



# Multi-view informed attention-based model for Irony and Satire detection in Spanish variants

Reynier Ortega-Bueno<sup>a,\*</sup>, Paolo Rosso<sup>a</sup>, José E. Medina Pagola<sup>b</sup>

<sup>a</sup> PRHLT Research Center, Universitat Politècnica de València, Valencia, Spain

<sup>b</sup> Universidad de Ciencias Informáticas, Havana, Cuba

## ARTICLE INFO

### Article history:

Received 25 January 2021

Received in revised form 5 October 2021

Accepted 12 October 2021

Available online 22 October 2021

### Keywords:

Irony and satire

Attention mechanism

Linguistic features

Contextualized pre-trained embedding

Fusing representation

Spanish variants

Figurative language

## ABSTRACT

Making machines understand language and reasoning on it has been one of the most challenging problems addressed by Artificial Intelligent researchers. This challenge increases when figurative language is used for communicating complex meanings, intentions, emotions and attitudes in creative and funny ways. In fact, sentiment analysis approaches struggle when facing irony, satire and other figurative languages, particularly those where the explanation of a prediction might arguably be as necessary as the prediction itself. This paper describes a new model MvAttLSTM based on deep learning for irony and satire detection in tweets written in distinct Spanish variants. The proposed model is based on an attentive-LSTM informed with three additional views learned from distinct perspectives. We investigate two strategies to pass these views into MvAttLSTM. We perform an extensive evaluation on three corpora, one for irony detection and two for satire detection. Moreover, in order to study the robustness of our proposed model, we investigate its performance on humor recognition. Experiments confirm that the proposed views help our model to improve its performance. Moreover, they show that affective information benefits our model to detect irony and satire. In particular, a first analysis of the results highlights the discriminating power of emotional features obtained from SenticNet and SEL lexicon. Overall, our system achieves the state-of-the-art performance in irony and satire detection in Spanish variants and competitive results in humor recognition.

© 2021 Elsevier B.V. All rights reserved.

## 1. Introduction

Language itself is a perfect illustration of human creativity, and it achieves its splendor when some semantics rules and maxims of human communication [1,2] are disrespected to create expressions whose real meaning diverges from what it is apparently said. This peculiar usage of language with creative and funny purposes has been coined with the term *Figurative Language* [3–5]. Irony, satire, sarcasm, humor, puns, simile, hyperbole, metonym and metaphor are forms (or devices) of figurative language. While it is true that all these forms are used to communicate complex meanings, not all of them are used by common people. Some forms are relegated only to literary and poetry usages [6].

Irony and satire are pervasive and popular in everyday communication. As human beings, we appeal to these devices as effective ways through which literal meanings are intentionally deviated in favor of secondary interpretations. Particularly, both devices are acknowledged to express an attitude that is generally negative and implicit behind an apparent positive message. Thus, they are frequently used to criticize, complain, ridicule or mock.

Even when there are commonalities between both phenomena, irony seem to be more primitive and universal than satire.

The most common types of irony used in social media are situational and verbal irony. On the one hand, situational irony refers to specific events that fail to meet expectations [7] e.g. “*The fire station burns down while the firemen are out on a call*”. On the other hand, verbal irony has been traditionally identified as figurative device where enunciated words imply something other than their literal meanings. In other words, their real meaning is opposite to the literal one and it needs to be inferred through interpretation e.g. “*A burned tongue is a lovely way to start the day*”. Sarcasm is often considered a specific type of verbal irony which has a more aggressive tone [8], is directed toward an individual or a group [9–12], and is used intentionally [13,14]. An example of sarcasm would be the exclamation “*You’re really brilliant!*” about someone who has done a foolish act.

Satire is an interesting concept which is strongly related with irony and humor. It takes advantage of indirect speech and negative attitude implicit in irony. This device also appeals to features of humor such as: parody, exaggeration, juxtaposition, comparison, analogy, and double entendres with censoring purpose. Satirical messages may be aggressive and offensive, but they always have a deeper meaning and a social signification beyond that of

\* Corresponding author.

E-mail address: [rortega@prhlt.upv.es](mailto:rortega@prhlt.upv.es) (R. Ortega-Bueno).

the humor [15]. Satire does not make sense when the reader does not understand the real intent hidden in the ironic/funny dimension; like in irony, the real meaning of a satirical message lays in the figurative interpretation of the content. Satire can be separated in two distinct directions: *Juvenalian* or *Horatian* styles [16]. On the one hand, the *Juvenalian* style of satire is based on ridicule and sarcasm. On the other hand, the *Horatian* style contains tease and humor.

Irony and satire have been studied from many disciplines such as Linguistics, Psychology, Rhetoric, Pragmatics, Semantics, etc., however, they are not only enclosed to these theoretical studies. Nowadays, both devices are typically used in social media platforms to favor social interactions, evoking humor [17], diminishing or enhancing criticism [18,19], and getting the attention of the readers by means of the creativity [20]. These forms of figurative language have great impact on several other Natural Language Processing (NLP) tasks that aiming at monitoring social media content. In some cases the presence of ironic message plays a specific role: “*implicit polarity reversal*”. This means, that a message seems to be positive but its real meaning is negative (or vice-versa). Due to this peculiarity of ironic and sarcastic expressions, the sentiment analysis approaches decline when facing irony in social media texts [21–26]. In fact, this problem become more challenging in sentiment analysis approaches where the explanation of the results is more important than the decision itself [27,28]. Sarcasm as a specific sub-type of verbal irony has implications to cope with the bad phenomenon of miscommunication, particularly: hate, aggressiveness, and nastiness speech [29]. Ignoring the presence of sarcasm causes that the implicit meaning, generally hurtful and offensive, be misunderstood with the results of exposing people to toxic information. In this context also it results crucial to understand what pieces of message are relevant. Recently, interesting evidence about the use of satire to disguise fake news has been discussed in [30,31]. For people, understanding satire as fake messages may deprive them of desirable entertainment content, while recognizing fake information as legitimate satire may expose them to disinformation or misinformation.

Following the timeline of computational methods for irony and satire detection, it is possible to envisage two distinguishable approaches namely, classical features-based machine learning approach, and deep-learning approach. In the literature many works explored several linguistic, stylistic, content, affective, and contextual features to address the problem in a shallow supervised way [32]. The handcraft features derived from this approach have proved to be feasible when small dataset are provided.

In the last five years, deep learning techniques became very popular in NLP, and applied to irony and satire detection. These methods show a better performance than classical feature-based machine learning models. In this scenario, a vast number of works are based on attentive-Recurrent Neural Networks (att-RNN), particularly by using of Long Short Term Memory (LSTM) [33] and Gated Recurrent Unit (GRU) [34] with attention mechanisms, which have been effective to capture complex dependencies among words within the text and pay more attention to those words that increase the effectiveness of these networks in several tasks of NLP [35–38].

Recently, a groundbreaking advance in NLP has been marked by using transformer-based network architectures [39] which opened a new avenue for training robust and contextual-aware word embeddings in unsupervised manner. The use of these pre-trained contextualized models has been widely spread by means of Bidirectional Encoder Representations from Transformers (BERT) [40] and BERT’s family architectures [41–44]. Clearly, this spreading has reached figurative language processing, and cause that these new computational methods outperformed the state of the art by a substantial margin [45–47].

The nested non-linear functions of deep learning algorithms provoke these models are usually applied in a black-box manner, that is, no interpretable knowledge is provided about what exactly causes them to arrive at their predictions. In this sense, the attention mechanisms became a (albeit narrow) way for dealing with the problem of model interpretability. Recently, attention as a way of explainable deep learning method became a very popular and controversial topic. Some works claimed that attention weights do not provide meaningful “explanations” for supporting the final predictions [48,49], however other works state that they are able for discovering how neural models capture several linguistic notions of syntax, semantic and coreference [50–52]. Despite diverse views on the matter, empirical results on the task of binary irony classification show that attention mechanisms are able for capturing ironic cues, word polarity and explicit and implicit sentiment incongruity [45,53].

Even when, remarkable advances have been observed in irony and satire detection, these advances have showed an asymmetric development with respect to languages other than English. This is the case of Spanish and its variants, where few researchers have addressed the problem. In the case of satire, only the works [54, 55] have studied the phenomenon in the Spanish language by using a machine learning-based approach. Irony detection in Spanish variants has been surveyed in [56], but few methods relied on deep learning approaches [57–60]. According to our knowledge only the works [45,60] take advantages of contextualized pre-trained word embedding. For that, more efforts must be paid to study irony and satire in Spanish language. In this work we propose an attentive-based deep learning method in order to investigate further the detection of irony and satire in Spanish:

- Irony and satire are pragmatic phenomena, hence both are contextual-dependents. Additional knowledge such as: language and its variety, sociolinguistics and cultural background are crucial for precisely recognize and understand these forms of figurative language.
- Many studies on irony and satire detection have been conducted from three directions: linguistics features with machine learning approaches, deep learning techniques based on att-RNN, and recently by means of contextualized pre-trained word embeddings with a single-modality representation of the texts. However, there are no works that pay attention to explore irony and satire in Spanish from a fusion information perspective where these three approaches are fused aiming to outperform the state of the art.

To overcome these challenges, we aim at addressing the following research questions:

- RQ1. Could irony and satire detection methods take advantage of combining multiple representations (views) of text in terms of linguistic-based representation, universal sentence encoder-based representation, and contextualized pre-trained embeddings?
- RQ2. How the proposed views can be effectively combined to proper inform an attentive recurrent model?
- RQ3. Are multiple heads of attention (multi-head) more feasible than single attention (self) for capturing multiples and complex relations among words in ironic and satirical texts?

With the aim to answer the formulated research questions, in this work we propose a new model (*MvAttLSTM*) which relies on a multi-view informed attentive-LSTM neural network. We consider an attentive recurrent model due to the attention mechanisms allow the model to focus and place more “attention” on the relevant parts of the text sequence in order to capture complex syntactic and semantic properties used in ironic and satirical messages. Specifically, we learn three independent views

for each text, and we pass them to our MvAttLSTM. The first one (*Linguistic-view*) is based on several linguistics features which have proved to be strong cues for discriminating both irony and satire. The second one considers a deep dense encoding of the text by means of Multilingual Universal Sentence Encoder (*MUSE-view*). And, finally the last one (*BERT-view*) considers a contextualized pre-trained embedding obtained after a tuning of the BERT model. We evaluate the effectiveness of our method on one corpus for irony detection and on two distinct corpora for satire detection. For irony detection, the corpus is the one proposed for the *IroSvA'19* shared task: *Irony detection in Spanish Variants* [56] was used, whereas, in the case of satire detection task the corpora introduced in [54,55] were employed. Our proposal outperformed both previous systems participating in the *IroSvA'19* shared task and recent methods [45,61]. Also, for satire detection our proposal outperformed previous methods by a substantial margin in both corpora. Additionally, we provide several analyses in order to evaluate two strategies for fusing the learned views into MvAttLSTM and investigate the impact of each view on the performance of MvAttLSTM. Finally, an interesting analysis is carried out on the attention mechanism to observe how our proposal takes into account some features related with irony and satire such affective content. In short, the major contributions of this paper are summarized below:

- To investigate the problem of computational irony and satire detection in Spanish variants in three widely used corpora. Moreover, taking into account the closed relation among irony, satire and humor, we evaluate the robustness of the proposed model on humor recognition.
- To propose a novel approach (MvAttLSTM) based on representation fusion. Particularly, efficient representations from three distinct perspectives are computed and combined to inform an attentive-LSTM model. The proposed method outperforms the state of the art approaches in satire and irony detection in Spanish and obtains competitive results in humor recognition.
- To investigate distinct forms of combining the proposed views, also to evaluate the impact of each view on the proposed model MvAttLSTM.
- To study the impact of two kinds of attention mechanisms, self attention vs. multi-head attention, on the proposed model.

The rest of this paper is structured as follows. Section 2 introduces the state of the art for both irony and satire detection, with special interest in those approaches proposed for Spanish. Section 3 formalizes our proposal based on a multi-view attention based model. Particularly, we describe each one of the representation used and two distinct ways of fusing these views for irony and satire detection. In Section 4 a detailed description of the corpora, resources, preprocessing and the experimental setup is introduced. Also, an exhaustive evaluation of our proposal and a comparison with other approaches is presented. Moreover, considering the closed relation among irony, satire and humor, we evaluate the robustness of our model to recognize humor. Finally, we draw some conclusions and discuss future work.

## 2. State of the art

There is a considerable amount of literature on computationally irony and satire detection [see 62–65]. In general, approaches to deal with these forms of figurative language can be classified into: features-based machine learning approach and deep learning-based approach. Initially, machine learning method combined with feature-based representations (lexical, contextual, stylistic, affective, discursive, etc.) received the most attention.

But, recently, deep learning-based approaches are gaining interest due to the capacity of these models to automatically learn feature representations that are omitted in hand-craft extraction or simply have abstraction levels beyond of human bounds. In this line, the next two subsections survey the relevant works for irony and satire detection. Also, considering the imbalanced number of studies in other languages than English, a third subsection discusses irony and satire in a multilingual setting, with special focus on Spanish.

### 2.1. Machine learning based approaches

*Irony.* Computational irony detection has been addressed by the NLP community from different perspectives. In preliminary works, the role of textual-based features obtained from the text (such as n-grams, punctuation marks, part-of-speech tags, among others) has been widely explored for its detection [66–70]. Other works drew attention to theoretical aspects of irony such as incongruity and opposition. Based on these aspects, features derived from semantic ambiguity, synonyms, antonyms and polarity contrast have been studied in [25,71–73]. Many theories seem to agree that an implicit attitude is expressed when being ironic. Aiming to capture the relation between irony and subjectivity in language, several approaches have focused on affective information for improving irony detection [74–78]. Verbal irony is without doubt a pragmatic phenomenon, hence, contextual and extra-linguistic information result crucial for its detection and comprehension. In this sense, information regarding the context surrounding a given text has been exploited in order to determine whether a text has an ironic or sarcastic intention [79–82]. Discovering new features with discriminative power and topic-independency have been the most active directions of machine learning approaches. Regarding machine learning algorithms, the most used have been Random Forest (RF), Decision Trees (DT), Naïve Bayes (NB) and Support Vector Machines (SVM). Recently, in [83] the impact of the imbalanced distribution of classes in irony and sarcasms detection has been studied from a machine learning perspective.

*Satire.* Machine learning has been the most used approach for satire detection [54,55,84–86]. In the seminal paper of Burfoot and Baldwin [84], the problem of detecting satire was explored with simple bag of words features (BoW) using two feature-weighting methods: (i) binary feature weighting and (ii) binormal separation (BNS) features scaling. Further, lexical (headlines, profanity, slang) and semantic features were added to enrich text representation. To compute the semantic feature they identify the named entities in a given document and query the web for the conjunction of those entities. In this direction, [85] proposed to extent the BNS features scaling method with the *tf-idf* weighting schema to improve satire detection in news genre.

In [54] studied the problem of satire detection in tweets. Linguistic differences between satirical and factual content were explored by mean of frequency, ambiguity, synonyms, part of speech (PoS) tags, sentiments, characters, and slang words as features. Experiments showed that some linguistic features are topic-independent and hence useful clues to address the problem. In a same fashion, a psycho-linguistics approach was introduced in [55] to identify satirical tweets. A wide variety of psychological and linguistic features from Linguistic Inquiry and Word Count lexicon (LIWC) [87] were evaluated. Results confirmed the usefulness of emotional, social, and psychological dimension for satire detection. In [30] were considered five predictive features: absurdity, humor, grammar, negative affect, and punctuation, and applied an SVM method. After, combining three out of five features (absurdity, grammar, and punctuation), the authors observed that the BNS feature scaling is suitable for satire detection and the model obtains good results.



