

# Towards Arabic Multimodal Dataset for Sentiment Analysis

Abdelhamid Haouhat

*Lab. d'Informatique et de Mathématiques*  
*Université Amar Telidji, Laghouat, Algérie*  
a.haouhat@lagh-univ.dz

Slimane Bellaouar

*Dept. of Mathematics and Computer Science*  
*Lab. des Mathématiques et Sciences Appliquées (LMSA)*  
*Université de Ghardaia, Algérie*  
bellaouar.slimane@univ-ghardaia.dz

Attia Nehar

*Université Ziane Achour - Djelfa, Algérie*  
*Lab. d'Informatique et de Mathématiques (LIM)*  
*Université Amar Telidji, Laghouat, Algérie*  
neharattia@univ-djelfa.dz

Hadda Cherroun

*Lab. d'Informatique et de Mathématiques*  
*Université Amar Telidji, Laghouat, Algérie*  
hadda\_cherroun@lagh-univ.dz

**Abstract**—Multimodal Sentiment Analysis (MSA) has recently become a centric research direction for many real-world applications. This proliferation is due to the fact that opinions are central to almost all human activities and are key influencers of our behaviors. In addition, the recent deployment of Deep Learning-based (DL) models has proven their high efficiency for a wide range of Western languages. In contrast, Arabic DL-based multimodal sentiment analysis (MSA) is still in its infantile stage due, mainly, to the lack of standard datasets. In this paper, our investigation is twofold. First, we design a pipeline that helps building our Arabic Multimodal dataset leveraging both state-of-the-art transformers and feature extraction tools within word alignment techniques. Thereafter, we validate our dataset using state-of-the-art transformer-based model dealing with multimodality. Despite the small size of the outcome dataset, experiments show that Arabic multimodality is very promising.

**Index Terms**—Sentiment Analysis, Multimodal Learning, Transformers, Arabic Multimodal Dataset.

## I. INTRODUCTION

The field of Multimodal Machine Learning (MML) has been growing rapidly over the past few decades, driven by the increasing availability of multimodal data and the need for more sophisticated and effective machine learning models. MML entails integrating and modeling data from various modalities (text, audio, image, video).

Multimodal Sentiment Analysis (MSA), for instance, is an important and growing area of MML that aims to automatically determine the sentiment expressed in various modalities [1]. Early works in MSA deal with feature extraction and fusion processes in straightforward ways using standard machine learning algorithms [2], [3]. Over time, more complex methods, mainly deep learning models, were developed, such as CNN [4], RNN and its architectural variants [5], and Multimodal Multi-Utterance models [6].

Arabic MSA is a promising area for academic research and practical applications due to the widespread use of the Arabic language and the increasing popularity of multimedia content. Furthermore, Arabic MSA is challenging due, on the one hand, to the complexity and the richness of the Arabic language, and on the other hand, to the significant cultural and linguistic variety of the Arab world. Therefore, it is still in its infancy [7].

Despite these challenges, there has been some limited work in Arabic MSA that has shown promising results [8], [9]. Hence, there is still much room for improvement in accuracy, efficiency, flexibility, and ability to handle diverse modalities.

This paper has two main investigations:

- 1) First, we design a pipeline that facilitates the construction of a novel Arabic multimodal dataset. We accomplish this by leveraging state-of-the-art transformers and feature extraction tools alongside word alignment methods.
- 2) To assess the effectiveness of our Arabic multimodal dataset, we employ cutting-edge transformer models that are intended to handle multimodality.

The remainder of this paper is structured as follows. Section II introduces some basic concepts concerning SA and MML required to understand the rest of the paper. Section III provides an overview of previous research on English and Arabic MSA. In Section IV, we describe the proposed methodology for multimodal dataset collection. We also present the models that were used to evaluate the designed dataset. Section V deals with experiments and interpretation of the empirical findings. Finally, Section VI outlines the conclusions and future works.

## II. PRELIMINARIES

Before diving into the details of our approach, we start with the terminologies and background concepts that concern Sentiment Analysis (SA) and Multimodal Machine Learning

(MML) elements: data representation, modality fusion methods, alignment, and pre-trained models.

#### A. Sentiment Analysis

Sentiment Analysis (SA), also referred to as opinion analysis, is the process of obtaining and examining the views, ideas, and perceptions of the public on a wide range of topics, products, subjects, and services. Corporations, governments, and people may all benefit from public opinion when gathering data and making choices based on it. [10].

Let us mention that the words *emotion* and *sentiment* are usually used interchangeably in daily life. While they are two different concepts. *Emotion* is defined as a complex psychological state. There are six basic emotions, i.e., happiness, sadness, anger, fear, surprise, and disgust. This list is enriched by adding emotions such as pride, excitement, embarrassment, contempt, and shame. On the other hand, *Sentiment* describes a mental attitude that is founded on emotion [11]. Positive, neutral, and negative are the three fundamental polarities. The SA also makes reference to a polarity categorization. There are several methods for performing SA, including rule-based methods, machine learning-based methods, and hybrid approaches. Some popular machine learning-based methods include Naive Bayes, Support Vector Machines (SVM), and Deep Learning-based (DL) models. However, DL-based models have proved their efficiency as SOTA approaches.

