNLP 2008, Proceedings of the Conference, Honolulu, Hawaii, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*, pp. 553–561, ACL, 2008.
- [65] A. Balahur and M. Turchi, "Improving sentiment analysis in twitter using multilingual machine translated data," in *Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP 2013*, pp. 49–55, 2013.
- [66] F. N. Ribeiro, M. Araújo, P. Gonçalves, M. A. Gonçalves, and F. Benevenuto, "Sentibench-a benchmark comparison of state-of-the-practice sentiment analysis methods," *EPJ Data Science*, vol. 5, no. 1, pp. 1–29, 2016.
- [67] M. Araújo, P. Gonçalves, M. Cha, and F. Benevenuto, "iFeel: A system that compares and combines sentiment analysis methods," in *Proceedings of the 23rd International Conference on World Wide Web*, pp. 75–78, 2014.
- [68] M. L. D. Araujo, J. P. Diniz, L. Bastos, E. Soares, M. Júnior, M. Ferreira, F. Ribeiro, and F. Benevenuto, "iFeel 2.0: A multilingual benchmarking system for sentence-level sentiment analysis," in *Tenth International AAAI Conference on Web and Social Media*, 2016.
- [69] J. Pan, G.-R. Xue, Y. Yu, and Y. Wang, "Cross-lingual sentiment classification via bi-view non-negative matrix tri-factorization," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 289–300, Springer, 2011.
- [70] X. Wan, "Bilingual co-training for sentiment classification of chinese product reviews," *Computational Linguistics*, vol. 37, no. 3, pp. 587–616, 2011.
- [71] B. Wei and C. Pal, "Cross lingual adaptation: an experiment on sentiment classifications," in *Proceedings of the ACL 2010 conference short papers*, pp. 258–262, Association for Computational Linguistics, 2010.
- [72] C. Tho, Y. Heryadi, L. Lukas, and A. Wibowo, "Code-mixed sentiment analysis of indonesian language and javanese language using lexicon based approach," *Journal of Physics: Conference Series*, vol. 1869, pp. 012–084, apr 2021.
- [73] L. Zhang and C. Chen, "Sentiment classification with convolutional neural networks: An experimental study on a large-scale chinese conversation corpus," in *2016 12th International Conference on Computational Intelligence and Security (CIS)*, pp. 165–169, IEEE, 2016.
- [74] M. Cliche, "Bb\_twt at semeval-2017 task 4: Twitter sentiment analysis with cnns and lstms," *arXiv preprint arXiv:1704.06125*, 2017.
- [75] H. Nguyen and M.-L. Nguyen, "A deep neural architecture for sentence-level sentiment classification in twitter social networking," in *International Conference of the Pacific Association for Computational Linguistics*, pp. 15–27, Springer, 2017.
- [76] L. Medrouk and A. Pappa, "Deep learning model for sentiment analysis in multi-lingual corpus," in *International Conference on Neural Information Processing*, pp. 205–212, Springer, 2017.
- [77] E. Boiy and M.-F. Moens, "A machine learning approach to sentiment analysis in multilingual web texts," *Information Retrieval Journal*, vol. 12, pp. 526–558, 10 2009.
- [78] X. Zhou, X. Wan, and J. Xiao, "Attention-based lstm network for cross-lingual sentiment classification," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 247–256, 01 2016.
- [79] J. Wehrmann, W. Becker, H. E. L. Cagnini, and R. C. Barros, "A character-based convolutional neural network for language-agnostic twitter sentiment analysis," in *2017 International Joint Conference on Neural Networks (IJCNN)*, pp. 2384–2391, 2017.
- [80] X. Dong and G. de Melo, "Cross-lingual propagation for deep sentiment analysis," in *AAAI*, 2018.
- [81] L. Medrouk and A. Pappa, "Do deep networks really need complex modules for multilingual sentiment polarity detection and domain classification?," in *2018 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–6, 2018.