Sensibility of lexical, linguistic and n-gram based features across three textual genres was reported in [88]. Specifically, the impact of features associated on affective words, acts of the speech, sensorial words, and shallow clues of figurative device (alliteration, grammatical inversion, hyperbole, onomatopoeia and imaginary) was evaluated. Results showed that n-grams and features related with the act of the speech were good as genre-independent and hence they resulted suitable for satire detection in multiples genres. In a similar fashion, an emotions and sentiments based representation was proposed in [89] for satire detection in newswires, Amazon product reviews and an in-house Twitter corpus. Experiments were performed using the SVM and RF methods. Results concluded the usefulness of the proposed features for satire detection. In [90] the impact of emotions on recognizing satirical texts from other figurative forms (humor, irony, sarcasm) and factual language was analyzed.

Recently, satire has received more attention due to the commonalities with the undesirable phenomenon of misinformation in social media, and particularly with fake news spreading [31, 91,92]. In order to reduce the exposure to misinformation in social media, publishers of fake news have begun to masquerade as satire sites to avoid being demoted. For users, incorrectly recognizing satire as fake news may deprive them of desirable entertainment content, while identifying a fake news story as legitimate satire may expose them to misinformation.

## 2.2. Deep learning-based approaches

*Irony.* Recently, many deep learning-based approaches for addressing irony detection have been proposed. Word embeddings, Convolutional Neural Networks (CNNs), att-RNNs, and Transformers-based models have been exploited for capturing the presence of irony in social media content [45,46,93–103]. The semantic and syntactic properties of pre-trained word embeddings have been highlighted in several studies [96,97,104]. For instance, word embeddings have been explored to capture incongruity in text with non-affective words [96]. In the study [104], irony detection was addressed using a multi-faceted representation which fuses psycho-linguistic features with word embedding vectors that were obtained by using Doc2Vec [105]. In [97] the generalization capabilities of an unsupervised topic model trained for irony detection showed a substantial increasing when the word embedding information was incorporated.

Several methods exploited the advantages of CNN for discovering local features that result useful for irony detection. An interesting idea was proposed in [98], which introduced a framework for learning irony features from a corpus using CNN. This approach investigated whether features extracted using pre-trained sentiment CNN, emotion CNN and personality CNN models can improve the overall performance. In another direction, the role of the content and contextual information for sarcasm detection taking advantages of a multi-view model were presented in [99]. For that purpose, two CNN models were trained to generate stylometric and personality embeddings for each user's comments. Later, both embedding were fused in a multi-view setting using Canonical Correlation Analysis (CCA) [106]. A content-based sentence representation was extracted using another CNN and appended with context vectors to obtain the final decision. In another study [93], a model that combines dense neural networks (DNNs) with time-convolution and LSTM (CNN-LSTM) was proposed for detecting sarcasm in tweets. These existing studies use the convolutional network to automatically derive deep features from texts for irony detection. Results of these deep learning-based models are generally better than those obtained with classical feature based machine learning methods.

RNNs have been used for addressing irony detection due to their abilities for capturing long and short dependencies among

words within texts. In [94] studied the role played by the conversational context in a sarcasm reply. Particularly, results proved that LSTMs that can model both the context and the sarcastic reply achieve better performance than LSTMs that read only the reply. In another direction, many approaches studied the impact of attentive-based representation with linguistics features [100, 107]. Experiments have concluded that considering hand-crafted features help models to increase their effectiveness. From another point of view, the model introduced in [53], proposed strategies to improve irony detection by transferring knowledge from sentiment resources. This work proposed three different attentive-LSTM approaches that differ in the way of including the sentiment resources, either injecting the sentiment directly to the attention mechanisms or merging the output of different networks specialized on sentiment analysis and irony detection. In a similar fashion, in [108] a multi-task learning approach was proposed to leverage the knowledge in sarcasm detection and sentiment analysis task. Experiments showed that these two tasks are correlated, and training a deep neural network that models this correlation in a multi-task learning setting improves the performance of both tasks. Moreover, in [109] a multi-task learning framework for multi-modal sarcasm, sentiment, and emotion analysis was proposed. The authors take advantage of the sentiment and emotions of the speaker to predict sarcasm. In the multi-task framework, sarcasm was considered as the main task, whereas emotion and sentiment detection were used as secondary tasks. Results confirmed that the multi-task framework achieves better performance for the primary task, i.e. sarcasm detection, with the help of emotion and sentiment analysis tasks.

The use of transformer-based models [39] has changed the way of modeling and working with textual data in an unprecedented way. In fact, these models have been widely spread by means of BERT [40] and other BERT's related architectures [41–44]. Clearly, this spreading has reached very fast also FL processing: new methods based on transformer models outperformed the state of the art in irony detection by a substantial margin [45–47]. In this line, in [46] the RoBERTa model [42] was used to encode the sentences, that was further contextualized by means of a Recurrent Convolutional Neural Network to address irony and sarcasm detection. This model outperformed state of the art on four benchmark datasets for irony and sarcasm detection in English. In [45] a simplification of the BERT architecture was proposed to contextualize pre-trained word embeddings. Specifically, this work contextualized Word2Vec word embeddings, trained with several millions of tweets both for English and Spanish. This strategy, opposite to the use of pre-trained BERT, aimed to train the proposed model from in-domain data using the same powerful backbone architecture as BERT. This model outperforms previous models for irony detection in Spanish short texts.

*Satire.* Notwithstanding a vast amount of deep learning-based methods have been proposed for irony detection, and the commonalities between irony and satire, few methods have addressed the problem of satire detection from a deep learning perspective. Recent works in this direction have been presented in [110–112]. In [110] a four-level hierarchical network with attention mechanism was presented to differentiate satirical news from true ones. Psycholinguistics, writing stylistic, structural and readability-based features were included to the model at both paragraph and document level. The evaluation suggested that readability features supported the overall classification while psycholinguistic features, writing stylistic features, and structural features are beneficial at paragraph level. The analysis of individual features reveals that satirical news tend to be emotional and imaginative. Another idea was explored in [111] which proposed to use CNN, LSTM, and GRU to detect satire at both sentence and document levels. They concluded that fine-grained sentence-level analysis provides an in-depth insight into the phenomenon of satire.

### 2.3. Multilingual approaches

Most of the works on irony and satire detection have investigated the problem in English. Notwithstanding, there have been some efforts to investigate it in other languages such as: Chinese [113], Czech [70], Dutch [69], French [114,115], Italian [24,116,117], Portuguese [66], Spanish [56,118,119], and Arabic [120,121]. Even when in closely related tasks like sentiment analysis have emerged an increasing number of works addressing the multilinguality issue [122–128], where few works explored this in the context of irony and satire detection. Taking advantage of the finding achieved for multilingual sentiment analysis would be an interesting direction to improve satire a irony in this scenario.

From a multilingual perspective, the approaches for irony and satire detection can be analyzed in two main directions: (i) multilingual setting, where the model is trained and evaluated separately on each language, (ii) cross-lingual setting, where the models is trained in one or more language and evaluated on another different one. Multilingual setting has been the most investigated. Prior works were presented in [70] and [113] for Czech–English and Chinese–English languages respectively. In [129] a novel fine-grained annotation schema was proposed to annotate irony categories, activators and markers in French, English and Italian language. The role played by dependency-based syntactic features on irony detection from a multilingual perspective (English, Spanish, French and Italian) was investigated in [130]. In the case of satire, in [55] the authors investigated the impact of psycholinguistic features in two distinct variants of the Spanish (Mexican and Castilian). Irony detection from a cross-lingual perspective (Arabic, French and English) was investigated in [131]. Results showed that, although irony is contextual, language and cultural-dependent pragmatic phenomenon, several features are universal and can be useful for addressing irony detection in languages which lack of annotated data. In the same line, in [86] the authors presented a set of language independent features that describe lexical, semantic and usage-related properties of the words in the tweets. The proposed features were evaluated in a cross-lingual setting. Results highlighted the complexity of modeling satirical texts in a cross-lingual setting, due to satire aims at criticizing social and moral behaviors which often are social and cultural-dependent.

### 3. Multiview informed attention-based models

In this section we introduce MvAttLSTM, our multi-view informed attentive LSTM model for irony and satire detection in Spanish. We addressed both tasks as binary classification problems applying a model based on LSTMs endowed with an attention mechanism. LSTM is an RNN that uses gating mechanisms to overcome the problem of the vanishing gradient. This type of neural networks can capture long-range relationships and hidden patterns in sequential data. In terms of architecture, the Bidirectional LSTM (BiLSTM) [132] is widely used, which has two LSTM units processing sequences forward and backward respectively. This property of BiLSTM is useful for language processing because the meaning of the words in texts can be inferred not only by previous words, but also considering other words after them can help to determine their meanings. Moreover, attention mechanisms have endowed the RNNs with a powerful strategy to enhance their performance and achieve better results. Our model considers multiple representations learned from three distinct perspectives: linguistic-based representation, universal sentence encoder-based and contextualized pre-trained embeddings. We introduce additional knowledge into MvAttLSTM model aiming at reinforcing linguistics and semantics properties which can result beneficial for detecting irony and satire. Concretely, our

model is compounded by an embedding layer which is fed into a BiLSTM layer. Later, the hidden states sequence returned by the BiLSTM is fed into an attention layer. Next, on the output of this layer are staked two LSTM layers. Finally, we incorporate a feed forward neural network for final prediction. As explained before, we inform the model with three additional views. Particularly, we investigate two different strategies for fusing these views into our MvAttLSTM. In the next subsections we present in detail the preprocessing carried out on the datasets, the additional representations, the main parts of the MvAttLSTM's architecture and the strategies for informing the model.

#### 3.1. Preprocessing

Social media texts, particularly those from Twitter are informal and noisy. The length constraints, and the free writing style present in this form of online communication provoke that texts have plenty of grammar and spelling mistakes. Particularly, length constrains caused that users use shortenings, abbreviation, homophonic encoding to save characters, and grammar and spelling misuses such as: character flooding, word repetition and wrong use of uppercase letters to denote emphasis. Twitter also offers to the users reserved symbols to mark explicitly important concepts in tweets (*# hashtag*), to refer o mention other users (*@ mention*), to reply the message of other users (*RT retweets*) or simply to mark texts as favorite (*FAV*). Aiming to inject emotional states tone and body language into tweets, emoticons and emojis became very popular. These symbols are an ultra-concise way to enrich writing language with visual information. All these issues impact sentence structure, content, word forms and increase the difficulty of their automatic processing and comprehension. In order to mitigate the effects of these problems, in our model we applied a basic preprocessing phase for cleaning the texts. Firstly, we applied a tokenization process on the tweets by using the TokTokTokenizer from NLTK [133]. Later, emoticons, emojis, URLs, hashtags, mentions are recognized and replaced by a corresponding wildcard which encodes the meaning of these special words. In the case of hashtag, we replaced the reserved symbol (*#*) by the word *topic\_* and retain the remaining characters. Emoticons and emojis were replaced by the word *emoji\_* concatenated with an integer value associated to each emoji. We have included these changes in order to reduce the impact of noisy and inconsistent writing on the processing of the text with the FreeLing tool [134]. Moreover, we replaced each mention and URL by the words *author\_token* and *url\_token* respectively. Finally, Twitter-reserved words like RT (for retweet) and FAV (for favorite) were removed. It is worth to notice that emoticons and emojis are a valuable source of information to take into account in social media content analysis [135–141]. Nevertheless, in this work we used emojis and emoticons to create features for capturing the frequency of positive emojis, negative emojis and neutral emojis as well as detecting polarity contradiction between the words and emojis in the text. In the second stage, and used only to obtain some linguistic features that were considered in the *Linguistic-view*, a more complex language analysis was carried out. For that, flooding tokens were normalized allowing the same character to appear only twice consecutively in a token (e.g. *hoolaaa* becomes *hoolaa*). Afterward, tweets were morphologically analyzed with the FreeLing tool. In this way, for each resulting token, its lemma and part-of-speech were considered.

#### 3.2. Addition knowledge to inform the model

Our MvAttLSTM relies on fusing multiple representations which are learned from distinct perspective. Specifically, we learn three independent views for each text which are introduced to

the model. The first one (*Linguistics-view*) consists in several linguistics features which have proved to be strong cues for discriminating both irony and satire. The second one considers a deep dense encoding of the message using Multilingual Universal Sentence Encoding (*MUSE-view*) [142,143]. And, finally the last one (*BERT-view*) considers a contextualized pre-trained embedding obtained after a tuning of the BERT model [40]. Next, we describe the three ways in which each view was learned.

### 3.2.1. Linguistics-view

Hand-crafted features, often linguistic-based, have proved to be effective for processing figurative language, particularly in case of irony, satire and humor [32,62–65]. From our perspective, the linguistics-based representation is able to capture certain types of irony, satire and other figurative devices disregarding textual genres and topics, and it makes this representation content independent and genre-unbiased. To represent each text, we use different group of features: stylistic and structural, semantics, affective, incongruity and psycho-linguistic. Many features are extracted to identify stylistic patterns in the structure of the ironic or satirical texts (e.g., type of punctuation, length, emoticons, distribution of nouns, adjectives, adverbs and verbs). Other features are extracted to consider affective information (e.g., polarity, sentiments, emotions, attitudes, etc.) by using several word-based lexicons resources. Moreover, features are extracted for considering semantics properties of texts (e.g., co-occurrence of synonyms and antonyms, maximum, minimum and mean of synsets, etc.). Finally, some features are designated to capture contrast and opposition in texts (e.g., polarity contrast, semantic incongruity, etc.). Specifically we use the features proposed in [144,145], the incongruity features used in [146] but using BabelSenticNet [147] as default polarity lexicon and including the emotional dimensions and the polarity feature in this resource as other affective features. BabelSenticNet is an extension of SenticNet [148] to 40 other languages, including Spanish (henceforth we refer the Spanish version of BabelSenticNet as SenticNet). For more details about the linguistic features considered in this work please see [Appendix A](#).

### 3.2.2. Task-independent embedding view

Our second representation aims at encoding the whole meaning of the text into a single dense vector based on deep learning models. Particularly, in vectors which capture rich semantic information that can be useful for recognizing semantics proprieties of the ironical and satirical texts. A contextual approach for creating the embedding vectors is proposed in [142], where complete sentences, instead of words, are mapped into a latent vector space. The approach provided two variations of Universal Sentence Encoder (USE) with some trade-offs in computation and accuracy. The first one consists of a computationally intensive transformer that resembles a transformer network [39], proved to achieve a higher performance. In contrast, the second one provides a lightweight model that averages input embedding weights for words and bi-grams by utilizing of a Deep Average Network (DAN) [149]. The output of DAN is passed through a feed-forward neural network in order to produce the sentence embedding. Both approaches take as input lower-cased strings and output a 512-dimensional sentence embedding. Although there are several methods like Doc2Vec [150], Sent2Vec [151], FastText [152] and InferSent [153] to generate sentence embeddings we used Multilingual Universal Sentences Encoding (MUSE)<sup>1</sup> [143] which is an extension of USE trained for 16 languages including Spanish. The most salience characteristic of this model is that it was

trained using multi-task learning to integrate semantic information. Particularly, sentence embeddings are learned across several languages and using multiple semantic tasks like sentiment analysis, semantic textual similarity, etc. This enables the learning process to dynamically accommodate a wide variety of knowledge in a single vector which is interesting to transfer to related tasks like irony and satire. Based on the MUSE model we transform the texts of the training dataset into dense vectors of 512 dimension (henceforth,  $H_{MUSE}$ ). It is important to highlight that  $H_{MUSE}$  is a completely task-independent representation, because we do not apply any parameters tuning of the model on the training data.

### 3.2.3. Task-dependent embedding view

The transformer-based neural network architectures [39] paved the way for training robust and contextual-aware language models in an unsupervised manner. The use of these pre-trained contextualized models have been widely spread through BERT [40] and BERT's family architectures [41–44]. To study the linguistic and semantic nuances of irony and satire in Spanish, we decide to incorporate BERT as another representation (view). BERT relies on bidirectional representation from transformers and achieves the state of the art for contextual language modeling and contextual pre-trained embeddings. This model is trained on a large text corpus and then used for downstream NLP tasks. While other word embedding like Word2Vec [154], Glove [155] and FastText [152] are context-free models that produce a single word embedding for each word in the vocabulary, BERT computes a representation of each word that is based on the other words in the context. It was built upon recent works in pre-training contextual representations, such as ELMo [156] and Universal Language Model Fine-tuning (ULMFiT) [157], and is deeply bidirectional. BERT represents each word using both its left and right contexts. Moreover, it is possible to fine-tune BERT for many downstream NLP tasks, including the tasks we are interested in. This goal can be achieved by removing the language modeling output layer (masked word prediction) and replacing it with a new layer appropriate for the target task (in our case, binary classification). Particularly, in this work we use the pre-trained multilingual versions of BERT<sup>2</sup> (mBERT, henceforth) and carried out a fine-tuned on it, using the training datasets of irony and satire. Our idea is not to use this model for as a classification method; instead, we considered it for the representation purpose.