Human natural perception refers to our ability to perceive and understand information from multiple modalities in a seamless and integrated way, such as seeing a picture and hearing a sound simultaneously to understand a concept. Multimodal SA aims to replicate this natural perception by combining information from multiple modalities (text, audio and image/video, and more) to improve the accuracy and efficiency of learning systems.

#### B. Multimodal Machine Learning

Multimodal Machine Learning (MML) involves integrating and modeling multiple communicative modalities, Such as linguistic (text), acoustic (sound), and visual messages(image and video) of data, from a variety of diverse and interconnected sources [12]. By leveraging the strengths of different modalities, multimodal learning can help overcome the limitations of individual modalities and enhance overall learning performance.

Liang et al. proposed a taxonomy of six core features in MML: Modality representation, alignment, reasoning, generation, transference, and quantification [13] that are understudied in conventional unimodal machine learning. Considering their importance for our study, we focus on two features: i) Representation: where we focus mainly on which adequate representation is suitable for each modality and then how and when to fuse and integrate information from two or more modalities, effectively reducing the number of separate representations. ii) Alignment: Alignment between modalities is also challenging and involves identifying connections between modality elements.

1) *Fusion Methods*: Basically, We have two main methods to make a fusion of modalities. The first is *Early Fusion*, which happens when we mix the modalities before making decisions with concatenation, summation, or cross-attention mechanism. While the second is *Late Fusion* method which makes a prediction based on each modality alone and then combines decisions to get a final prediction [14]. In our approach, We deploy an early fuse approach.

2) *Pre-trained Models*: The deployment of semantic and Deep Learning based approaches leads us generally to use some pre-trained models such as GloVe multimodal bi-transformer model (MMBT) models [15], CLIP [16], BERT [17] and AraBERT [18]. For the purpose of this paper, we have deployed AraBERT.

BERT, which represents the basis of ArabBert, is a Bidirectional Encoder Representation from Transformers developed by Google [16]. BERT large encompasses 24 encoders with 16 bidirectional self-attention heads trained from unlabeled data extracted from the BooksCorpus and English Wikipedia.

AraBERT [18] is a pre-trained BERT transformer built for Arabic NLP tasks. It is trained on 70 million sentences, corresponding to 24GB of Arabic text. AraBERT uses the same configuration as Bert. It has 12 encoder blocks, 768 hidden dimensions, 12 attention heads, 512 maximum sequence length, and 110M parameters.

### III. RELATED WORK

This section reviews the relevant studies in the field of multimodal sentiment analysis (MSA), including traditional, machine learning, and deep learning approaches for both English and Arabic languages.

#### A. English Multimodal Sentiment Analysis

The study of Zadeh [2] is considered one of the pioneering works in the field of MSA. It is the first work to tackle the challenge of tri-modal (visual, audio, and textual features) sentiment analysis. The author creates a dataset of 47 videos from YouTube. Each input in the dataset was annotated with either a positive, negative, or neutral label. Moreover, the paper identifies specific subset of audio-visual features relevant to sentiment analysis and presents some instructions for integrating these features. In experiment, author uses the Hidden Markov Model (HMM) classifier. The findings demonstrate the promise of MSA despite the small size of the dataset and the straightforward text analysis method.

Poria et al. [3] propose a novel methodology for performing MSA based on sentiment extraction from online videos. They deploy the dataset initially created by [2]. The authors discuss features extracting process from various modalities (text, audio, and visual). These features are fused by incorporating different techniques (feature-level and decision-level). The authors used multiple supervised machine learning classifiers (Support Vector Machine (SVM), Extreme Learning Machine (ELM), Naive Bayes (NB), and Neural Networks) to validate their approach. Finally, a comparative study was carried out on the selected dataset, revealing that their proposed MSA system

outperforms the current state-of-the-art systems. The best performance was achieved with Extreme Learning Machine (ELM) method.

The study in [4] provides a detailed review that explores the applicability, challenges, and issues for textual, visual, and MSA using CNNs. Several enhancements have been proposed, such as combining CNN and long short-term memory (LSTM) techniques.

Tembhurne and Diwan [5] study the role of sequential deep neural networks in MSA. They thoroughly examined applicability, problems, issues, and methodologies for textual, visual, and MSA based on RNN and its architectural variants.

Recently, Abdu et al. [6] draw up a survey on MSA using deep learning. They have categorized 35 cutting-edge models, recently suggested for the video sentiment analysis field, into eight categories, based on the specific architecture employed in each model. After a detailed examination of the results, authors conclude that the *Multimodal Multi-Utterance* based architecture is the most powerful in the task of MSA.

Before concluding this section, we point out that the two transformer-based models known as the Multimodal transformer (Mult) [19] and LS-LSTM [20] that we have used to evaluate our dataset are described in Section IV-B.

### B. Arabic Multimodal Sentiment Analysis

In contrast to the MSA studies made for the English language, the one for the Arabic language encompasses a limited number of works.