- [82] M. Giatsoyglou, M. G. Vozalis, K. Diamantaras, A. Vakali, G. Sarigiannidis, and K. C. Chatzisavvas, "Sentiment analysis leveraging emotions and word embeddings," *Expert Systems with Applications*, vol. 69, pp. 214–224, 2017.
- [83] Y. K. Lal, V. Kumar, M. Dhar, M. Shrivastava, and P. Koehn, "De-mixing sentiment from code-mixed text," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pp. 371–377, 2019.
- [84] A. Gupta, S. Menghani, S. K. Rallabandi, and A. W. Black, "Unconscious self-training for sentiment analysis of code-linked data," *ArXiv*, vol. abs / 2103.14797, 2021.
- [85] N. Sabri, A. Edalat, and B. Bahrak, "Sentiment Analysis of Persian-English Code-mixed Texts," *2021 26th International Computer Conference, Computer Society of Iran (CSICC)*, pp. 1–4, 2021.
- [86] B. R. Chakravarthi, R. Priyadarshini, V. Muralidaran, N. Jose, S. Suryawanshi, E. Sherly, and J. P. McCrae, "Dravidiancodemix: sentiment analysis and offensive language identification dataset for dravidian



- languages in code-mixed text,” *Language Resources and Evaluation*, vol. 56, pp. 765–806, 2022.
- [87] K. S. B. S. Saikrishna and C. N. Subalalitha, “Sentiment analysis on telugu–english code-mixed data,” in *Intelligent Data Engineering and Analytics* (S. C. Satapathy, P. Peer, J. Tang, V. Bhateja, and A. Ghosh, eds.), (Singapore), pp. 151–163, Springer Nature Singapore, 2022.
- [88] S. S. V. Kusampudi, P. Sathineni, and R. Mamidi, “Sentiment analysis in code-mixed Telugu-English text with unsupervised data normalization,” in *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, (Held Online), pp. 753–760, INCOMA Ltd., Sept. 2021.
- [89] S. Thara and P. Poornachandran, “Social media text analytics of malayalam–english code-mixed using deep learning,” *Journal of Big Data*, vol. 9, 2022.
- [90] C. Dos Santos and M. Gatti, “Deep convolutional neural networks for sentiment analysis of short texts,” in *Proceedings of COLING 2014*, the 25th International Conference on Computational Linguistics: Technical Papers, pp. 69–78, 2014.
- [91] O. Isroy and C. Cardie, “Opinion mining with deep recurrent neural networks,” in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 720–728, 2014.
- [92] K. S. Tai, R. Socher, and C. D. Manning, “Improved semantic representations from tree-structured long short-term memory networks,” *arXiv preprint arXiv:1503.00075*, 2015.
- [93] C. Baziotis, N. Pelekis, and C. Doulkeridis, “Datastories at semeval-2017 task 4: Deep lstm with attention for message-level and topic-based sentiment analysis,” in *Proceedings of the 11th international workshop on semantic evaluation (SemEval-2017)*, pp. 747–754, 2017.
- [94] S. R. Shah and A. Kaushik, “Sentiment analysis on indian indigenous languages: A review on multilingual opinion mining,” *arXiv preprint arXiv:1911.12848*, 2019.
- [95] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, “Deep contextualized word representations,” in *NAACL*, 2018.
- [96] T. Tomihira, A. Otsuka, A. Yamashita, and T. Satoh, “Multilingual emoji prediction using bert for sentiment analysis,” *International Journal of Web Information Systems*, vol. 16, pp. 265–280, 09 2020.
- [97] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov, “Unsupervised cross-lingual representation learning at scale,” *CoRR*, vol. abs/1911.02116, 2019.
- [98] B. R. Chakravarthi, V. Muralidaran, R. Priyadarshini, and J. P. McCrae, “Corpus creation for sentiment analysis in code-mixed Tamil-English text,” in *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)*, (Marseille, France), pp. 202–210, European Language Resources association, May 2020.