For fine-tuning mBERT, we add a layer that receives as the input the vector in the first position (the CLS token). On this layer, we stacked an output layer that makes the final prediction for the targeted task. For that purpose, we follow the strategy proposed in ULMFiT [157]. For each layer of mBERT, a different learning rate is set up, increasing it using a multiplier while the neural network gets deeper. This multiplier increases 0.1 points from a layer  $L_i$  to another  $L_{i+1}$ . We use this dynamic learning rate to keep most information from the pre-training at shallow layers and biasing the deeper ones to learn about the specific tasks. For all corpora, the same hyperparameters were used. Concretely, we defined the *batch\_size* = 32 and the sequence length was limited to 50 tokens. The optimizer used is Adam [158] with an initial learning rate of 0.001,  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$  and a *weight\_decay* = 0.01. The model was trained during 15 epochs and using the ModelCheckpoint callback for obtaining the model that has achieved the best performance on the validation subset.

After tuning mBERT, we pass again the training dataset, but this time, we get the deep representation of each text of the training dataset (henceforth,  $H_{BERT}$ ). This view is task-dependent because we refine the parameters learned by mBERT in order to capture semantics and pragmatics characteristics, which results crucial for understanding and recognizing ironic and satirical intents.

<sup>1</sup> <https://tfhub.dev/google/universal-sentence-encoder-multilingual/3>

<sup>2</sup> <https://huggingface.co/bert-base-multilingual-cased>



### 3.3. MvAttLSTM model

Let us describe the architecture of the MvAttLSTM model. We give details about each layer, starting from the embedding layer to the loss function.

#### 3.3.1. Embedding layer

In this layer each word  $w_i$  is map into a highly dimensional feature space for capturing the meaningful semantic and syntactic information. Given an input text  $T$  which consists of at most  $N$  words  $w_i$ , where  $i \in [1, N]$ . For each word into  $T$ , we examine the embedding matrix  $E \in \mathbb{R}_{B \times d}$ , where  $B$  is the length of the vocabulary, and  $d$  is the dimension of word embedding vectors. The matrix  $E$  can be initialized randomly or by means of a word embedding matrix. In this work we decided to initialize the embedding layer with context-free pre-trained word representations. For that, we learned the embedding matrix  $E$  by using the FastText model trained on the Spanish Billion Words Corpus<sup>3</sup> and an in-house background corpus of 9 millions of Spanish tweets. We aim to join both corpora for obtaining robust word representations taking advantage of the peculiar writing style used in Twitter. For training the FastText model we used the setting reported in [152], except for the value of the vector size, which was defined as a 300-dimensional. In this layer, each word  $w_i$  is transformed into a vector  $x_i \in \mathbb{R}^d$ :

$$H^0 = \text{Embedding}(E, M) \quad (1)$$

Thus, every text  $T$  can be converted in a sequence of vectors, in the form of a 2d-matrix  $H^0 = [x_1, x_2, x_3, \dots, x_N]^T$  with shape  $N \times d$ . The matrix  $H^0$  is given as input to the next layer. It is worth noting that in the model the weights in  $E$  are fixed. We aim at making the model to be trained faster and mitigate the impact of the overfitting due to the reduction of parameters that must be learned. Moreover, we consider that the recurrent and attention layers are feasible to take advantage of the semantic and syntactic properties of the vectors in  $E$  for classifying irony and satire.

#### 3.3.2. BiLSTM layer

After passing the sequence of word  $T$  to the Embedding layer, each word  $w_i$  is encoded by a vector  $x_i$  which captures the semantic and syntactic properties of  $w_i$  out of context. In other words, the representation of each  $w_i$  is independent of the other words in the text  $T$ . In this layer, a new representation for each word is learned by summarizing the contextual information, previous and after to the word in the text. For achieving this goal we use a BiLSTM layer. This, type of neural network consists of two LSTM units which process the sequential input in both directions forward and backward simultaneously.

$$H^1 = \text{BiLSTM}(H^0) \quad (2)$$

The output of this layer is a sequence of hidden states  $H^1 = [h_1^1, h_2^1, \dots, h_N^1]$  where each  $h_i^1 \in \mathbb{R}^{2 \times d_h}$  is the concatenation of the hidden state of each LSTM (right and left), specifically, the  $h_i = [\vec{h}_i^1, \overleftarrow{h}_i^1]$ , and  $d_h$  is the number of hidden neuron into the LSTM unit. Standard LSTM receives sequentially (in a left to right order) at each time step a word vector  $x_i$  and produces a hidden state  $h_i$ . For that, this neural network relies on a cell of memory and a gating mechanism consisting of an input gate, forget gate, and output gate. These gates help to determine whether the information in the previous state should be retained or forgotten in the current state. Hence, the gating mechanism helps the LSTM

to cope with long-term information preservation. Each hidden state  $h_i$  is determined as follows:

$$I_t = \sigma(W^i x_t + U^i h_{t-1} + b^i) \quad (3)$$

$$F_t = \sigma(W^f x_t + U^f h_{t-1} + b^f) \quad (4)$$

$$O_t = \sigma(W^o x_t + U^o h_{t-1} + b^o) \quad (5)$$

$$\tilde{C}_t = \sigma(W^u x_t + U^u h_{t-1} + b^u) \quad (6)$$

$$C_t = i_t \odot \tilde{C}_t + F_t \odot C_{t-1} \quad (7)$$

$$h_t = O_t \odot \tanh(C_t) \quad (8)$$

Where all  $W^*$ ,  $U^*$  and  $b^*$  are parameters of the recurrent layer which are learned during the training phase and the  $x_t$  is the pre-trained vector of the word in the time-step  $t$  and it is not trained in the model. The operator  $\sigma$  is the sigmoid function and the operator  $\odot$  stands for element-wise vector multiplication. The  $I_t$ ,  $F_t$  and  $O_t$  are the input, forget and output gates in the time step  $t - 1$  whereas  $\tilde{C}_t$ ,  $C_t$  and  $h_t$  are the new cell, the updated cell memory and the final hidden state in the time step  $t$ . Notice that the BiLSTM initial hidden states and cells memory are set to 0 in both directions  $\vec{c}_0^1 = \overleftarrow{c}_0^1 = \vec{0}$  and  $\vec{h}_0^1 = \overleftarrow{h}_0^1 = \vec{0}$ . We highlight this detail because we use these states as the way to incorporate additional information into the MvAttLSTM model.

#### 3.3.3. Attention layer

The BiLSTM Layer has two major problems. Since the meaning of the message cannot be encoded in one fixed-size vector, there is some information loss. Hence, the performance of this type of models for representation learning decreases when the length of inputs become large. Another concern is that LSTMs aggregate information word-by-word in sequential order, but there is no explicit mechanism to make inferences over the structure and modeling relations among tokens. To overcome these limitations, the output of the BiLSTM Layer  $H^1 \in \mathbb{R}^{2 \times d_h \times N}$  is fed into an Attention Layer. This layer helps BiLSTM in deciding which parts of the sequence pay more interest. In this work, we investigate the performance of two attention mechanisms in our model: self-attention and multi-head attention [39].

Self-attention mechanisms can capture the explicit and latent relations among words beyond their sequential order. While attention mechanism [159] allows the outputs for attending some parts of the inputs. Self-attention also allows the inputs for interacting each other, hence amplifying the importance of each one plays in determining the meaning of others. Moreover, is it beneficial for discovering word relations which can be crucial for understanding ironic and satirical texts such as, oppositions and incongruities. Given the matrices  $A$ ,  $B$  and  $C$ , mathematically self-attention is formulated as follows:

$$\text{Att}(A, B, C) = \text{Attention}(AW^Q, BW^K, CW^V) \quad (9)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (10)$$

Where  $W^Q \in \mathbb{R}^d$ ,  $W^K \in \mathbb{R}^d$ ,  $W^V \in \mathbb{R}^d$  are the projection matrices for the query  $Q$ , key  $K$  and value  $V$ . In the case of self-attention, the matrices  $A$ ,  $B$  and  $C$  are the same. Thus, given the output  $H^1$  of the BiLSTM layer, the new sequence of weighted hidden states  $H^2 \in \mathbb{R}^{2 \times d_h \times N}$  which is the output of the Attention layer is computed by Eq. (11) as:

$$H^2 = \text{Att}(H^1, H^1, H^1) \quad (11)$$

<sup>3</sup> <https://crsccardellino.github.io/SBWCE/>

In [39] another attention mechanism was introduced. The multi-head attention mechanism uses multiple individual attention functions (heads) for obtaining different contexts and paying attention simultaneous to distinct aspects in the sequences. This can jointly pay attention to information from different representation sub-spaces at different positions. Like in self-attention, the attention function takes as input a matrix for the query  $A$ , a matrix for the keys  $B$  and a matrix for values  $C$ . The multi-head attention model first transforms  $A$ ,  $B$  and  $C$  into  $\mathbb{C}$  sub-spaces, with different, trainable linear projections:

$$\text{MultiHead}(A, B, C) = [\text{head}_1, \text{head}_2, \dots, \text{head}_r] * W^0 \quad (12)$$

$$\text{head}_c = \text{Attention}(AW_c^Q, BW_c^K, CW_c^V) \quad (13)$$

Where  $W_c^Q \in \mathbb{R}^{d \times d_k}$ ,  $W_c^K \in \mathbb{R}^{d \times d_k}$ ,  $W_c^V \in \mathbb{R}^{d \times d_v}$  are projection matrices for the inputs  $A, B, C$  with respect to  $\text{head}_c$ , and  $W^0 \in \mathbb{R}^{r \times d_k \times d}$ . The parameter  $r$  is the number of heads for the multi-head attention mechanism; and  $\text{head}_c \in \mathbb{R}^{N \times d_k}$  is the output of the  $c$ th head. Notice that, for each  $\text{head}_c$ , the weights of  $W_c^Q, W_c^K, W_c^V$  are independently learned during the training phase. The attention for each head  $c$  (see Eq. (13)), like in self-attention, is computed by the formula in Eq. (10). Thus, given the output of the BiLSTM layer  $H^1$ , the output of the multi-head attention is computed as follows:

$$H^2 = \text{MultiHead}(H^1, H^1, H^1) \quad (14)$$

### 3.3.4. LSTM layers

Even when, it is not theoretically clear what is the additional power gained by the deeper recurrent architectures, it was observed empirically that a deep LSTM works better than shallower ones in some tasks [100,160]. Taking this into account, on the output  $H^2$  of the Attention Layer we stacked two layers of LSTMs to deep contextualize the previously learned representation. This means that the output of the first LSTM layer is given as input to the second one. The two LSTM layers are defined as follows:

$$H^3 = \text{LSTM}(H^2) \quad (15)$$

$$H^4 = \text{LSTM}(H^3) \quad (16)$$

Where  $H^3 \in \mathbb{R}^{d_{k1} \times N}$ ,  $H^4 \in \mathbb{R}^{d_{k1} \times N}$  are the outputs of the first and second LSTM layer and  $d_{k1}$  is the number of hidden neurons into LSTM cells. The last LSTM layer output the hidden representation of the text. Particularly, we only consider the last hidden state ( $h_N^4$ ) into the matrix  $H^4$ . Let us redefine it as  $h_{last}^4$  henceforth. Like in the BiLSTM layer, the initial hidden state and cell memory of both LSTMs are set to 0,  $h_0^3 = c_0^3 = \vec{0}$ , and  $h_0^4 = c_0^4 = \vec{0}$ .

### 3.3.5. Fusion strategies

Our model aiming at improving irony and satire detection in Spanish by incorporating multiple views into the MvAttLSTM. For this purpose, we investigate two strategies for passing the views to the model. The first one, *Early Fusion* method, aiming at enriching the representation learned by the LSTMs with additional knowledge using the last dense layers. The second one, *Contextual Fusion* method, which aims to condition the learning process of the LSTMs with prior knowledge injected in the initial cell memory. Next, we give details about both strategies.

#### Early Fusion

The main idea behind this strategy is to separately learn different features spaces from the training data. These feature spaces (views) capture distinct characteristics of the same texts. We aim to jointly use these representations to retain discriminant information while reducing the redundant one. The overall architecture of the MvAttLSTM with *Early fusion* is showed in Fig. 1.

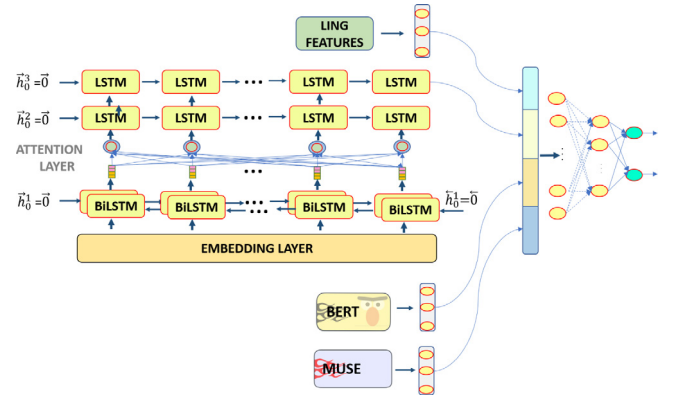


Fig. 1. MvAttLSTM: Multi-view informed attentive LSTM deep neural network using an early fusion strategy.

Firstly, let us define  $H_{\text{LING}}$ ,  $H_{\text{MUSE}}$  and  $H_{\text{BERT}}$  as the Linguistic-view, MUSE-view and BERT-view respectively (see Section 3.2). These views differ from each other in the way by which were learned and the number of features used for encoding the text. Thus, we pass each view to a dense layer in order to reduce and unify the views' dimensionality using the Eq. (17), (18), (19):

$$g_l(H_{\text{LING}}) = \sigma(W^l H_{\text{LING}} + b^l) \quad (17)$$

$$g_m(H_{\text{MUSE}}) = \sigma(W^m H_{\text{MUSE}} + b^m) \quad (18)$$

$$g_b(H_{\text{BERT}}) = \sigma(W^b H_{\text{BERT}} + b^b) \quad (19)$$

Where  $W^l, W^m, W^b$  and  $b^l, b^m, b^b$  are parameters of the model to be learned during the training process, and  $\sigma$  is the sigmoid function. After having reduced representations for each view ( $g_l, g_m, g_b$ ), then we concatenate them with a deep representation learned by the attentive LSTM based architecture  $h_{last}^4$  (Eq. (20)). Later, the merged representation denoted as  $F_0$  is fed into a dense layer with sigmoid activation for fusing all views into a new non-linear space using Eq. (21). Finally, the output of this layer denoted as  $F_1$  is a multi-view encoding of the texts, and it is fed into a feed-forward neural network for the final classification of the texts in ironic vs. non-ironic or satirical vs. non-satirical.

$$F_0 = \text{Concat}(g_l, g_m, g_b, h_{last}^4) \quad (20)$$

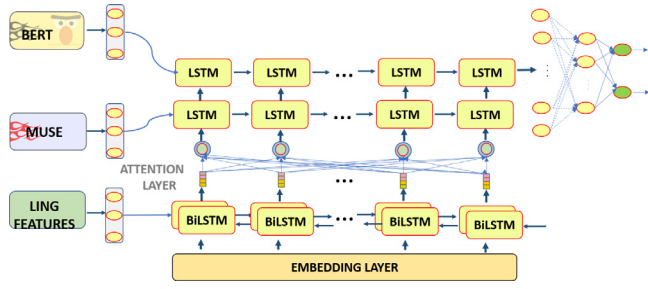
$$F_1 = \sigma(W^0 F_0 + b^0) \quad (21)$$

#### Contextual Fusion

In this strategy, we enrich our MvAttLSTM with additional external knowledge to take advantage of the initial memory cell in the LSTMs. We experiment with a strategy similar to the conditional encoding model introduced in [161] for the task of recognizing textual entailment and applied later in [82] for modeling conversation context for improving sarcasm detection. Conversely, to the approach presented by Ghosh et al. [82], we do not learn contextual information by using LSTMs, instead, we learn independently three distinct views with the aim to capture syntactic, semantic, and pragmatics aspects used in ironical and satirical texts. In Fig. 2 is showed the overall architecture of our MvAttLSTM model using the *Contextual Fusion* strategy. As can be observed, the learning process of each LSTM is conditioned to prior information passed to the initial memory cell.