Najadat and Abushaqra [8] aim to address the issue of MSA for Arabic. They start by building their dataset from YouTube. They extract different features (linguistic, audio, and visual) from the collected videos. Also, they augment data using *Weka* re-sampling option. For training and testing purposes, the authors annotate their dataset by positive, negative, and neutral polarities. In the experiment stage, the authors use different machine learning classifiers (Decision Trees, Support Vector Machine (SVM), k-Nearest Neighbor (KNN), Naive Bayes (NB), and Neural Networks). Obtained results reveal that the Neural Network classifier performs best when using only the audio modality. However, obtained results can be enhanced by feeding the dataset with more features.

In their paper [9], Alqarafi et al. try to tackle the problem of sentiment analysis in online opinion videos for modern standard Arabic. They begin by constructing their Arabic Multimodal Dataset (AMMD) from 40 different YouTube videos. First, they used the extracted features (text, video) to feed their dataset. After that, they add metadata about the videos, including audio, transcription, visual motions, and sentiment polarities. Authors use, to conduct experiments, the Support Vector Machine (SVM) classifier. Despite the limited size of the dataset, the experimental results demonstrate the validity of the constructed dataset. Additionally, the results indicate that for several sentiment analysis tasks, including subjectivity and polarity classifications, the fusion of different features (utterance, visual) improves the performance compared to using utterance features alone.

## IV. METHODOLOGY

In order to tackle the problem of Arabic Multimodal Sentiment Analysis *AMSA*, in our study, we conduct two main investigations.

First, we design a pipeline that eases building Multimodal dataset for sentiment analysis that respects dataset collection engineering and harnesses transformers and SOTA feature extraction tools( IV-A).

Second, we assess our built dataset using SOTA transformer-based models that deal with multimodality. The transformers are chosen to leverage inherent semantics while multimodality is deployed to improve the sentiment learner (Section IV-B).

### A. Multimodal Dataset collection Methodology

As mentioned above, our targeted multimodal dataset for Arabic Sentiment Analysis, involving dataset collection engineering, aims to leverage transformers and SOTA feature extraction techniques. Indeed, we have proposed this generic pipe to build the Multimodal dataset:

- 1) Data Inventory, Collection, and Preprocessing.
- 2) Annotation.
- 3) Data representation.

Our methodology is inspired by both MOSEI [21] and CMU-MOSI [22] dataset-building processes while taking into account Arabic specificities.

1) *Data Inventory, Collection, and Preprocessing*: We rely mainly on videos on Youtube and Social Media platforms that include various information about the videos, such as audio, visual gestures, metadata, and probably transcripts. First, we have to identify sources that guarantee encompassing subjective information passages such as video-Bloggers, political analysts, and influencers channels. We also rely on some Tv' Talk Shows. To get a large size of video we automatically draw some search lists and API to ease scraping information and its related metadata.

For our targeted NLP task, pre-processing the collected video include objective segment removal, speech extraction, text extraction, and video/audio segmentation. All these processing are semi-automatic, using open source tools and quiet manual intervention.

2) *Annotation*: Annotating the polarity of video segments, as well as their associated text and speech, is the most challenging and resource-intensive task in our work. It requires a significant amount of time and resources to accurately annotate. We have opted to rely on manual crowdsourcing and manual annotation through a homemade platform. A guideline is devised in order to uniform the annotation. In this step, we use the classic polarities  $[-1, 0, 1]$  for negative, neutral, and positive sentiments, respectively. The annotation evaluation is performed through the standard automatic Inter-Annotator Agreement method.

3) *Data Representation*: The three targeted modalities are represented in such a way that they exhibit more information on the inherent sentiment.

### Text

Concerning the text, one can use either word-embedding or pretrained transformers. However, the latter allows learning contextual relationships between words in a sentence through a bidirectional attention mechanism. This means that we represent words taking into account both the left and right context of each word in a sentence, giving it a more comprehensive understanding of the semantic meaning and providing more accurate representations of textual modality.

### Visual Features

The combination of body gestures and facial features can convey a more nuanced range of sentiments and emotions.

Body gestures refer to physical characteristics of a person's body, such as Open/crossed arms, nodding, shaking head, and Shrugging shoulders. Facial features refer to the characteristics of a face, such as facial expressions and movements, that are used to represent emotions and sentiments of an individual. These features include commonly: smiling, frowning, raised eyebrows, squinted eyes, lip biting, tears, and blushing.

For this version of our pipe, we have opted to rely on facial features as they are the most commonly used features. In addition, they are also easier to capture from videos compared with body gestures. In fact, facial features can be extracted using computer vision techniques. It may include measurements of facial landmarks, facial action units, and head movements. The main descriptors we extract are:

- Action Units (AUs): these are facial muscle movements that are associated with various facial expressions such as brow raise, lip stretch, and eye closure.
- Head pose: it estimates the orientation of the head in three dimensions, including pitch, yaw, and roll.
- Eye gaze: it captures the direction of the eye gaze, including the location of the gaze and the direction of the gaze vector.
- Facial Action Coding System (FACS): FACS is a system that describes facial expressions based on AUs.
- Facial symmetry: it informs about the symmetry of the face by comparing the left and right sides of the face.