- [99] J. O. Alabi, D. Ifeoluwa Adelani, M. Mosbach, and D. Klakow, “Multilingual Language Model Adaptive Fine-Tuning: A Study on African Languages,” *arXiv e-prints*, p. arXiv:2204.06487, Apr. 2022.
- [100] I. Mozetič, M. Grčar, and J. Smailović, “Multilingual twitter sentiment classification: The role of human annotators,” *PLoS ONE*, vol. 11, no. 5, 2016. Cited by: 100; All Open Access, Gold Open Access, Green Open Access.
- [101] J. Wehrmann, W. E. Becker, and R. C. Barros, “A multi-task neural network for multilingual sentiment classification and language detection on twitter,” *Proceedings of the ACM Symposium on Applied Computing*, p. 1805–1812, 2018. Cited by: 14.
- [102] X. Dong and G. de Melo, “A robust self-learning framework for cross-lingual text classification,” *EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing*, *Proceedings of the Conference*, p. 6306–6310, 2019. Cited by: 17.
- [103] J. Kranjc, J. Smailovic, V. Podpecan, M. Grcar, M. Znidari, and N. Lavrac, “Active learning for sentiment analysis on data streams: Methodology and workflow implementation in the clowdflows platform,” *Information Process Management*, vol. 51, pp. 187–203, 2015.
- [104] A. Balahur and J. M. Perea-Ortega, “Sentiment analysis system adaptation for multilingual processing: The case of tweets,” *Information Processing & Management*, vol. 51, no. 4, pp. 547–556, 2015.
- [105] G. Liu, X. Huang, X. Liu, and A. Yang, “A novel aspect-based sentiment analysis network model based on multilingual hierarchy in online social network,” *The Computer Journal*, vol. 63, pp. 410–424, 2020.
- [106] S. Hassan Muhammad, D. Ifeoluwa Adelani, S. Ruder, I. Said Ahmad, I. Abdulmumin, B. Shehu Bello, M. Choudhury, C. Chinenye Emezue, S. Salahudeen Abdullahi, A. Arenu, A. George, and P. Brazdil, “Nai-jaSenti: A Nigerian Twitter Sentiment Corpus for Multilingual Sentiment Analysis,” *arXiv e-prints*, p. arXiv:2201.08277, Jan. 2022.
- [107] K. Dashtipour, M. Gogate, J. Li, F. Jiang, B. Kong, and A. Hussain, “A Hybrid Persian Sentiment Analysis Framework: Integrating Dependency Grammar Based Rules and Deep Neural Networks,” *Neurocomputing*, vol. 380, pp. 1–10, 2020.
- [108] K. Rakshitha, H. M. Ramalingam, M. Pavithra, A. H. D., and M. Hegde, “Sentimental analysis of indian regional languages on social media,” *Global Transitions Proceedings*, vol. 2, no. 2, pp. 414–420, 2021. *International Conference on Computing System and its Applications (ICCSA-2021)*.
- [109] S. S. Mukku and R. Mamidi, “ACTSA: Annotated corpus for Telugu sentiment analysis,” in *Proceedings of the First Workshop on Building Linguistically Generalizable NLP Systems*, (Copenhagen, Denmark), pp. 54–58, Association for Computational Linguistics, Sept. 2017.
- [110] M. E. Basiri and A. Kabiri, “Sentence-level sentiment analysis in persian,” *2017 3rd International Conference on Pattern Recognition and Image Analysis (IPRIA)*, pp. 84–89, 2017.
- [111] B. Roshanfekr, S. Khadivi, and M. Rahmati, “Sentiment analysis using deep learning on persian texts,” in *2017 Iranian Conference on Electrical Engineering (ICEE)*, pp. 1503–1508, 2017.
- [112] R. Ghasemi, S. A. A. Asli, and S. Montazi, “Deep persian sentiment analysis: Cross-lingual training for low-resource languages,” *Journal of Information Science*, vol. 48, no. 4, pp. 449–462, 2022.