Like in the *Early fusion* strategy, firstly each view is fed into a dense layer with non-linear activation, specifically using the Eq. (17), (18), (19). Later, each reduced representation ( $g_l, g_m, g_b$ )





**Fig. 2.** MvAttLSTM: Multi-view informed deep attentive LSTM neural network using contextual fusion strategy.

is used to inform the LSTM in the MvAttLSTM. The initial memory cell and hidden states of the LSTMs are used as input to pass the prior knowledge of each view as defined in Eq. (22), (23), (24). Where  $c_0^1$  and  $h_0^1$  are the initial memory cells and hidden states of the BiLSTM. And,  $h_0^3$ ,  $c_0^3$ ,  $h_0^4$  and  $c_0^4$  are the initial memory cell and hidden states of the second and last LSTM respectively. The order in which each view is assigned to the LSTMs was empirically defined. We decide to introduce low-level linguistic features for reinforcing the BiLSTM layer which aims at capturing language generalization. In the second LSTM, we propose to introduce the MUSE-view for incorporating a high-level semantic representation to encode the global meaning of the text. Finally, in the last LSTM, we introduce the BERT-view to incorporate semantics and pragmatics abstractions useful for the task to solve, considering that in this view the BERT model is tuned using the same training data available for the task. This introduces a task-dependent bias in the language representation learned by the original model:

$$c_0^1 = h_0^1 = g_l(H_{LING}) \quad (22)$$

$$h_0^2 = c_0^2 = g_m(H_{MUSE}) \quad (23)$$

$$h_0^3 = c_0^3 = g_b(H_{BERT}) \quad (24)$$

Notice that in this case, the final multi-view representation of the text  $F_1$  is the same that the last hidden state of the last LSTM layer  $h_{last}^4$  in our MvAttLSTM, hence  $F_1 = h_{last}^4$ . And, we pass this representation into the feed forward neural network for the classification of the texts in ironic vs. non-ironic or satirical vs. non-satirical.

### 3.3.6. Feed-forward layer for final classification

For achieving the final classification we fed the multi-view encoding of the text  $F_1$  into a Dropout layer to prevent the model's over-fitting. Subsequently, the output of the Dropout layer  $F_2$  is passed to a dense layer with ReLU activation, and finally, the output of this layer  $F_3$  is given as input to another dense layer with two neurons, but this time with the *softmax* function for the prediction:

$$F_2 = \text{Dropout}(F_1) \quad (25)$$

$$F_3 = \max(0, W^2 F_2 + b_2) \quad (26)$$

$$O = \text{softmax}(W^3 F_3 + b_3) \quad (27)$$

The MvAttLSTM model can be trained in an end-to-end way by the back-propagation method, and we use categorical cross-entropy as the loss function. This function can be observed in Eq. (28), where  $\mathcal{D}$  is the dataset,  $\mathcal{L}$  is the loss function,  $f$  is our

model parameterized by  $\theta$  and  $\mathbb{G} = \{1, 0\}$  is the set of labels in the task.

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathcal{D}}[\mathcal{L}(f(x, \theta), y)] = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{|\mathbb{G}|} y_{ij} * \log(f(x_i, \theta)_j) w_j \quad (28)$$

## 4. Experiments and results

### 4.1. Datasets description

In order to validate our proposed model for irony and satire detection in short texts written in distinct Spanish variants, we used three corpora, one for irony detection, and two for satire detection. Moreover, we also tested the robustness of our model on humor recognition on another corpus. They have been extensively used with the aim of training and evaluating state-of-the-art systems for irony, satire and humor detection in Spanish.

#### Irony Corpus

For what concerns irony detection, we decided to use the corpus proposed in the *IroSvA'19* shared task [56]. This is the first public available corpus for irony detection in Spanish. The *IroSvA'19* shared task, framed in the Iberian Languages Evaluation Forum (*IberLEF'19*)<sup>4</sup> and co-located within *SEPLN 2019*<sup>5</sup> aimed at investigating whether a short message, written in Spanish, is ironic or not within a given context. For that, three corpora with short texts from Spain, Mexico and Cuba were proposed with the purpose of exploring the way irony changes in Spanish variants. In particular, the Castilian and Mexican corpora consist of ironic tweets about 10 controversial topics for Spanish and Mexican users. In the case of the Cuban corpus, it consists of ironic news comments which were extracted from 113 controversial news about social, economic, and political issues concerning the Cuban people. It is worthy to notice that, for each text a context is provided, consisting of a short description about the topic, which defines its scope. The distribution of the texts is showed in Table 1.

As can be observed in Table 1 all subcorpora are composed of 3000 texts split into 2400 and 600 texts for training and testing respectively. The training set is separated into 800 ironic and 1600 non-ironic texts, whereas the testing partition is divided into 200 ironic and 400 non-ironic texts. Notice that both training and testing sets maintain the ratio of 2/3 vs. 1/3 between non-ironic and ironic text. In order to assess the performance of the systems, the evaluation metrics used by the organizers were precision (P), recall (R), and F1 score. These metrics were calculated per class and macro-averaged. Due to the imbalance between the non-ironic and ironic classes, the macro-averaged F1 score was used as the overall metric to rank the participating systems.

#### Satire Corpora

For satire detection, we used the corpora proposed in [55] and [86]. Both corpora were created using a self-annotation strategy. Specifically, Barbieri et al. [86] retrieved tweets from popular satirical news accounts and from legitimate news sources in three languages: Spanish, English, and Italian. In this work, we were interested in the Spanish subset of this corpus, and we refer to this as *Barbieri'15-es* henceforth. The Spanish tweets (Castilian variant) were gathered from two satirical Twitter's accounts *El Mundo Today* and *El Jueves* whereas non-satirical tweets were retrieved from the legitimate newspaper Twitter's accounts *El Mundo* and *El Pais*. Later, a shallow cleaning process was carried out on data for filtering those tweets that were not relevant to satire analysis. As can be observed in Table 2, the corpus is

<sup>4</sup> <https://sites.google.com/view/iberlef-2019>

<sup>5</sup> <http://www.hitzeus/sepln2019/?language=es>

**Table 1**  
IroSvA'19 distribution for ironic and non-ironic classes.

Corpus	Variant	Training			Testing		
		Non-Ironic	Ironic	Total	Non-Ironic	Ironic	Total
IroSvA'19	Castilian (es)	1600	800	2400	400	200	600
	Mexican (mx)	1600	800	2400	401	199	600
	Cuban (cu)	1600	800	2400	400	200	600

composed of 10888 uniformly distributed in 5444 satirical tweets 5444 non-satirical ones.

The corpus introduced in [55] was guided by the same methodology presented in [86]. The most salience difference relies on the study of satirical tweets in two variants of the Spanish. Particularly, tweets from Mexican and Castilian Twitter's accounts were retrieved. For investigating how satire is realized in Mexican tweets, data from four Mexican Twitter accounts were retrieved. The satirical tweets were obtained from *El Deforma* and *El Dizque* satirical accounts whereas the non-satirical tweets were gathered from legitimate newspaper accounts *El Universal* and *Excelsior*. We refer to this subset of data as *Salas'17-mx* henceforth. The tweets in the Castilian corpus of [55], *Salas'17-es* henceforth, were retrieved using the four Twitter's accounts proposed in [86]. It is important to note that even when the Twitter's accounts used to obtain the tweets were the same, the tweets in each collection are different.

An automatic cleaning process was carried out on the data. Specifically, retweets, duplicates, tweets only with URLs, and tweets written in a language other than Spanish were removed. Moreover, a manual inspection was performed in order to ensure that the tweets obtained were relevant for satire detection. In Table 2 can be observed that both corpora contain 5000 tweets, which are uniformly distributed in 2500 satirical and 2500 non-satirical ones. The different characteristics of the *Barbieri'15-es* and *Salas'17-es* will allow us to validate the robustness of our *MvAttLSTM* model.

#### Humor Corpus

For further investigating the robustness of our model we decided to evaluate it on humor recognition in Spanish. We considered the corpus proposed in the *HAHA'19* shared task [162,163] organized at *IberLEF'2019* and co-located within *SEPLN 2019*. Two subtasks were proposed, one for humor binary classification (*Humor Recognition*) and another for predicting how funny is a tweet into 5-star ranking (*Funniness Score Prediction*), considering that the tweets present humorous content. The organizers provided a human-annotated corpus of 30000 Spanish tweets separated into 24000 for training and 6000 for testing. The training subset consists of 9253 humorous and 14747 non-funny tweets, whereas the testing subset consists of 2342 humorous tweets and 3658 non-funny ones. In Table 3 we summarize the distribution of the tweets within *HAHA'19*. Taking into account the scope of this work, we are only interested in the first subtask. As can be noted, in both training and testing subsets the distribution of the classes are slightly unbalanced, hence a difficulty is added to the learning algorithm. The performance metrics used to rank the participated systems in the *Humor recognition* subtask were F1 score for the *humorous* class and accuracy (*Acc*).

#### 4.2. Experimental setting

We use the same architecture for all tasks, but we calibrated the hyper-parameters independently for each corpus. Specifically, we defined the number of hidden neurons in the BiLSTM layer to 64, the number of hidden neurons in the last two LSTM layer to 128, the maximum number of epochs and length of the sequence  $ep = 50$  and  $N = 50$  respectively. For the remainder of hyperparameters, we experimented with distinct values.

Specifically, we defined the search space as follows: batch size  $batch \in [32, 64, 128, 256]$ , dropout  $dp \in [0.25, 0.30, 0.35, 0.4]$ , optimizer update rules  $op \in (adam, rmsprop)$ , learning rate  $lr \in [1 \times 10^{-2}, 1 \times 10^{-3}, 1 \times 10^{-4}]$ . When the model uses multi-head attention we evaluated distinct number of heads  $h \in [2, 4, 8, 16]$ . In an intent to prevent the overfitting in the training step, an early stopping with the patience of 10 epochs was used as a stopping criterion. We explored the search space by means of the Grid Search strategy. Analyzing the best hyperparameters obtained for each corpus and model, we observed that our models are sensitive to the hyperparameter setting. Particularly, we noted that the learning rate  $lr = 1 \times 10^{-2}$  achieved the best performance across all corpora. However, the remainder of the hyperparameters has a distinct behavior. Roughly speaking, we appreciated that the *MvAttLSTM<sub>self</sub>Contextual* and *MvAttLSTM<sub>multi</sub>Contextual* are the most sensitives models. One possible reason for that is the small number of ironic examples in each corpus which makes the model generalization more complex. In the case of the *MvAttLSTM<sub>self</sub>Early* model, it was observed that it performs well on large corpora using short batches ( $batch_{size} = 32$ ) whereas *MvAttLSTM<sub>multi</sub>Early* requires longer batches ( $batch_{size} = 128$ ). Also, we observed that 4 heads of attention were enough to achieve good results. Concerning the optimizer, generally, the *Adam* rule obtained the best result in 9 settings out of 16. The best hyperparameters for each corpus and model are summarized in Appendix B.

In order to evaluate the performance of the distinct settings of our model, we define a baseline method (*Bert-baseline*). Concretely, we use the mBERT method fine-tuned on each corpus separately. For that, we adopt the same hyperparameter and tuning strategy proposed in Section 3.2.3.

Regarding the *Linguistic-view*, we experimented with distinct views to investigate whether some groups of features are more feasible than others to detect irony cues. In this sense, we defined three views for considering affective information: *Aff\_All*, *Aff\_Emo*, *Aff\_App*. In *Aff\_All* we considered all features related to polarity, emotions (categorical, and dimensional), and attitudes. Whereas in *Aff\_Emo* we only used those features related to emotions, and in *Aff\_App* we only use the attitude words. The features that capture polarity oppositions were included in the group *Contrast*. We evaluated two groups (*LIWC*, *Sverb*) based on the psycholinguistic dimensions in the LIWC dictionary and the semantic classes of verbs in the ADDESE lexicon<sup>6</sup> respectively. Moreover, other groups of features were obtained by using a feature selection method, specifically the Wilcoxon rank-sum test [164] were explored. With this statistical test, the features were ranked considering their *p*-value, and three groups were defined. In the groups *W\_64*, *W\_128* and *W\_All* we considered the subsets of 64 and 128 best-ranked features and all features with *p*-value  $\leq 0.05$ . Finally, a group with all the linguistic features *LingAll* was considered.

#### 4.3. Results in irony detection in spanish variants

In this section, we present an exhaustive evaluation of distinct settings of the *MvAttLSTM* model in the task of irony detection.

<sup>6</sup> <http://adesse.uvigo.es/data/clases.php>

**Table 2**  
Distribution for satirical and non-satirical classes in *Salas'17* and *Barbieri'15* datasets.

Corpus	Variant	Data		
		Non-Satirical	Satirical	Total
<i>Salas'17</i>	Castilian (es)	2500	2500	5000
	Mexican (mx)	2500	2500	5000
<i>Barbieri'15</i>	Castilian (es)	5444	5444	10888

**Table 3**  
HAHA'19 distribution for humorous and non-humorous classes.

Corpus	Language	Training			Testing		
		Non-Humorous	Humorous	Total	Non-Humorous	Humorous	Total
HAHA'19	Spanish	14747	9253	24000	3658	2342	6000

Our first experiment aimed at investigating the impact of different types of linguistic views on the model. For that, we analyzed what subsets of features are most relevant to irony detection. The second aspect that we considered relevant to explore was the impact of the fusion strategies to inform the model (*Early* vs. *Contextual*) and the attention mechanism used by the MvAttLSTM model (*self* vs. *multi-head*). Lastly, we investigated the impact of each proposed view. For that, we ignored one view and fed the other two into the MvAttLSTM model.

To evaluate the effectiveness of our proposal in each experiment we computed the F1 score for the two classes ( $F1_{iro}$  and  $F1_{no-iro}$ ), along with their macro-averaged and micro-averaged versions of F1 ( $F1_{Micro}$  and  $F1_{Macro}$ ). We split the training data into 80% and 20% for training and validation purposes and evaluated the generalization of our model on the official test provided by the organizers. The results obtained for *IroSvA'19* corpus on the test dataset in the three Spanish variants are shown in Table 4. We only included the results using the *Contextual fusion* strategy for the Castilian (es), the Mexican (mx) and the Cuban (cu) variants due to another fusion strategy achieved worse results in the three variants.

It can be observed in Table 4 that the model  $MvAttLSTM_{Contextual}^{multi}$  obtained slightly better results than  $MvAttLSTM_{Contextual}^{self}$  in the three variants (es, mx and cu) for all the evaluation metrics. Concretely, for the Castilian variant, the best results were obtained by  $MvAttLSTM_{Contextual}^{multi}$  using all views, but in the case of *Linguistic-view* only considering emotional *Aff\_Emo* or attitudinal features *Aff\_App*. Moreover, the model  $MvAttLSTM_{Contextual}^{self}$  achieves competitive results but considering the views  $W_{128}$  or  $W_{64}$ . Regarding the Mexican variant, both models  $MvAttLSTM_{Contextual}^{multi}$  and  $MvAttLSTM_{Contextual}^{self}$  obtained the best results when the *Aff\_All* and *Aff\_Emo* views are used respectively. Also, it is important to notice that the second better results for each model are achieved when the views *LIWC* and *Aff\_App* are considered. In the case of the Cuban variant, both models obtain similar results. However,  $MvAttLSTM_{Contextual}^{multi}$  using all linguistic features *LingAll* slightly outperform  $MvAttLSTM_{Contextual}^{multi}$  with the linguistic view *Aff\_All*.