### Acoustic Features

The speech extracted from the video can be characterized at different levels acoustic, phonology, or prosody. The acoustic features are the most effective ones as they are language-independent since they rely on the physical features of the signal. However, also prosody features are essential as they capture features related to the emotion related to the speaker's speech. These features are commonly used in speech recognition, speaker identification, and sentiment and emotion recognition systems. In our study, we extract those main features:

- Mel-Frequency cepstral coefficients (MFCCs): A set of coefficients that represent the spectral envelope of a speech *signal*.
- Prosody: These include fundamental frequency (F0), speaking rate, and energy.
- Voice quality: These include jitter, shimmer, harmonic-to-noise ratio (HNR), and glottal waveform features.

- Emotion-related: These include pitch slope, pitch variance, and various modulation features.
- Spectral: These include spectral centroid, spectral flux, and spectral roll-off.
- Formant features: These include the first three formants, which are resonant frequencies of the vocal tract.
- Timing features: These include various measures of speech timing, such as pause duration and speech rate.

4) *Alignment Techniques*: One of the crucial aspects of the MML is the alignment of multimodal data, which involves synchronizing the different modalities, such as text, audio, and video. This is typically done by aligning the timestamps of each modality and mapping them onto a common timeline.

Considering text as an important modality, we use, in our work, two stages in achieving this alignment. In the first stage, we perform Text and Audio alignment, where data are aligned at the word level. In the second stage, we perform video and text alignment. Thus, we get a global alignment over the text common modality. For both stages, we use forced alignment techniques.

### Forced Alignment Text-Audio

Within this alignment, a transcript is synchronized with an audio recording by mapping each speech segment to its corresponding words. This process of forced alignment typically involves breaking down the audio and transcript into smaller segments and using algorithms to compare the speech and text segments to determine their correspondence. The algorithms consider various factors, such as speech timing, pronunciation of words, and speech sounds. After the forced alignment process, we get a time-stamped representation of speech.

### Pivot-Based Multimodal Alignment

For effective handling of multimodal time series data featuring multiple views at different frequencies, it is crucial to align them to a designated "pivot" modality, which is typically done through textual modality. This involves grouping feature vectors from other modalities into bins based on the timestamps of the pivot modality and then applying a specific processing function, known as the "collapse function", to each bin. This function, often a pooling function, merges multiple feature vectors from another modality into a single vector, resulting in sequences of equal lengths across all modalities (matching the length of the pivot modality) in all-time series.

### B. Models

In this section, we provide a detailed explanation of the selected models used to validate our dataset. Our explanation includes a discussion of crucial multimodal learning (MML) techniques used in these models, such as fusion, modeling, and alignments.

The initial state-of-the-art model, known as the Multimodal transformer (Mult) [19], is a transformer-based model that deploys an attention mechanism. It allows each element of the input sequence  $X_i$  to attend to all the other elements, resulting in a new weighted sequence  $\hat{X}_i$ . This process is referred to as

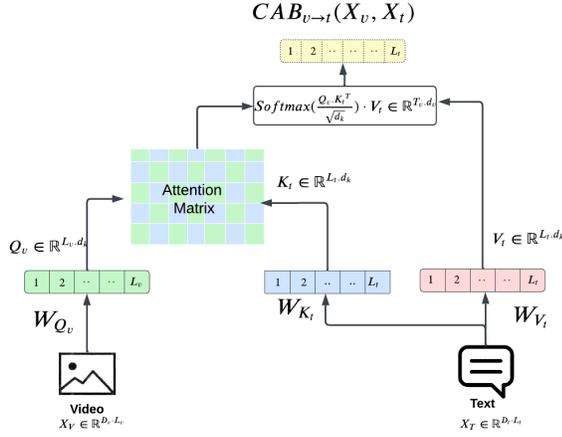


Fig. 1. Cross-modal with CAB lead model to integrate input modalities. In this illustration, CAB combines video and text information ( $X_v$  and  $X_t$ ) through an attention mechanism.

self-attention, as it enables the elements to focus on the most relevant ones,  $i$  representing the modality among {text, video, audio}.

Mult integrates these  $\hat{X}_t$ ,  $\hat{X}_a$ , and  $\hat{X}_v$  by utilizing a pairwise feed-forward approach. This is achieved through the implementation of a deep Cross Attention Block (CAB), which is defined below, where each pair-modality is fed into two CABs that alternate between Query, Key, and Values matrices. Before making predictions, the Multimodal Transformer model concatenates the  $Z_{i \rightarrow j}$  and  $Z_{k \rightarrow j}$  matrices obtained from the CABs, where  $i$ ,  $j$ , and  $k$  represent distinct modalities. This process results in the  $Z_i$  matrix. As previously mentioned, the prediction classifier is applied to the concatenated or summed  $\hat{Y}$ , which is composed of the  $Z_t, Z_a$ , and  $Z_v$  matrices where  $t$ ,  $a$ , and  $v$  are referred to text, audio, and video respectively.