- [113] S. Sazed, “Improving sentiment classification in low-resource bengali language utilizing cross-lingual self-supervised learning,” in *Natural Language Processing and Information Systems: 26th International Conference on Applications of Natural Language to Information Systems, NLDB 2021, Saarbrücken, Germany, June 23–25, 2021, Proceedings*, (Berlin, Heidelberg), p. 218–230, Springer-Verlag, 2021.
- [114] Gati Lothar Martin and Medard Edmund Mswahili and Young-Seob Jeong, “Sentiment Classification in Swahili Language Using Multilingual BERT,” *African NLP Workshop, EACL 2021*, vol. abs/2104.09006, 2021.
- [115] V. Marivate, T. Sefara, V. Chabalala, K. Makhaya, T. B. Mokgonyane, R. Mokoena, and A. Modupe, “Low resource language dataset creation, curation and classification: Setswana and sepedi - extended abstract,” *CoRR*, vol. abs/2004.13842, 2020.
- [116] K. Ogueji, Y. Zhu, and J. Lin, “Small data? no problem! exploring the viability of pretrained multilingual language models for low-resourced languages,” in *Proceedings of the 1st Workshop on Multilingual Representation Learning*, (Punta Cana, Dominican Republic), pp. 116–126, Association for Computational Linguistics, Nov. 2021.
- [117] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” 2019.
- [118] H. Xu, B. Van Durme, and K. Murray, “Bert, mbert, or bibert? a study on contextualized embeddings for neural machine translation,” *arXiv*, 2021.
- [119] L. Ein-Dor, A. Halfon, A. Gera, E. Shnarch, L. Dankin, L. Choshen, M. Danilevsky, R. Aharonov, Y. Katz, and N. Slonim, “Active Learning for BERT: An Empirical Study,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, (Online), pp. 7949–7962, Association for Computational Linguistics, Nov. 2020.



KOENA RONNY MABOKELA is the Head of the Technopreneurship Centre within the School of Consumer Intelligence and Information Systems. He is also a lecturer in the Department of Applied Information Systems at the University of Johannesburg. He earned his MSc degree in Computer Science from the University of Limpopo, where he obtained his undergraduate degrees in Computer Science and Mathematics. Prior to joining the University of Johannesburg in 2019, he worked for over 5 years in the telecommunications sector, thus acquiring industry experience. He has presented his research work on numerous platforms nationally and internationally and has a keen research interest in NLP, speech technologies, multilingual sentiment analysis for under-resourced languages, among other areas. He is currently pursuing his PhD studies at the University of the Witwatersrand. He is a member of the South African Institute for Computer Scientists & Information Technologists and also serves on various boards.



TURGAY CELIK is a Professor of Digital Transformation Engineering in the School of Electrical and Information Engineering and Director of the Wits Institute of Data Science at the University of the Witwatersrand. He received his second PhD degree from the University of Warwick, UK in 2011. He actively reviews and publishes in various international journals and conferences. Professor Celik is an Associate Editor of IEEE Access, IET Electronics Letters, IEEE Geoscience and Remote Sensing Letters, IEEE JSTARS and Springer Signal, Image and Video Processing Journal. His research interests are in the areas of signal and image processing, computer vision, machine learning, artificial intelligence, robotics, data science and remote sensing.



MPHO RABORIFE. is an Associate Professor and Deputy Director at the Institute for Intelligent Systems, University of Johannesburg. An NRF-rated researcher, she specialises in theoretical computer science and computational phonology. She has a PhD in Computer Science from the University of the Witwatersrand. Her achievements include the L'OREAL-UNESCO sub-Saharan regional Women in Science Fellow, the M&G 200 Award and a DST Women in Science alumni. In addition, Prof. Raborife received numerous scholarships and bursaries for her undergraduate through to her post-PhD studies. She has successfully supervised master's and PhD students. She has worked on numerous projects (multidisciplinary) with project partners across three continents. Prof. Raborife also contributes to scientific citizenship by being a member of scientific bodies.

• • •