To sum up, we found that some subsets of features are the most relevant in our model for irony detection in the three variants; particularly, those related to affective information such as *Aff\_Emo* and *Aff\_App*. This fact indicates the discriminatory property of the emotional dimensions in SenticNet and the emotional categories in Spanish Emotions Lexicon (SEL) [165]. Also, the attitude-based features obtained from Appraisal Lexicon (LAM) [166] were relevant. This result is in line with the findings presented in [75] which investigated the role of affective information in irony detection using machine learning models. Furthermore, we found that the *Bert-baseline* method performs significantly worse than the MvAttLSTM model in the three variants. One possible explanation for that is the small number of ironic examples in the training dataset that make more complex the learning process.

In a second direction, we investigated the importance of each proposed view (*Bert-view*, *Muse-view* and *Linguistic-view*) on the performance of MvAttLSTM. For that, we evaluated the model ignoring one view and including the remaining two. In this experiment, the *Linguistic-view* (*Ling*) represents the subsets of features that achieved the best  $F1_{Macro}$  (see Table 4). Notice that *Ling* is different for each MvAttLSTM setting and dataset. The results obtained are summarized in Table 5.

As can be shown in Table 5, the best  $F1_{Macro}$  in all corpora was obtained when all views were used together. This fact confirms that informing our model with the proposed views helps the model to detect irony. However, we observed that for the Castilian variant, ignoring *Muse-view* caused the most significant drop in performance of  $MvAttLSTM_{Contextual}^{self}$  whereas omitting the *Bert-view* produced the worse performance in  $MvAttLSTM_{Contextual}^{multi}$ . In the case of the Mexican variant, both settings of  $MvAttLSTM$  achieved the worse performance when the *Muse-view* was removed from the model. However, for the Cuban variant, we found that the model drops its performance when the *Linguistic-view* was omitted. Analyzing the results in Tables 4 and 5 together, the linguistic view (*LingAll*) was found to produce a lower F1-macro than when no linguistic features are introduced in the model. In this sense, we considered that a deeper analysis would be necessary to explain the reasons for the negative result achieved when including all the linguistic features whether it is due to noisy features or the fusion strategy used to feed this view into the model. Further efforts need to be made for investigating why the attention mechanisms (*self* vs. *multi*) attend different linguistic views for obtaining better effectiveness.

Following, we present a comparison of our best results on the three corpora with other state-of-the-art systems. In Table 6 we show how the results of the participating systems in the *IroSvA'19* shared task ranked according to the official evaluation measure  $F1_{Macro}$  average. It is important to highlight that the participants were not restricted to submit the same system for each corpus. Thus, the F1-AVG means the average of the results of the team instead of evaluating the performance of one model on the three corpora. Our best model for the Castilian variant  $esMvAttLSTM_{Contextual}^{multi}$  outperforms the results achieved by ELiRF\_UPV [45,58] and CIMAT [59] on the Castilian and Cuban corpora. However, our model drops its performance on the Mexican corpus. Regarding our best model for the Mexican variant  $mxMvAttLSTM_{Contextual}^{multi}$ , it outperforms the results obtained by ELiRF\_UPV and CIMAT on all corpora. The  $cuMvAttLSTM_{Contextual}^{multi}$  model achieved better results than ELiRF\_UPV and CIMAT on the Cuban corpus but dropped its effectiveness on the Castilian and Mexican variants. It is important to remark that ELiRF\_UPV is based on a deep learning model; particularly it proposed a simplification of the BERT model, and CIMAT proposed a combination of deep learning-based representations with n-gram features. From an overall point of view, our proposed models are placed in the



**Table 4**  
MvAttLSTM for irony detection in IroSvA'19.

Views	$F1_{iro}$	$F1_{no-iro}$	$F1_{Micro}$	$F1_{Macro}$	$F1_{iro}$	$F1_{no-iro}$	$F1_{Micro}$	$F1_{Macro}$
	<i>MvAttLSTM<sup>Contextual-self</sup></i> IroSvA'19-es				<i>MvAttLSTM<sup>Contextual-multihead</sup></i> IroSvA'19-es			
Bert+Muse+Aff_All	0.487	0.821	0.734	0.654	0.596	0.835	0.766	0.716
Bert+Muse+Aff_Emo	0.619	0.815	0.751	0.717	<b>0.668</b>	<b>0.842</b>	<b>0.786</b>	<b>0.755</b>
Bert+Muse+Aff_App	0.538	0.82	0.741	0.679	<b>0.651</b>	<b>0.835</b>	<b>0.776</b>	<b>0.743</b>
Bert+Muse+Contrast	0.59	0.833	0.762	0.712	0.587	0.83	0.759	0.708
Bert+Muse+LIWC	0.549	0.801	0.724	0.675	0.626	0.831	0.767	0.728
Bert+Muse+SVerb	0.576	0.798	0.726	0.687	0.627	0.824	0.761	0.726
Bert+Muse+W_64	<b>0.629</b>	<b>0.832</b>	<b>0.769</b>	<b>0.731</b>	0.583	0.834	0.762	0.708
Bert+Muse+W_128	<b>0.656</b>	<b>0.835</b>	<b>0.777</b>	<b>0.746</b>	0.642	0.82	0.761	0.731
Bert+Muse+W_All	0.595	0.811	0.742	0.703	0.55	0.808	0.731	0.679
Bert+Muse+LingAll	0.553	0.802	0.726	0.678	0.589	0.831	0.761	0.710
Bert-baseline	0.302	0.741	0.594	0.521	0.302	0.741	0.594	0.521
	<i>MvAttLSTM<sup>Contextual-self</sup></i> IroSvA'19-mx				<i>MvAttLSTM<sup>Contextual-multihead</sup></i> IroSvA'19-mx			
Bert+Muse+Aff_All	0.516	0.79	0.708	0.654	<b>0.647</b>	<b>0.817</b>	<b>0.759</b>	<b>0.732</b>
Bert+Muse+Aff_Emo	<b>0.642</b>	<b>0.821</b>	<b>0.761</b>	<b>0.732</b>	0.508	0.785	0.701	0.647
Bert+Muse+Aff_App	<b>0.644</b>	<b>0.794</b>	<b>0.739</b>	<b>0.719</b>	0.585	0.774	0.708	0.68
Bert+Muse+Contrast	0.601	0.812	0.744	0.706	0.506	0.792	0.708	0.649
Bert+Muse+LIWC	0.506	0.792	0.708	0.649	<b>0.611</b>	<b>0.826</b>	<b>0.759</b>	<b>0.718</b>
Bert+Muse+SVerb	0.564	0.765	0.694	0.664	0.596	0.828	0.759	0.712
Bert+Muse+W_64	0.529	0.788	0.708	0.659	0.515	0.786	0.703	0.65
Bert+Muse+W_128	0.567	0.777	0.706	0.672	0.591	0.591	0.689	0.670
Bert+Muse+W_All	0.512	0.79	0.706	0.651	0.599	0.790	0.724	0.695
Bert+Muse+LingAll	0.606	0.809	0.743	0.707	0.469	0.808	0.718	0.638
Bert-baseline	0.293	0.771	0.611	0.532	0.293	0.770	0.611	0.532
	<i>MvAttLSTM<sup>Contextual-self</sup></i> IroSvA'19-cu				<i>MvAttLSTM<sup>Contextual-multihead</sup></i> IroSvA'19-cu			
Bert+Muse+Aff_All	<b>0.604</b>	<b>0.807</b>	<b>0.741</b>	<b>0.706</b>	0.534	0.796	0.716	0.665
Bert+Muse+Aff_Emo	0.574	0.809	0.736	0.691	0.56	0.796	0.721	0.678
Bert+Muse+Aff_App	0.53	0.803	0.723	0.666	0.563	0.793	0.719	0.678
Bert+Muse+Contrast	0.556	0.827	0.751	0.692	0.557	0.793	0.718	0.675
Bert+Muse+LIWC	0.536	0.8	0.721	0.668	0.577	0.788	0.718	0.683
Bert+Muse+SVerb	0.546	0.819	0.741	0.683	0.568	0.79	0.718	0.679
Bert+Muse+W_64	0.554	0.792	0.716	0.673	<b>0.596</b>	<b>0.816</b>	<b>0.748</b>	<b>0.706</b>
Bert+Muse+W_128	0.556	0.791	0.716	0.674	0.529	0.8	0.719	0.665
Bert+Muse+W_All	<b>0.582</b>	<b>0.819</b>	<b>0.748</b>	<b>0.701</b>	0.473	0.825	0.738	0.649
Bert+Muse+LingAll	0.517	0.804	0.721	0.661	<b>0.602</b>	<b>0.817</b>	<b>0.749</b>	<b>0.709</b>
Bert-baseline	0.472	0.803	0.693	0.638	0.472	0.803	0.693	0.638

**Table 5**

The impact of the views on MvAttLSTM for irony detection. The ignored view is denoted by (×) symbol whereas the included views are denoted by (✓) symbol.

Model	Ling	Muse	Bert	$F1_{iro}$	$F1_{no-iro}$	$F1_{Micro}$	$F1_{Macro}$
<i>IroSvA'19-es</i>							
Contextual-self	×	✓	✓	0.627	0.835	0.771	0.730
	✓	×	✓	0.593	0.804	<b>0.736</b>	<b>0.699</b>
	✓	✓	×	0.611	0.819	0.753	0.715
Contextual-multi	×	✓	✓	0.611	0.827	0.761	0.719
	✓	×	✓	0.607	0.788	0.724	0.697
	✓	✓	×	0.592	0.780	<b>0.714</b>	<b>0.686</b>
<i>IroSvA'19-mx</i>							
Contextual-self	×	✓	✓	0.546	0.824	0.746	0.685
	✓	×	✓	0.511	0.788	<b>0.704</b>	<b>0.650</b>
	✓	✓	×	0.620	0.804	0.741	0.712
Contextual-multi	×	✓	✓	0.530	0.776	0.696	0.653
	✓	×	✓	0.503	0.793	<b>0.708</b>	<b>0.648</b>
	✓	✓	×	0.565	0.800	0.726	0.682
<i>IroSvA'19-cu</i>							
Contextual-self	×	✓	✓	0.557	0.793	<b>0.718</b>	<b>0.675</b>
	✓	×	✓	0.559	0.803	0.728	0.681
	✓	✓	×	0.576	0.824	0.751	0.700
Contextual-multi	×	✓	✓	0.528	0.805	<b>0.724</b>	<b>0.667</b>
	✓	×	✓	0.570	0.786	0.714	0.678
	✓	✓	×	0.520	0.827	0.746	0.674

**Table 6**  
Comparison with state-of-the-art methods for irony detection in Spanish variants (IroSvA'19).

Ranking	Team	IroSvA'19-es $F1_{Macro}$	IroSvA'19-mx $F1_{Macro}$	IroSvA'19-cu $F1_{Macro}$	IroSvA'19 $F1_{AVG}$
(*)	$mxMvAttLSTM_{Contextual}^{multi}$	<b>0.716</b>	<b>0.732</b>	<b>0.665</b>	<b>0.704</b>
(**)	$esMvAttLSTM_{Contextual}^{multi}$	<b>0.755</b>	0.674	<b>0.678</b>	<b>0.702</b>
(***)	$cuMvAttLSTM_{Contextual}^{multi}$	0.710	0.638	<b>0.709</b>	<b>0.685</b>
1st	ELiRF-UPV	0.717	0.680	0.653	0.683
2nd	CIMAT	0.645	0.671	0.660	0.659
3th	JZaragoza	0.661	0.67	0.616	0.649
4th	ATC	0.651	0.645	0.594	0.630
...	...	...	...	...	...
14th	UO	0.511	0.489	0.499	0.499

first positions in the ranking. This fact shows the effectiveness of our model in addressing the problem of irony detection in multiple variants of Spanish.

#### 4.4. Results in satire detection in spanish variants

Irony and satire are both indirect forms of communication that are strongly related to each other. These forms aim at communicating in implicit ways complex meanings which often aim at criticizing, offending or hurting a victim. The major differences between them are based on the intention of the author and the linguistic resources used to effectively communicate the real meaning. In this section, we present an evaluation of our model on two corpora of satirical tweets (*Salas'17* and *Barbieri'15*) for analyzing the feasibility of our model for satire detection in two Spanish variants (Castilian, and Mexican). Conversely, to the *IroSvA'19* corpus, these corpora are not explicitly divided into train and test, then we use 5-fold cross-validation to compute the generalization capability of our model in each corpus. In each iteration 80% of the data was used for training meanwhile the remainder 20% was considered for testing purpose. Also, to calibrate the hyperparameters of the model, the training set was split into two subsets 90% to train the model and 10% for validation purpose. For each corpus, the hyperparameters were tuned independently.

The results of our model on *Salas'17* and *Barbieri'15* are summarized in Table 7. In this table only the results using the *Early fusion* strategy are reported, due to the other fusion method obtained relatively worse results. At a first glance, in Table 7 can be observed that both settings of the model  $MvAttLSTM_{self}^{Early}$  and  $MvAttLSTM_{multi}^{Early}$  achieved similar results, even when the model which uses multi-head attention showed a slight improvement at the expense of more trainable parameters.

Concretely, for the Castilian variant in both corpora *Barbieri'15* and *Salas'17-es* the model with self-attention  $MvAttLSTM_{self}^{Early}$  obtained good results when the *Linguistic-view* is used to inform the model. Particularly, appraisal features (*Aff\_App*) was the most relevant for satire detection in *Salas'17-es* and the second better in *Barbieri'15*. Moreover, the features obtained by using the Wilcoxon test showed a good performance, resulting the second more relevant *W\_64* and *W\_All* in *Salas'17-es*, and *W\_128* the most relevant in *Barbieri'15-es*. In the case of  $MvAttLSTM_{multi}^{Early}$  the best results for both Castilian corpora were achieved using the 64 best-ranked features *W\_64* according to the Wilcoxon test. Regarding the results on the Mexican variant, they were different to those achieved on the Castilian tweets. Particularly, the model  $MvAttLSTM_{multi}^{Early}$  perform better when the features related to polarity opposition *Contrast* were used whereas  $MvAttLSTM_{multi}^{Early}$  obtained the best performance when psycho-linguistic *LIWC* features were considered. Moreover, it can be observed that *Bert-baseline* obtained very competitive results on the three corpora.

This fact, confirmed that the representations leaned by the BERT model are good enough to discriminate between satirical and non-satirical tweets.

In a second direction, we aim at exploring the role of the three views proposed to inform  $MvAttLSTM$ . For that, we evaluated the model ignoring one view and including the remaining two. In this experiment, the linguistic view (*Ling*) represents the subsets of features that achieved the best  $F1_{macro}$  (see Table 7). As can be shown in Table 8 the best  $F1_{Macro}$  in all corpora was obtained when all views were used together. In general, we observed that for the three variants, ignoring *Bert-view* caused the most significant drop in performance of both settings of  $MvAttLSTM$ .

In order to have a comparison with the performance obtained by other methods proposed in the literature, we compare our models with the results presented in [54,55]. According to our knowledge, these works are the only two that addressing the problem of satire detection in Spanish. In Table 9 we compare our model with three methods (*SMO+LIWC-ALL*, *BayesNet+LIWC-ALL*, *J48+LIWC-ALL*) proposed in [55] and three methods (*SVM+W-B*, *SVM+Intrinsic*, *SVM+ALL*) introduced in [54] for satire detection. The methods proposed in [55] are based on machine learning combined with hand-crafted features, particularly features derived from *LIWC*. The major difference among these methods is the machine learning algorithm used. The methods were evaluated using precision ( $P_{sat}$ ) recall ( $R_{sat}$ ) and  $F1$  score ( $F1_{sat}$ ) on the positive class (satirical tweets).