**Cross Attention Block (CAB):** In order to implement the fusion from multiple modalities, Cross-modal allows the model to focus on relevant information from each modality and weigh the contribution of different modalities in the final prediction. We can express CAB in the mathematical formulation below. Let us consider an MML model with two modalities:  $X, Y$ . The query, key, and value matrices for each modality can be represented as follows:

$$Q_X = W_X^Q X, \quad K_X = W_X^K X, \quad V_X = W_X^V X \quad (1)$$

$$CAB_{x \rightarrow y} = \text{Softmax}\left(\frac{Q_x \cdot K_y^T}{\sqrt{d_k}}\right) \cdot V_y = Z_{x \rightarrow y}, \quad (2)$$

where  $W_X^Q, W_X^K, W_X^V, W_Y^Q, W_Y^K, W_Y^V$  are the weight matrices for the query, key, and value computations for each modality. and  $Q, K, V \in \mathbb{R}^{l \times d}$ ,  $d$  is the modality dimension,  $l$  is the length of input token  $X$ .

Fig. 1 illustrates the cross attention block. The main purpose of using these two sub-layers before and after CAB layers is to make our model focus only on dependent features.

In light of Equations (1) and (2) described above, we can express for each modality latent representation  $Z_i$  by Equation (3).

$$Z_x = [CAB_{y \rightarrow x}; CAB_{k \rightarrow x}] \quad (3)$$

where  $x, y$ , and  $k$  are all possible modalities. The output  $\hat{Y}$  of Mult model by Equation (4) as follows:

$$\hat{Y} = \sum [Z_t; Z_a; Z_v] \quad (4)$$

In another hand, we select another deep learning model using another fusing approach. This model consists of three separate Long Short-Term Memory (LSTM) networks [20], one for each modality (text (t), visual (v), and acoustic (a)). Here we use a late-fusion where different modalities are processed separately to obtain their respective feature representations, and then these features are combined using a fusion mechanism to make the final prediction. In the first stage, the model takes as input three modalities  $X_{\{t,a,v\}}$  and extracts features  $h1_i, h2_i$  from each of them using the corresponding two LSTMs with a normalization layer between them as shown in equations below.

$$o_t = \delta(W_o \times X_{\{t/a/v\}} + h_{t-1} + b_o)$$

$$i_t = \sigma(W_i[X_{\{t/a/v\}}, h_{t-1}] + b_i)$$

$$g_t = \tanh(W_c[X_{\{t/a/v\}}, h_{t-1}] + b_c)$$

$$O_{lstm1}(X) = h1_t = o_t \odot \tanh(c_t)$$

$$O_{lstm2}(h1_t) = h2_t = o_t \odot \tanh(c_t)$$

These features are then concatenated and normalized, before being fed into a fully connected layer with a ReLU activation function and dropout regularization. Finally, the output is generated using another fully connected layer defined as:

$$Output\_model = O_{lstm2}(Norm\_Layer(O_{lstm1}(X_{\{t,a,v\}})))$$

$$Output\_model = Concat[h2_t; h1_t; h2_a; h1_a; h2_v; h1_v]$$

Where  $x_t$  represents the input at time step  $t$ ,  $h_{t-1}$  represents the hidden state at the previous time step,  $i_t, f_t$ , and  $o_t$  represent the input, forget, and output gates at time step  $t$ , respectively. The symbol  $\odot$  denotes element-wise multiplication, and  $\sigma$  and  $\tanh$  are the sigmoid and hyperbolic tangent activation functions, respectively.

## V. EXPERIMENTS

In this section, first, we start by describing the details related to the implementation of our proposed Arabic Multimodal dataset pipe as well as the description of the collected dataset. In the second step, we empirically evaluate our built dataset through both Mult and LS-LSTM models.

### A. Data collection

Following the above-designed pipe, our dataset is gathered from video-blogging’ videos on mainly YouTube and some other social media platforms. The videos are retrieved and scraped automatically using a predefined list of keywords. In fact, by means of this latter, we ensure the existence of subjective information related to Arabic content.

All videos have been checked manually to keep the most convenient ones for our study. Then using a homemade collaborative front-end tool <sup>1</sup>, We segmented each video by placing *start* and *end* markers so that each video segment encompasses one subjective information. Then we extracted from each video segment its related Arabic transcription and speech. For those purposes, we use both Klaam tool <sup>2</sup> and *Almufaragh* tool <sup>3</sup> for Arabic speech recognition. The automatically extracted transcripts are also checked manually to avoid and fix any transcription errors. The forced and pivot alignments are performed on the fly thanks to Audacity tool <sup>4</sup>.

Let us mention that word alignment is a challenging task. The quality of alignment can be negatively affected when dealing with speeches in which words are not fully enunciated by the speaker.

Concerning the annotation process, each segment is labeled by 5 in lab annotators. A guideline is designed to reach more similar annotations. The Inter Agreement Annotator method is applied to assign a final label.

The resulting three modalities (Video, text, and audio) are then preprocessed to exhibit more information about their inherent sentiment, as described below.