In the same fashion, the methods proposed in [54] differ from each other in the features employed to describe the satirical texts: *SVM+W-B* considers features based on word n-grams whereas *SVM+Intrinsic* employs linguistic features which are topic-independent, and finally *SVM+ALL* combines both subgroups of features. In this case, the methods were evaluated using  $F1_{Macro}$ . To establish a fair comparison with the previous works, we evaluated the performance of our models that achieved the best results on each corpus independently (*salasEsMvAttLSTM\_{multi}^{Early}*, *salasMxMvAttLSTM\_{multi}^{Early}* and *barbEsMxMvAttLSTM\_{multi}^{Early}*) and we reevaluated the models on the two remaining corpora. As can be observed in Table 9 the three settings of our model  $MvAttLSTM_{multi}^{Early}$  outperformed by almost 10% points the results achieved by the best methods reported in [54,55].

#### 4.5. Discriminating between irony and satire

Some figurative languages devices like irony and satire are difficult to distinguish from each other due to they share several characteristics, even can be nested. For instance, satire can appeal to irony for communicating indirect and complex meanings often aiming at censoring or criticizing peoples, things, social and moral norms in an ironical way. Furthermore, to evaluate our model beyond *irony vs. non-irony* and *satire vs. non-satire* scenarios we evaluated the capability of our model for discriminating between both phenomena.

**Table 7**  
MvAttLSTM for satire detection in Spanish variants *Salas'17* and *Barbieri'15*.

Views	$F1_{sat}$	$F1_{no-sat}$	$F1_{Micro}$	$F1_{Macro}$	$F1_{sat}$	$F1_{no-sat}$	$F1_{Micro}$	$F1_{Macro}$
	$MvAttLSTM_{self}^{Early}$ <i>Salas'17-es</i>				$MvAttLSTM_{multihead}^{Early}$ <i>Salas'17-es</i>			
Bert+Muse+Aff_All	0.958	0.958	0.958	0.958	0.956	0.956	0.956	0.956
Bert+Muse+Aff_Emo	0.957	0.958	0.958	0.958	0.959	0.96	0.959	0.959
Bert+Muse+Aff_App	<b>0.96</b>	<b>0.961</b>	<b>0.961</b>	<b>0.961</b>	0.958	0.959	0.958	0.958
Bert+Muse+Contrast	0.956	0.958	0.957	0.957	0.959	0.96	0.959	0.959
Bert+Muse+LIWC	0.941	0.944	0.943	0.943	<b>0.959</b>	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>
Bert+Muse+SVerb	0.959	0.959	0.959	0.959	0.959	0.959	0.959	0.959
Bert+Muse+W_64	<b>0.959</b>	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>	<b>0.964</b>	<b>0.965</b>	<b>0.964</b>	<b>0.964</b>
Bert+Muse+W_128	0.955	0.956	0.955	0.955	0.953	0.954	0.954	0.954
Bert+Muse+W_All	<b>0.959</b>	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>
Bert+Muse+LingAll	0.957	0.958	0.958	0.958	0.953	0.954	0.954	0.954
Bert-baseline	0.938	0.925	0.924	0.924	0.938	0.925	0.924	0.924
	$MvAttLSTM_{self}^{Early}$ <i>Salas'17-mx</i>				$MvAttLSTM_{multihead}^{Early}$ <i>Salas'17-mx</i>			
Bert+Muse+Aff_All	0.964	0.964	0.964	0.964	0.956	0.955	0.956	0.956
Bert+Muse+Aff_Emo	0.948	0.947	0.948	0.948	0.956	0.955	0.956	0.955
Bert+Muse+Aff_App	0.947	0.946	0.947	0.947	0.952	0.951	0.952	0.952
Bert+Muse+Contrast	<b>0.97</b>	<b>0.969</b>	<b>0.969</b>	<b>0.969</b>	0.956	0.954	0.955	0.955
Bert+Muse+LIWC	0.965	0.964	0.965	0.965	<b>0.969</b>	<b>0.969</b>	<b>0.969</b>	<b>0.969</b>
Bert+Muse+SVerb	<b>0.967</b>	<b>0.967</b>	<b>0.967</b>	<b>0.967</b>	<b>0.967</b>	<b>0.966</b>	<b>0.966</b>	<b>0.966</b>
Bert+Muse+W_64	0.961	0.96	0.96	0.96	0.96	0.959	0.959	0.959
Bert+Muse+W_128	<b>0.967</b>	<b>0.966</b>	<b>0.967</b>	<b>0.967</b>	0.957	0.957	0.957	0.957
Bert+Muse+W_All	0.966	0.966	0.966	0.966	0.964	0.963	0.963	0.963
Bert+Muse+LingAll	0.961	0.960	0.961	0.961	0.966	0.965	0.966	0.965
Bert-baseline	0.941	0.950	0.951	0.951	0.941	0.950	0.951	0.951
	$MvAttLSTM_{self}^{Early}$ <i>Barbieri'15-es</i>				$MvAttLSTM_{multihead}^{Early}$ <i>Barbieri'15-es</i>			
Bert+Muse+Aff_All	0.953	0.952	0.953	0.953	<b>0.955</b>	<b>0.954</b>	<b>0.955</b>	<b>0.955</b>
Bert+Muse+Aff_Emo	0.954	0.953	0.954	0.954	0.953	0.952	0.952	0.952
Bert+Muse+Aff_App	<b>0.956</b>	<b>0.955</b>	<b>0.955</b>	<b>0.955</b>	0.952	0.95	0.951	0.951
Bert+Muse+Contrast	0.951	0.95	0.951	0.951	0.953	0.952	0.952	0.952
Bert+Muse+LIWC	0.954	0.953	0.954	0.954	0.954	0.953	0.954	0.954
Bert+Muse+SVerb	0.948	0.948	0.948	0.948	0.954	0.954	0.954	0.954
Bert+Muse+W_64	0.954	0.954	0.954	0.954	<b>0.958</b>	<b>0.957</b>	<b>0.957</b>	<b>0.957</b>
Bert+Muse+W_128	<b>0.956</b>	<b>0.956</b>	<b>0.956</b>	<b>0.956</b>	0.954	0.952	0.953	0.953
Bert+Muse+W_All	0.95	0.949	0.949	0.949	0.948	0.946	0.947	0.947
Bert+Muse+LingAll	0.953	0.952	0.952	0.952	0.954	0.954	0.954	0.954
Bert-baseline	0.942	0.947	0.947	0.947	0.942	0.947	0.947	0.947

**Table 8**

The impact of the views on MvAttLSTM for satire detection. The ignored view is denoted by (×) symbol whereas the included views are denoted by (✓) symbol.

Model	Ling	Muse	Bert	$F1_{sat}$	$F1_{no-sat}$	$F1_{Micro}$	$F1_{Macro}$
<i>Salas'17-es</i>							
Early-self	×	✓	✓	0.958	0.959	0.959	0.959
	✓	×	✓	0.955	0.956	0.955	0.955
	✓	✓	×	0.710	0.720	<b>0.721</b>	<b>0.715</b>
Early-multi	×	✓	✓	0.952	0.954	0.953	0.953
	✓	×	✓	0.952	0.954	0.953	0.953
	✓	✓	×	0.603	0.685	<b>0.654</b>	<b>0.644</b>
<i>Salas'17-mx</i>							
Early-self	×	✓	✓	0.958	0.959	0.959	0.959
	✓	×	✓	0.958	0.957	0.958	0.958
	✓	✓	×	0.719	0.858	<b>0.822</b>	<b>0.788</b>
Early-multi	×	✓	✓	0.956	0.956	0.956	0.956
	✓	×	✓	0.936	0.931	0.934	0.934
	✓	✓	×	0.649	0.683	<b>0.681</b>	<b>0.667</b>
<i>Barbieri'15-es</i>							
Early-self	×	✓	✓	0.95	0.949	0.95	0.95
	✓	×	✓	0.952	0.950	0.951	0.951
	✓	✓	×	0.743	0.728	<b>0.736</b>	<b>0.735</b>
Early-multi	×	✓	✓	0.939	0.935	0.937	0.937
	✓	×	✓	0.951	0.950	0.951	0.951
	✓	✓	×	0.743	0.702	<b>0.725</b>	<b>0.722</b>

In this sense, the first step was building the satire-irony corpus. For that purpose, we merged the 1000 ironic tweets of the *IroSvA'19* corpus with the 2500 satirical tweets of *Salas'17* for the Mexican and Castilian variants of the Spanish independently. After that, we reevaluated the best model for irony detection ( $iroMvAttLSTM_{multi}^{Contextual}$ ) and the best model for satire detection

( $satMvAttLSTM_{multi}^{Early}$ ) in each variant (see Section 4.4). In Table 10, we present the results obtained. As can be observed, the results show that our model is able to effectively discriminate satire from irony in both variants (*es* and *mx*) with an effectiveness  $F1_{Macro} = 0.987$  for Castilian tweets and  $F1_{Macro} = 0.978$  for



**Table 9**

Comparison with state-of-the-art methods for satire detection in Spanish variants.

Method	Salas'17-mx			Salas'17-es			Barbieri'15-es
	$P_{sat}$	$R_{sat}$	$F1_{sat}$	$P_{sat}$	$R_{sat}$	$F1_{sat}$	$F1_{Macro}$
SMO+LIWC-ALL	0.855	0.855	0.855	0.846	0.84	0.84	–
BayesNet+LIWC-ALL	0.757	0.756	0.756	0.734	0.734	0.734	–
J48+LIWC-ALL	0.752	0.752	0.752	0.774	0.774	0.774	–
SVM+W-B	–	–	–	–	–	–	0.738
SVM+Intrinsic	–	–	–	–	–	–	0.816
SVM+ALL	–	–	–	–	–	–	0.852
salasEsMvAttLSTM <sup>Early multi</sup>	0.966	0.973	<b>0.969</b>	0.969	0.950	<b>0.959</b>	<b>0.954</b>
salasMxMvAttLSTM <sup>Early multi</sup>	0.951	0.969	<b>0.960</b>	0.973	0.955	<b>0.964</b>	<b>0.957</b>
barbEsMvAttLSTM <sup>Early multi</sup>	0.951	0.969	<b>0.960</b>	0.973	0.955	<b>0.964</b>	<b>0.957</b>

**Table 10**

MvAttLSTM for irony vs. satire detection in Castilian and Mexican tweets.

Views	$F1_{iro}$	$F1_{sat}$	$F1_{Micro}$	$F1_{Macro}$	$F1_{iro}$	$F1_{sat}$	$F1_{Micro}$	$F1_{Macro}$
	sat – MvAttLSTM <sup>Early multihead</sup> es				sat – MvAttLSTM <sup>Early multihead</sup> mx			
Bert+Muse+Aff_All	0.950	0.980	0.972	0.965	0.953	0.981	0.973	0.967
Bert+Muse+Aff_Emo	<b>0.979</b>	<b>0.992</b>	<b>0.988</b>	<b>0.985</b>	<b>0.968</b>	<b>0.987</b>	<b>0.982</b>	<b>0.978</b>
Bert+Muse+Aff_App	<b>0.981</b>	<b>0.992</b>	<b>0.989</b>	<b>0.987</b>	0.955	0.982	0.975	0.969
Bert+Muse+Contrast	0.957	0.983	0.976	0.970	<b>0.964</b>	<b>0.986</b>	<b>0.979</b>	<b>0.975</b>
Bert+Muse+LIWC	0.976	0.990	0.986	0.983	0.952	0.981	0.972	0.966
Bert+Muse+SVerb	0.886	0.957	0.938	0.921	0.948	0.980	0.971	0.964
Bert+Muse+W_64	0.953	0.982	0.974	0.968	0.948	0.979	0.970	0.964
Bert+Muse+W_128	0.953	0.982	0.974	0.967	0.950	0.980	0.971	0.965
Bert+Muse+W_All	0.956	0.983	0.975	0.970	0.888	0.963	0.945	0.926
Bert+Muse+LingAll	0.920	0.966	0.952	0.943	0.952	0.981	0.973	0.966
Bert-baseline	0.983	0.991	0.987	0.984	0.963	0.975	0.964	0.956
	iro – MvAttLSTM <sup>Contextual multihead</sup> es				iro – MvAttLSTM <sup>Contextual multihead</sup> mx			
Bert+Muse+Aff_All	<b>0.966</b>	<b>0.987</b>	<b>0.981</b>	<b>0.976</b>	0.962	0.985	0.979	0.974
Bert+Muse+Aff_Emo	0.944	0.975	0.966	0.959	<b>0.968</b>	<b>0.987</b>	<b>0.982</b>	<b>0.978</b>
Bert+Muse+Aff_App	<b>0.966</b>	<b>0.986</b>	<b>0.981</b>	<b>0.976</b>	0.963	0.985	0.979	0.974
Bert+Muse+Contrast	0.963	0.986	0.980	0.975	0.963	0.985	0.979	0.974
Bert+Muse+LIWC	0.965	0.986	0.980	0.975	0.955	0.982	0.974	0.968
Bert+Muse+SVerb	0.963	0.986	0.979	0.974	0.964	0.985	0.979	0.975
Bert+Muse+W_64	0.947	0.980	0.971	0.963	0.963	0.985	0.979	0.974
Bert+Muse+W_128	0.964	0.986	0.980	0.975	0.964	0.986	0.979	0.975
Bert+Muse+W_All	0.929	0.965	0.953	0.947	0.964	0.986	0.980	0.975
Bert+Muse+LingAll	<b>0.966</b>	<b>0.987</b>	<b>0.981</b>	<b>0.976</b>	0.962	0.985	0.979	0.974
Bert-baseline	0.983	0.991	0.987	0.984	0.963	0.975	0.964	0.956

Mexican tweets. Also, *Bert-baseline* showed very high results on both corpora.

Concretely, the model *satMvAttLSTM<sup>Early multi</sup>* achieves, in general, the best performance in both variants. All these results make evident that those views such as *Linguistic-view* and *Muse-views* have a low impact on the model. A possible reason is that these views were learned without any supervision related to the specific task. Conversely, *Bert-view*, which is a task-dependent view, has a major impact on the model effectiveness, particularly due to this view was learned in a supervised way and it is strongly related to the specific task dataset. Regarding the linguistics views, the best results of *satMvAttLSTM<sup>Early multi</sup>* in the Mexican and Castilian variants were *Aff\_Emo* and *Aff\_App* respectively. According to these results, we could appreciate that affective information was, in general, the most relevant to inform our model for capturing useful information to detect irony and satire in Spanish variants.

Table 11 shows the impact of the views to inform our model. The obtained results are aligned with the results achieved in the task of satire detection. As can be observed, ignoring *Bert-view* caused the most significant drop in the performance of the *MvAttLSTM* model whereas the model is less sensitive to exclude *Ling-view* and *Muse-view*.

Analyzing the results presented in Tables 7 and 10, we could appreciate that our model is better discriminating irony from satire than irony from no-irony and satire from no-satire. This behavior is caused by the nature of the dataset. Satirical tweets were retrieved from different topics than ironic tweets. This fact

introduces a bias with respect to the topics discussed in the ironic and satirical tweets.

#### 4.6. Validating the robustness of the models in humor recognition

Our final experiment aims at investigating the robustness of our model for recognizing humorous tweets written in Spanish. We are intrigued by the fact that irony and satire are two phenomena that are strongly related to humor. Particularly, some theoretical works comments about the relation between humor-irony [167,168] and humor-satire [19].

In this sense, we evaluated our model with the corpus *HAHA'19* and the results are shown in Table 12. In this case, only the results achieved by our model using the *Early fusion* strategy are presented due to the *Contextual fusion* method obtained worse results. At a first glance, we can appreciate that our model achieves very similar results for both attention mechanisms, although the model *MvAttLSTM<sup>Early multi</sup>* shows a slight improvement in terms of  $F1_{humor}$ . Another important aspect to notice is regarding the linguistic views, particularly the model *MvAttLSTM<sup>Early self</sup>* that performs better when affective features such as *Aff\_Emo* and *Aff\_App* are used. However, the model *MvAttLSTM<sup>Early multi</sup>* obtains its best results when more features are considered, particularly those best-ranked according to the Wilcoxon test  $W_{128}$  and  $W_{All}$ . This behavior is aligned with the results presented in [145]. Also, in this corpus *Bert-baseline* achieved competitive results in comparison with our proposed models. With respect to the impact of

**Table 11**

The impact of the views on MvAttLSTM for irony and satire distinguishing. The ignored view is denoted by (×) symbol whereas the included views are denoted by (✓) symbol.

Model	Ling	Muse	Bert	$F1_{iro}$	$F1_{sat}$	$F1_{Micro}$	$F1_{Macro}$
<i>Castilian variant</i>							
Early-multi (sat)	×	✓	✓	0.957	0.983	0.976	0.970
	✓	×	✓	0.966	0.987	0.981	0.977
	✓	✓	×	0.942	0.977	<b>0.967</b>	<b>0.959</b>
Contextual-multi(iro)	×	✓	✓	0.961	0.985	0.978	0.973
	✓	×	✓	0.967	0.987	0.981	0.976
	✓	✓	×	0.873	0.930	<b>0.911</b>	<b>0.901</b>
<i>Mexican variant</i>							
Early-multi(sat)	×	✓	✓	0.951	0.981	0.973	0.966
	✓	×	✓	0.948	0.979	0.970	0.964
	✓	✓	×	0.407	0.896	<b>0.824</b>	<b>0.651</b>
Contextual-multi(iro)	×	✓	✓	0.964	0.986	0.979	0.975
	✓	×	✓	0.966	0.986	0.981	0.976
	✓	✓	×	0.532	0.730	<b>0.740</b>	<b>0.631</b>

**Table 12**

MvAttLSTM for humor recognition in Spanish tweets (HAHA'19)

Views	$F1_{hum}$	$F1_{no-hum}$	$F1_{Micro}$	$F1_{Macro}$	$F1_{hum}$	$F1_{no-hum}$	$F1_{Micro}$	$F1_{Macro}$
	$MvAttLSTM_{self}^{Early}$				$MvAttLSTM_{multithread}^{Early}$			
Bert+Muse+Aff_All	0.802	0.876	0.848	0.839	0.794	0.879	0.848	0.837
Bert+Muse+Aff_Emo	<b>0.804</b>	<b>0.881</b>	<b>0.852</b>	<b>0.842</b>	0.799	0.877	0.847	0.838
Bert+Muse+Aff_App	<b>0.804</b>	<b>0.878</b>	<b>0.85</b>	<b>0.841</b>	0.798	0.879	0.849	0.838
Bert+Muse+LIWC	0.796	0.881	0.85	0.839	0.798	0.88	0.849	0.839
Bert+Muse+SVerb	0.798	0.88	0.849	0.839	0.797	0.877	0.847	0.837
Bert+Muse+W_64	0.803	0.881	0.852	0.842	0.801	0.879	0.85	0.84
Bert+Muse+W_128	0.798	0.881	0.85	0.839	<b>0.806</b>	<b>0.879</b>	<b>0.851</b>	<b>0.842</b>
Bert+Muse+W_All	0.795	0.88	0.849	0.838	<b>0.804</b>	<b>0.882</b>	<b>0.853</b>	<b>0.843</b>
Bert+Muse+LingAll	0.802	0.872	0.845	0.837	0.796	0.88	0.849	0.838
Bert-baseline	0.802	0.864	0.84	0.833	0.802	0.864	0.84	0.833

**Table 13**

The impact of the views on MvAttLSTM for humor recognition. The ignored view is denoted by (×) symbol whereas the included views are denoted by (✓) symbol.

Model	Ling	Muse	Bert	$F1_{hum}$	$F1_{no-hum}$	$F1_{Micro}$	$F1_{Macro}$
<i>HAHA'19</i>							
Early-self	×	✓	✓	0.792	0.88	0.848	0.836
	✓	×	✓	0.79	0.874	0.843	0.832
	✓	✓	×	0.783	0.876	<b>0.842</b>	<b>0.829</b>
Early-multi	×	✓	✓	0.795	0.88	0.848	0.838
	✓	×	✓	0.784	0.878	0.844	0.831
	✓	✓	×	0.781	0.881	<b>0.845</b>	<b>0.831</b>

the views in our models, we observed in Table 13 that *Bert-view* and *Muse-view* are more important than *Linguistic-view*.

We compare the results of *MvAttLSTM<sup>Early</sup><sub>multi</sub>* with those of the participating systems in the shared task at HAHA'19 organized in the framework of IberLEF'19. In this task, the systems were ranked according to the official measure F1 score in the humor class, although also and Acc was reported. As can be observed in Table 14, the results obtained by our model are very competitive, obtaining the fourth position of the ranking according to  $F1_{humor}$  and the third position in terms of Acc out of 18 systems. The performance of our model is similar to the best-ranked system *Adilism* in terms of F1. However, the difference in terms of precision and recall shows that the *Adilism* system is better at detecting a major number of humorous tweets whereas our model is better at detecting the humorous tweets. As future work, a deeper study is required to analyze the low recall achieved by our model compared to the *Adilism* system.

## 5. Conclusions and future work

In this work, we have presented MvAttLSTM, a deep learning-based method for irony and satire detection in Spanish variants.

It is based on an Attentive-LSTM model informed with additional knowledge learned from three distinct perspectives: *Linguistic-view*, *MUSE-based view*, and *BERT-based view*. We observed that our model achieved better performance when it is enriched with the three proposed views. We have evaluated our model on the corpus *IroSvA'19* for irony and on the corpora *Salas'17* and *Barbieri'15* for satire detection in Spanish variants. In both tasks, the model outperforms the state-of-the-art results. Furthermore, we have evaluated our model on humor recognition using the corpus HAHA'19 showing a very competitive behavior. Particularly, linguistic information and deep sentence encoding were more feasible for irony detection whereas BERT views increased the performance of satire detection and satire vs. irony detection (RQ1). Interestingly, the results revealed that affective information helps in detecting irony and satire. Particularly, those related to emotions (*Aff\_Emo*) which are based on the resources SenticNet and SEL; and those related to attitude words (*Aff\_App*) based on the LAM lexicon. Experiments also confirmed that both fusion strategies are feasible. However *Contextual fusion* achieved better performance in relative small corpus like *IroSvA'19*, whereas the *Early fusion* takes advantage of large and self-annotated corpora

**Table 14**

Comparison with state of the art systems for humor recognition in Spanish (HAHA'2019).

Ranking	Team	$P_{hum}$	$R_{hum}$	$F1_{hum}$	Acc
1st	Adilism	0.791	0.852	0.821	0.855
2nd	Kevin & Hiromi	0.802	0.831	0.816	0.854
3th	Bfarzin	0.782	0.839	0.810	0.846
**	$MvAttLSTM_{multi}^{Early}$	<b>0.819</b>	0.792	0.806	<b>0.851</b>
4th	Jamestjw	0.793	0.804	0.798	0.842
5th	INGEOTEC	0.758	0.819	0.788	0.828
6th	BLAIR GMU	0.745	0.827	0.784	0.822
7th	UO_UPV2	0.78	0.765	0.773	0.824
...	...	...	...	...	...
18th	Amrita CEN	0.478	0.514	0.495	0.591

**Table A.15**

Description of the linguistic features used in the Linguistic view to inform the proposed models.

Group	Feature	Description
Stylistics-based	multiLines	It takes into account whether the tweet is composed of multiple lines or not (one vs. many lines).
	lengthW lengthC meanLengthW	Three different features are considered; (i) the number of words, (ii) the number of characters, and (iii) the means of words' length in the tweet.
	isDialog nDialogMark	Two distinct features are considered; (i) the tweet contains any line that starts with a long dash (dialog marker), (ii) the number of lines that start with long dashes.
	hashtagsFreq urlsFreq emojisFreq	These count the number of hashtags, URLs, and emojis in the tweet, respectively.
	exclMarkFreq	It counts the exclamation marks in the tweet.
	wordRep wordUpper wordCharRep wordWithExcl	Four distinct features are considered: (i) the number of words emphasized by word's repetition, (ii) the number of words with emphasis by uppercase, (iii) the number of words emphasized by character flooding, and (iv) the number of words emphasized by continues exclamation marks.
	alliter	It captures the occurrence of simple alliteration in the tweet. For that, we considered a fixed-length sequence of phonetic prefixes with size=3.
	quotation	It quantifies the phrases enclosed in a double quote.
	Q?A	It quantifies the question and answer structures in the tweet.
	person_p <sup>a</sup>	It quantifies the number of verbs conjugated in the first, second, third persons and the nouns and adjectives which agree with such verbal conjugations.
	tense_t <sup>b</sup>	It quantifies the usage of different verbal tenses in the tweet.
	posN posV posA posR	These count the nouns, verbs, adverbs and adjectives in the tweet.
	Punctuation	It counts the occurrence of dots, commas, semicolons, and question marks in the tweet.
Content-based	Animal centered-words	It counts the words that occur in a lexicon of animal names.
	Toponym words	It counts the words that occur in a lexicon of country's names, capital's names, city's names and nationalities.
	ObsceneSexual words	It counts the words that occur in an in-house lexicon of sexual and obscene words.

(continued on next page)

(RQ2). Unexpectedly we found no strong differences in the effectiveness of our model when self-attention or multi-head attention was considered. However, we appreciated that each attention type attends distinct linguistic features (RQ3). We are aware that our model has an important limitation which lies in the lack of explainability about why the model used some features from one setting to another and from one Spanish variant to the others. As future work, we will aim at investigating other methods for fusing the additional knowledge into our model. Moreover, we plan to carry out a fine-grained analysis on the impact of linguistic features joined with the information captured by the attention mechanism for irony and satire interpretability. Finally, we are interested in exploring our model in multilingual and cross-lingual settings.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgments

The work of the first two authors was in the framework of the research project MISIMIS-FAKENHATE on MISinformation and MIScommunication in social media: FAKE news and HATE speech (PGC2018-096212-B-C31), funded by Spanish Ministry of Science and Innovation, and DeepPattern (PROMETEO/2019/121), funded by the Generalitat Valenciana, Spain.



**Table A.15** (continued).

Group	Feature	Description
Semantic-based	Antonyms	It quantifies the pairs of antonyms that occurs in the tweet. This feature is based on the antonym' relations provided by WordNet [169], particularly, for the Spanish language we used the Multilingual Central Repository (MCR) [170].
	LexAmbiguity	Three different features are considered; (i) the average of the meanings associated with each word in the tweet, (ii) the number of meanings for the most ambiguous word in the tweet, (iii) the gap between the value of two previous features.
	DomAmbiguity	Conversely, to consider the meanings of the words, in these features we consider the number of domains assigned to the words. Particularly, three distinct features are considered; (i) the average of domains associated with each word in the tweet, (ii) the greatest number of domains that a single word has in the tweet, (iii) the gap between the value of the two previous features. For obtaining the domains of the words we used the WordNet Domains <sup>c</sup> and SUMO <sup>d</sup> each separately.
	SVerb_classes	These features capture distinct semantic frames of the verbs in the tweet based on ADDESE. <sup>e</sup>
	Negation	It counts the negation words in the tweet.
Affective-based	SSL_polarity ESL_polarity CriSol_polarity LAM11_polarity <sup>f</sup> SenticNet_polarity	These features count positive and negative words in many sentiment resources. Notice that, for each resource two features are computed. Particularly, we explore four distinct dictionaries: Spanish Sentiment Lexicon (SSL) [171], Elhuyar Sentiment Lexicon (ESL) [172], CriSol lexicon [171], and the lexicon LAM11 introduced in [166]. Moreover, the polarity score associated with the words and concepts in SenticNet was considered.
	emojiPol_pos emojiPol_neg	The number of positive and negative emoticons and emojis considering the resource Emoticons Sentiment [173].
	LAM11_attitude eCrisol_attitude <sup>g</sup>	These features count the number of words according to the three distinct attitude categories (affect, judgment, and appreciation) proposed in [166]. For that, we considered two lexicons, (i) the LAM11 lexicon introduced in [166] and an extended version of the CriSol lexicon, where all words were automatically annotated with attitudes (eCrisol) by using the method proposed in [174].
	EmoCat	These features count the number of words according to the six basic emotions provided by the resource SEL [165].
	EmoDim	These features are based on the four affective dimensions in SenticNet of the Cambria's hourglass of emotions model [148,175]: introspection, temper, attitude and sensitivity. <sup>h</sup>
Contrast-based <sup>i</sup>	wordPolCont	It computes the gap between the most positive and the most negative word in the tweet. This feature, consider the distance, in terms of tokens, between the words.
	emoTextPol-Cont	It computes the polarity difference between emoticons and words in the tweet.
	antConsPolCont	It considers the polarity contrast between two parts of the tweet when the tweet is split by a delimiter. In this work we consider as delimiter some adverbs and punctuation marks.
	meanPolPhrase	It is the mean of the polarities of the words that belong to phrases enclosed by quotes.
	polStandDev	It is the standard deviation of the polarities of the words that belong to phrases enclosed by quotes.
	prePastPolCont	It computes the polarity difference between the parts of the tweet written in present and past tenses.
	skipGPolRate	It computes the rate among skip-grams with polarity opposition on the total of candidate skip-grams. The candidate skip-grams are those composed of two words (nouns, adjectives, verbs, adverbs) with skip=1. The skip-grams with polarity opposition are those that match with the patterns positive-negative, positive-neutral, negative-neutral, and vise-versa.
	upperTextPol-Cont	It computes the polarity difference between capitalized words and the remainder words in the tweets.
Psycholinguistic-based	LIWC_cat <sup>j</sup>	These features count the frequency of words in each category provided by the resource Linguistic Inquiry and Word Count <sup>k</sup> dictionary [176].

<sup>a</sup>*p* is parametric to the three persons used in Spanish grammar.<sup>b</sup>*t* is parametric to the various tense in Spanish grammar i.e., present, past, future, etc.<sup>c</sup><http://wdomains.fbk.eu/hierarchy.html>.<sup>d</sup><http://www.adampease.org/OP/>.<sup>e</sup><http://adesse.uvigo.es/data/clases.php>.<sup>f</sup>*polarity* is parametric to the type of sentiment, positive and negative.<sup>g</sup>*attitude* is parametric to the type of attitudes affect, judgment, and appreciation.<sup>h</sup>It is worthy to note that for the Spanish language, we used BabelSenticNet [147]. In this extension of SenticNet, the affective dimensions are sensitivity, attention, aptitude and pleasantness.<sup>i</sup>With the aim of capturing some types of explicit polarity opposition, we used the features proposed in [177]. The Spanish version of SenticNet was used to determine the polarity contrast between different parts of the text.<sup>j</sup>*cat* is parametric to the 68 categories in the LIWC 2001 Spanish dictionary.<sup>k</sup><http://www.liwc.net>.**Appendix A. Linguistic features****Appendix B. Best MvAttLSTM hyperparameters for each corpus**

See Table A.15.

See Tables B.16 and B.17.

**Table B.16**

Hyperparameters for the Contextual MvAttLSTM model on the IroSvA'19 corpora.

Dataset	Model	Hyperparameters	Model	Hyperparameters
IroSvA'19-es	Contextual Self	batch=256 att=self h=1 dp=0.25 op=adam lr=0.01	Contextual Multi	batch=256 att=multihead h=2 dp=0.3 op=adam lr=0.01
IroSvA'19-mx	Contextual Self	batch=32 att=self h=1 dp=0.3 op=adam lr=0.001	Contextual Multi	batch=128 att=multihead h=8 dp=0.25 op=adam lr=0.01
IroSvA'19-cu	Contextual Self	batch=128 att=self h=1 dp=0.4 op=rmsprop lr=0.01	Contextual Multi	batch=256 att=multihead h=8 dp=0.3 op=rmsprop lr=0.01

**Table B.17**

Hyperparameters for the Early MvAttLSTM model on the satire and humor corpora.