Our method for extracting word vectors padded to max length from these transcripts is based on AraBERT [18] transformer. In fact, it is a BERT-based model that allows learning contextual relationships between words in a sentence through bidirectional attention mechanisms. That means that we represent words taking into account both the left and right context of each word in a sentence, giving it a more comprehensive understanding of the semantic meaning and providing more accurate representations of textual modality. Our text embeddings are in 768 dimensional vector.

Concerning the visual features, we opted for the facial features. Thanks to OpenFace toolkit [23], we extracted 45 facial features belonging to those described above.

The acoustic features are extracted using OpenSmile tool <sup>5</sup>. We extracted 52 features described previously.

Table I reports more details on the built dataset.

After that, we construct data formatted as a dictionary of multiple computational sequences using CMU-multimodal SDK [22].

A sample of our formatted dataset is available in <sup>6</sup>

TABLE I  
DATASET DETAILS.

Number of videos	60
Total number of segments	540
Total number of subjective segments	318
Number of unique words	2485
Total videos time	02h : 47min : 27s
Average length of segments	17.46 seconds
Number of Positive segment	130
Number of Negative segment	129
Number of Neutral segment	59
Text embedding dimensions	768
Visual Feature dimensions	45
Acoustic Feature	52
Number of speakers	23

### B. Results and Discussion

The deployed models are measured through three metrics: Accuracy, F1 score and Mean Absolute Error (MAE). MAE tells us the mean absolute difference between predicted sentiment scores and the true sentiment scores.

$$MAE = \frac{\sum_{i=1}^n |y_i - x_i|}{n}$$

Accuracy measures the proportion of true positives and true negatives out of total predictions.

F1 score measures the harmonic mean of precision and recall.

Figures 3 and 4 report the performances of the Mult model (respectively LF-LSTM) using our Arabic Multimodal dataset in terms of Accuracy, F1, and MAE metrics. For each model, four variants are learned. Three uni-modal models considered Text, Audio, and Video modalities alone. While TVA is the Multimodal that fuses the three modalities.

Let us mention that the uni-modal models for Mult are obtained by feeding the features of that specific modality with self-attention so that the CAB is replaced by a self-attention mechanism.

The results show that TVA Mult-based learner outperforms the uni-modal models regarding all metrics. It improves the accuracy by 15.15%, 19.63%, and 18.22% for Text, Audio, and Video-based uni-models, respectively, while MAE is improved by 2% to 4% compared to uni-modals.

Concerning the LS-LSTM-based models, the same result is observed. The TVA learner outperforms the uni-modal models regarding all metrics. However, with less improvement. Multimodality has enhanced the learner in terms of F1 score by 3.9%, 10.19%, and 5.6% for Text, Audio, and Video-based

<sup>1</sup><https://github.com/belgats/Arabic-Multimodal-Dataset/>

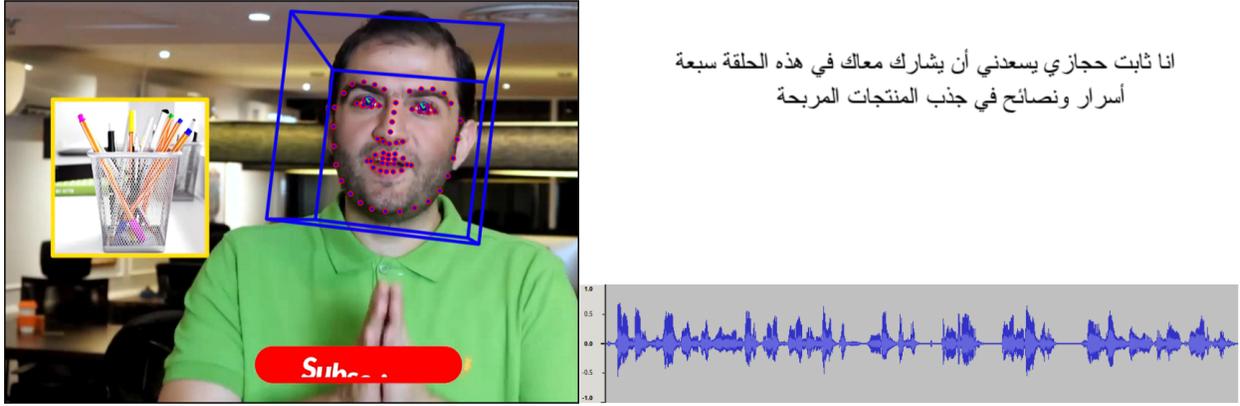
<sup>2</sup><https://github.com/ARBML/klaam>

<sup>3</sup><https://almufaragh.com/>

<sup>4</sup><https://www.audacityteam.org/>

<sup>5</sup><https://OpenSmilehttps://www.audeering.com/research/opensmile>

<sup>6</sup><https://github.com/belgats/Arabic-Multimodal-Dataset/>



انا ثابت حجازي يسعدني أن يشارك معاك في هذه الحلقة سبعة  
أسرار ونصائح في جذب المنتجات المربحة

Fig. 2. Illustration of an instance of our Dataset.

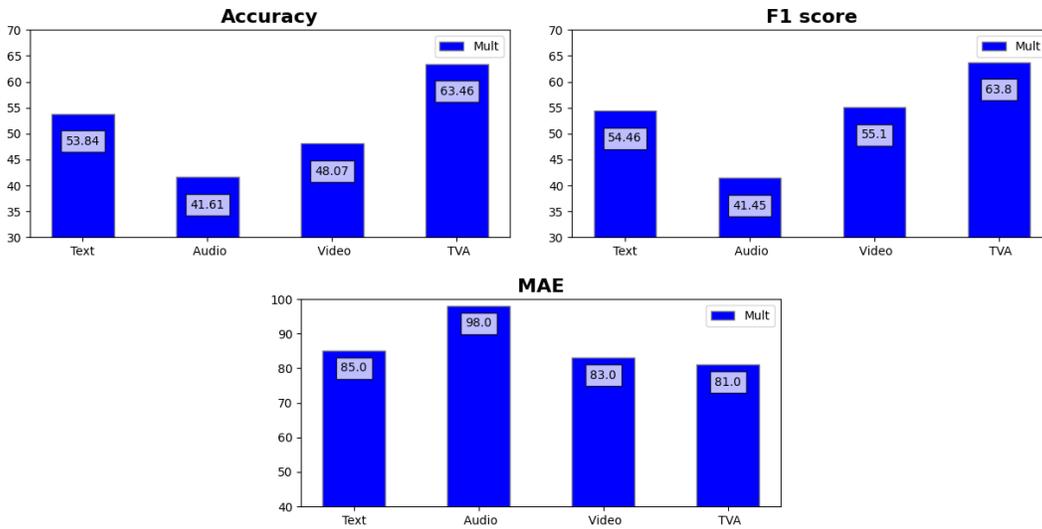


Fig. 3. Performances of Mult-based Models.

uni-models, respectively. Furthermore, the MAE has decreased by more than 8.64% for the text uni-modal model.

One can observe that the reached Multimodal based performances are not very high. F1 scores are about 63, 8% and 58, 9% for Mult and LF-LSTM-based models, respectively. However, these models show their superiority compared to uni-modal models.

One could argue that these results are impacted by two factors. Firstly, the dataset size is relatively modest and needs to be expanded to ensure greater accuracy. Secondly, the alignment process is highly challenging. As previously mentioned, we encountered significant difficulties when dealing with speeches where the speaker swallowed words, which negatively affected the word alignment.

Another result to be underlined is the superiority of modalities early fusion (Mult) compared to the late fusion (LF-LSTM), one at least for our built dataset. This result is expected as it is also confirmed for other Languages' Multimodal

models [19].

## VI. CONCLUSION AND FUTURE WORK

In this paper, we have addressed the topic of multimodal sentiment analysis, which has the potential to revolutionize our understanding and analysis of human emotions, opening up new avenues for research and practical uses. However, to address the issue of the scarcity of Arabic multimodal datasets, we have developed a methodology for creating such a dataset. Subsequently, we assessed the effectiveness of our constructed dataset using state-of-the-art transformer models designed to handle multimodality.

Despite the relatively small size of the constructed dataset, the findings show that considering multimodality is crucial for accurate Arabic sentiment analysis.

As further work, we intend to expand our Arabic multimodal dataset to meet the size requirements for deep learning algorithms. Furthermore, we conjecture that enhancing the

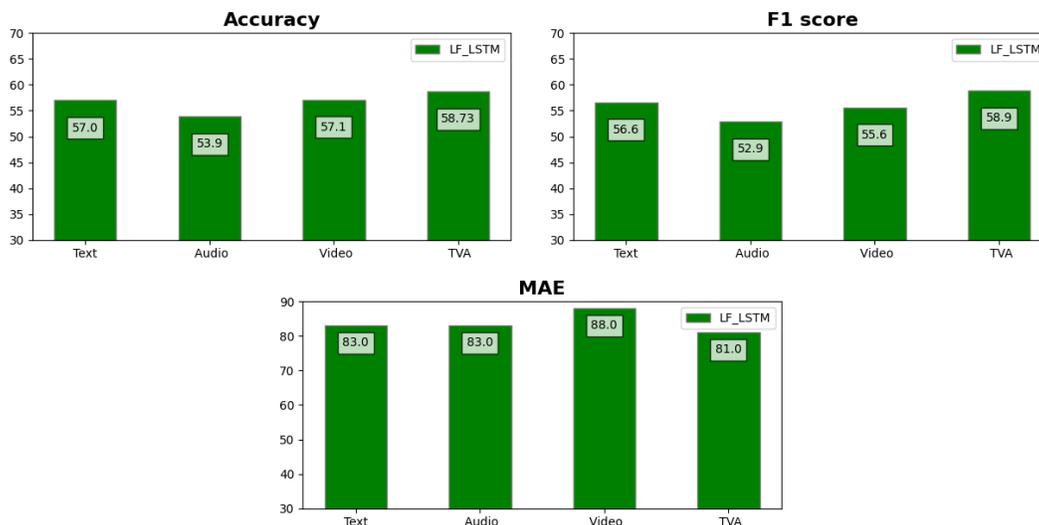


Fig. 4. Performances of LS-LSTM-based Models.

alignment techniques used in the dataset can considerably improve the accuracy and effectiveness of sentiment analysis.

#### REFERENCES

- [1] M. Soleymani, D. Garcia, B. Jou, B. Schuller, S.-F. Chang, and M. Pantic, "A survey of multimodal sentiment analysis," *Image and Vision Computing*, vol. 65, pp. 3–14, 2017.
- [2] A. Zadeh, "Micro-opinion sentiment intensity analysis and summarization in online videos," in *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, ser. ICMI '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 587–591. [Online]. Available: <https://doi-org.snd11.arn.dz/10.1145/2818346.2823317>
- [3] S. Poria, E. Cambria, N. Howard, G.-B. Huang, and A. Hussain, "Fusing audio, visual and textual clues for sentiment analysis from multimodal content," *Neurocomputing*, vol. 174, pp. 50–59, 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0925231215011297>
- [4] T. Diwan and J. V. Tembhurne, "Sentiment analysis: a convolutional neural networks perspective," *Multimedia Tools and Applications*, vol. 81, no. 30, pp. 44405–44429, Dec 2022. [Online]. Available: <https://doi.org/10.1007/s11042-021-11759-2>
- [5] J. V. Tembhurne and T. Diwan, "Sentiment analysis in textual, visual and multimodal inputs using recurrent neural networks," *Multimedia Tools and Applications*, vol. 80, no. 5, pp. 6871–6910, Feb 2021. [Online]. Available: <https://doi.org/10.1007/s11042-020-10037-x>
- [6] S. A. Abdu, A. H. Yousef, and A. Salem, "Multimodal video sentiment analysis using deep learning approaches, a survey," *Information Fusion*, vol. 76, pp. 204–226, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1566253521001299>
- [7] A. Ghallab, A. Mohsen, and Y. Ali, "Arabic sentiment analysis: A systematic literature review," *Applied Computational Intelligence and Soft Computing*, vol. 2020, pp. 1–21, 2020.
- [8] H. Najadat and F. Abushaqra, "Multimodal sentiment analysis of arabic videos," *Journal of Image and Graphics*, vol. 6, no. 1, 2018.
- [9] A. S. Alqarafi, A. Adeel, M. Gogate, K. Dashiypour, A. Hussain, and T. Durrani, "Toward's arabic multi-modal sentiment analysis," in *Communications, Signal Processing, and Systems: Proceedings of the 2017 International Conference on Communications, Signal Processing, and Systems*. Springer, 2019, pp. 2378–2386.
- [10] G. Vinodhini and R. Chandrasekaran, "Sentiment analysis and opinion mining: a survey," *International Journal*, vol. 2, no. 6, pp. 282–292, 2012.
- [11] L. Yue, W. Chen, X. Li, W. Zuo, and M. Yin, "A survey of sentiment analysis in social media," *Knowledge and Information Systems*, vol. 60, pp. 617–663, 2019.
- [12] L.-P. Morency, P. P. Liang, and A. Zadeh, "Tutorial on multimodal machine learning," in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Tutorial Abstracts*. Seattle, United States: Association for Computational Linguistics, Jul. 2022, pp. 33–38. [Online]. Available: <https://aclanthology.org/2022.naacl-tutorials.5>
- [13] P. P. Liang, A. Zadeh, and L.-P. Morency, "Foundations and recent trends in multimodal machine learning: Principles, challenges, and open questions," 2022. [Online]. Available: <https://arxiv.org/abs/2209.03430>
- [14] P. Xu, X. Zhu, and D. A. Clifton, "Multimodal learning with transformers: a survey," *arXiv preprint arXiv:2206.06488*, 2022.
- [15] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 18–24 Jul 2021, pp. 8748–8763. [Online]. Available: <https://proceedings.mlr.press/v139/radford21a.html>
- [16] 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.
- [17] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [18] W. Antoun, F. Baly, and H. Hajj, "Arabert: Transformer-based model for arabic language understanding," *arXiv preprint arXiv:2003.00104*, 2020.
- [19] Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov, "Multimodal transformer for unaligned multimodal language sequences," in *Proceedings of the conference. Association for Computational Linguistics. Meeting*, vol. 2019. NIH Public Access, 2019, p. 6558.
- [20] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [21] A. B. Zadeh, P. P. Liang, S. Poria, E. Cambria, and L.-P. Morency, "Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2236–2246.
- [22] A. Zadeh, R. Zellers, E. Pincus, and L.-P. Morency, "Mosi: multimodal corpus of sentiment intensity and subjectivity analysis in online opinion videos," *arXiv preprint arXiv:1606.06259*, 2016.
- [23] T. Baltrušaitis, P. Robinson, and L.-P. Morency, "Openface: An open source facial behavior analysis toolkit," in *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2016, pp. 1–10.