Dataset	Model	Hyperparameters	Model	Hyperparameters
Salas'17-es	Early Self	batch=32 att=self h=1 dp=0.25 op=adam lr=0.01	Early Multi	batch=128 att=multihead h=4 dp=0.25 op=adam lr=0.01
Salas'17-mx	Early Self	batch=32 att=self h=1 dp=0.4 op=adam lr=0.01	Early Multi	batch=128 att=multihead h=2 dp=0.4 op=rmsprop lr=0.01
Barbieri'15-es	Early Self	batch=32 att=self h=1 dp=0.25 op=rmsprop lr=0.01	Early Multi	batch=128 att=multihead h=4 dp=0.3 op=adam lr=0.01
HAHA'19	Early Self	batch=32 att=self h=1 dp=0.4 op=adam lr=0.01	Early Multi	batch=128 att=multihead h=4 dp=0.25 op=rmsprop lr=0.01

## References

- [1] H.P. Grice, Logic and conversation, in: P. Cole, J. Morgan (Eds.), *Syntax and Semantics 3: Speech Acts*, Academic Press, New York, 1975, pp. 41–58.
- [2] H.P. Grice, Further notes on logic and conversation, in: P. Cole (Ed.), *Syntax and Semantics 9: Pragmatics*, Academic Press, New York, 1978, pp. 113–127.
- [3] J. Raymond W. Gibbs, H.L. Colston, *Interpreting Figurative Meaning*, Cambridge University Press, 2012, <http://dx.doi.org/10.1080/10926488.2018.1407996>.
- [4] B. Dancygier, *Figurative Language*, Cambridge University Press, 2014.
- [5] H.L. Colston, *Using Figurative Language*, Cambridge University Press, 2015.
- [6] A. Reyes, *Linguistic-Based Patterns for Figurative Language Processing: The Case of Humor Recognition and Irony Detection* Antonio Reyes P (Ph.D.), Universitat Politècnica de València, 2012.
- [7] J. Lucariello, Situational irony: A concept of events gone awry, *J. Exp. Psychol. [Gen.]* 123 (1994) 129–145, <http://dx.doi.org/10.1037/0096-3445.123.2.129>.
- [8] S. Attardo, Irony as relevant inappropriateness, *J. Pragmat.* 32 (6) (2000) 793–826, [http://dx.doi.org/10.1016/S0378-2166\(99\)00070-3](http://dx.doi.org/10.1016/S0378-2166(99)00070-3).
- [9] R.J. Kreuz, S. Glucksberg, How to be sarcastic: The echoic reminder theory of verbal irony, *J. Exp. Psychol. [Gen.]* 118 (1989) 374–386, <http://dx.doi.org/10.1037/0096-3445.118.4.374>.
- [10] R.J. Kreuz, R.M. Roberts, On satire and parody: The importance of being ironic, *Metaphor Symbol. Act.* 8 (1993) 97–109, [http://dx.doi.org/10.1207/s15327868ms0802\\_2](http://dx.doi.org/10.1207/s15327868ms0802_2).
- [11] R.J. Kreuz, K.E. Link, Asymmetries in the use of verbal irony, *J. Lang. Soc. Psychol.* 21 (2002) 127–143, <http://dx.doi.org/10.1177/02627X02021002002>.
- [12] D. Sperber, D. Wilson, Irony and the use-mention distinction, in: P. Cole (Ed.), *Radical Pragmatics*, Academic Press, New York, 1981, pp. 295–318.
- [13] R.W. Gibbs, J.E. O'Brien, S. Doolittle, Inferring meanings that are not intended: Speakers' intentions and irony comprehension, *Discourse Processes* 20 (1995) 187–203, <http://dx.doi.org/10.1080/01638539509544937>.
- [14] J. Haiman, *Talk Is Cheap: Sarcasm, Alienation, and Evolution of Language*, Oxford University Press, New York, USA, 1998, <http://dx.doi.org/10.1017/s0047404500211032>.
- [15] L. Colletta, Political satire and postmodern irony in the age of Stephen Colbert and Jon Stewart, *J. Popul. Cult.* 42 (2009) 856–874, <http://dx.doi.org/10.1111/j.1540-5931.2009.00711.x>.
- [16] C. Condren, in: S. Attardo *Satire* (Ed.), *Encyclopedia of Humor Studies Chapter Satire*, SAGE Publications, Inc, USA, 2014, pp. 661–664.
- [17] D. Wilson, D. Sperber, On verbal irony, *Lingua* 87 (1992) 53–76, [http://dx.doi.org/10.1016/0024-3841\(92\)90025-E](http://dx.doi.org/10.1016/0024-3841(92)90025-E).
- [18] P. Brown, S.C. Levinson, *Politeness: Some Universals in Language Usage*, second ed., Cambridge University Press, 1987, <http://dx.doi.org/10.2307/3587263>.
- [19] P. Simpson, *On the Discourse of Satire: Towards a Stylistic Model of Satirical Humour*, Vol. 2, John Benjamins Publishing Company, 2003, <http://dx.doi.org/10.1177/0963947006060558>.

- [20] T. Veale, Y. Hao, Support Structures for Linguistic Creativity : A Computational Analysis of Creative Irony in Similes, in: *Proceedings of CogSci 2009*, the 31st Annual Meeting of the Cognitive Science Society, 2009, pp. 1376–1381.
- [21] D. Maynard, M.A. Greenwood, Who cares about sarcastic tweets? Investigating the impact of sarcasm on sentiment analysis, in: N. Calzolari, K. Choukri, T. Declerck, H. Loftsson, B. Maegaard, J. Mariani, A. Moreno, J. Odijk, S. Piperidis (Eds.), *Proceedings of the 9th International Conference on Language Resources and Evaluation, LREC 2014*, 423, European Language Resources Association (ELRA), Reykjavik, Iceland, 2014, pp. 8–4243.
- [22] A. Ghosh, G. Li, T. Veale, P. Rosso, E. Shutova, J. Barnden, A. Reyes, Semeval-2015 task 11: Sentiment analysis of figurative language in Twitter, in: *Proceedings of the 9th International Workshop on Semantic Evaluation, SemEval'15*, Association for Computational Linguistics, Denver, Colorado, 2015, pp. 470–478, <http://dx.doi.org/10.18653/v1/s15-2080>.
- [23] V. Basile, A. Bolioli, M. Nissim, V. Patti, P. Rosso, Overview of the Evalita 2014 SENTIMENT POLARITY classification task, in: *Proceedings of the 1st Italian Conference on Computational Linguistics (CLiC-it 2014) & the Fourth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian EVALITA 2014*, 2014, pp. 50–57.
- [24] F. Barbieri, V. Basile, D. Croce, M. Nissim, N. Novielli, V. Patti, Overview of the Evalita 2016 sentiment polarity classification task, in: *Proceedings of 3rd Italian Conference on Computational Linguistics (CLiC-It 2016) & 5th Evaluation Campaign of Natural Language Processing and Speech Tools for Italian, EVALITA 2016*, Napoli, Italy, December 5-7, 2016, vol. 1749, CEUR Workshop Proceedings. CEUR-WS.org, 2016.
- [25] C.V. Hee, *Can Machines Sense Irony?* (Ph.D. thesis), Universiteit Gent, 2017.
- [26] D.I. Hernández Farías, P. Rosso, Irony, sarcasm, and sentiment analysis, in: *Sentiment Analysis in Social Networks*, 2017, pp. 113–128, <http://dx.doi.org/10.1016/B978-0-12-804412-4.00007-3>.
- [27] C. Zucco, H. Liang, G.D. Fatta, M. Cannataro, Explainable sentiment analysis with applications in medicine, in: *Proceeding of the IEEE International Conference on Bioinformatics and Biomedicine, BIBM 2018*, 2019, pp. 1740–1747, <http://dx.doi.org/10.1109/BIBM.2018.8621359>.
- [28] F. Bodria, A. Panisson, A. Perotti, S. Piaggese, Explainability Methods for Natural Language Processing: Applications to Sentiment Analysis, in: *CEUR Workshop Proceedings*, Vol. 2646, 2020, pp. 100–107.
- [29] R. Justo, J.M. Alcaide, M.I. Torres, M. Walker, Detection of sarcasm and nastiness: New resources for spanish language, *Cognitive Comput.* 10 (2018) 1135–1151, <http://dx.doi.org/10.1007/s12559-018-9578-5>.
- [30] V.L. Rubin, N.J. Conroy, Y. Chen, S. Cornwell, Fake news or truth ? Using satirical cues to detect potentially misleading news, in: *Proceedings of the Workshop on Computational Approaches To Deception Detection at the 15th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-CADD2016*, Association for Computational Linguistics, California, USA, 2016, pp. 7–17, <http://dx.doi.org/10.18653/v1/W16-0802>.
- [31] J. Golbeck, M. Mauriello, B. Auxier, K.H. Bhanushali, C. Bonk, M.A. Bouzaghrane, C. Buntain, R. Chanduka, P. Chekalos, J.B. Everett, W. Falak, C. Geringer, J. Graney, K.M. Hoffman, L. Huth, Z. Ma, M. Jha, M. Khan, V. Kori, E. Lewis, G. Mirano, W.T. Mohn, S. Mussenden, T.M. Nelson, S. Mcwillie, A. Pant, P. Shetye, R. Shrestha, A. Steinheimer, A. Subramanian, G. Visnansky, Fake News vs Satire: A Dataset and Analysis, in: *Proceedings of the 10th ACM Conference on Web Science, WebSci 2018*, Amsterdam, Netherlands, 2018, pp. 17–21.
- [32] B.C. Wallace, Computational irony: A survey and new perspectives, *Artif. Intell. Rev.* 43 (2015) 467–483.
- [33] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural Comput.* 9 (1997) 1735–1780, <http://dx.doi.org/10.1162/neco.1997.9.8.1735>.
- [34] J. Chung, C. Gulcehre, K. Cho, Y. Bengio, Gated recurrent neural networks on sequence modeling, in: *Proceedings of the NIPS'2014 Deep Learning Workshop*, 2014, pp. 1–9, arXiv: [arXiv:1412.3555v1](https://arxiv.org/abs/1412.3555).
- [35] T. Luong, H. Pham, C.D. Manning, Effective approaches to attention-based neural machine translation, in: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Lisbon, Portugal, 2015, pp. 1412–1421, <http://dx.doi.org/10.18653/v1/D15-1166>, URL <https://www.aclweb.org/anthology/D15-1166>.
- [36] Y. Wang, M. Huang, L. Zhao, et al., Attention-based LSTM for aspect-level sentiment classification, in: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2016, pp. 606–615.
- [37] Z. Yang, D. Yang, C. Dyer, X. He, A. Smola, E. Hovy, Hierarchical attention networks for document classification, in: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics, Human Language Technologies*, 2016, pp. 1480–1489.
- [38] M. Yang, W. Tu, J. Wang, F. Xu, X. Chen, Attention based LSTM for target dependent sentiment classification, in: *AAAI*, 2017, pp. 5013–5014.
- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, in: *Proceedings of the 31st Conference on Neural Information Processing Systems, NIPS 2017*, Long Beach, CA, USA, 2017, pp. 1–11.
- [40] J. Devlin, M.W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, in: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Volume 1 (Long and Short Papers), Association for Computational Linguistics (ACL), Minneapolis, Minnesota, USA, 2019, pp. 4171–4186, URL <https://www.aclweb.org/anthology/N19-1423>, arXiv: [arXiv:1910.10219](https://arxiv.org/abs/1910.10219).
- [41] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, F. Guzmán, G. Wenzek, M. Ott, L. Zettlemoyer, V. Stoyanov, Unsupervised cross-lingual representation learning at scale, in: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, 2020, pp. 8440–8451, URL <https://www.aclweb.org/anthology/2020.acl-main.747>, arXiv: [arXiv:1911.02116v2](https://arxiv.org/abs/1911.02116).
- [42] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, RoBERTa: A robustly optimized BERT pretraining approach, 2019, arXiv: [arXiv:1907.11692](https://arxiv.org/abs/1907.11692).
- [43] V. Sanh, L. Debut, J. Chaumond, T. Wolf, DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter, 2019, pp. 1–5, arXiv: [arXiv:1910.01108](https://arxiv.org/abs/1910.01108).
- [44] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, R. Soricut, ALBERT: A lite BERT for self-supervised learning of language representations, in: *Proceedings of the International Conference on Learning Representation, ICLR 2020*, 2020, pp. 1–17, arXiv: [arXiv:1909.11942v6](https://arxiv.org/abs/1909.11942).
- [45] J.Á. González, L.F. Hurtado, F. Pla, Transformer based contextualization of pre-trained word embeddings for irony detection in Twitter, *Inf. Process. Manage.* 57 (2020) 1–15, <http://dx.doi.org/10.1016/j.ipm.2020.102262>.
- [46] R.A. Potamias, G. Siolas, A.G. Stafylopatis, A transformer-based approach to irony and sarcasm detection, *Neural Comput. Appl.* 32 (2020) 17309–17320, <http://dx.doi.org/10.1007/s00521-020-05102-3>, arXiv: [arXiv:1911.10401](https://arxiv.org/abs/1911.10401).
- [47] D. Ghosh, A. Vajpayee, S. Muresan, A report on the 2020 sarcasm detection shared task, in: *Proceedings of the 2nd Workshop on Figurative Language Processing*, Association for Computational Linguistics, 2020, pp. 1–11.
- [48] S. Serrano, N.A. Smith, Is attention interpretable?, in: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 2931–2951, arXiv: [arXiv:1906.03731](https://arxiv.org/abs/1906.03731).
- [49] S. Jain, B.C. Wallace, Attention is not explanation, in: *Proceedings of the 17th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019*, Association for Computational Linguistics (ACL), Minneapolis, Minnesota, USA, 2019, pp. 3543–3556.
- [50] J. Vig, Y. Belinkov, Analyzing the structure of attention in a transformer language model, in: *Proceeding of the ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, Association for Computational Linguistics (ACL), Florence, Italy, 2019, pp. 63–76, <http://dx.doi.org/10.18653/v1/w19-4808>, arXiv: [arXiv:1906.04284](https://arxiv.org/abs/1906.04284).
- [51] K. Clark, U. Khandelwal, O. Levy, C.D. Manning, What does BERT look at? An analysis of BERT's attention, 2019, <http://dx.doi.org/10.18653/v1/w19-4828>, arXiv: [arXiv:1906.04341](https://arxiv.org/abs/1906.04341).
- [52] I. Tenney, D. Das, E. Pavlick, BERT rediscovers the classical NLP pipeline, in: *Proceeding of the 57th Annual Meeting of the Association for Computational Linguistics*, ACL 2019, 2020, pp. 4593–4601, arXiv: [arXiv:1905.05950](https://arxiv.org/abs/1905.05950).
- [53] S. Zhang, X. Zhang, J. Chan, P. Rosso, Irony detection via sentiment-based transfer learning, *Inf. Process. Manage.* 56 (2019) 1633–1644, <http://dx.doi.org/10.1016/j.ipm.2019.04.006>.
- [54] F. Barbieri, F. Ronzano, H. Saggion, Is this tweet satirical? A computational approach for satire detection in Spanish, *Procesamiento de Lenguaje Natural* 55 (2015) 135–142.
- [55] M. d. P. Salas-Zárate, M.A. Paredes-Valverde, M.Á. Rodríguez-García, R. Valencia-García, G. Alor-Hernández, Automatic detection of satire in Twitter: A psycholinguistic-based approach, *Knowl.-Based Syst.* 128 (2017) 20–33, <http://dx.doi.org/10.1016/j.knsys.2017.04.009>.
- [56] R. Ortega, D.I. Rangel, P. Rosso, M. Montes, J.E. Medina, Overview of the task on irony detection in Spanish variants, in: *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019)*, Co-located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2019, CEUR Workshop Proceedings. CEUR-WS.org, Bilbao, Spain, 2019, pp. 229–256.
- [57] L. Seda Mut Altin, A. Bravo, H. Saggion, LaSTUS/TALN at IroSvA: Irony detection in Spanish variants, in: *Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019)*, Co-located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2019, CEUR Workshop Proceedings. CEUR-WS.org, Bilbao, Spain, 2019.

- [58] J.A. González, L.F. Hurtado, Ferran Pla, ELIRF-UPV at IroSvA: Transformer Encoders for Spanish Irony Detection, in: Proceedings of the Iberian Languages Evaluation FÓrum (IberLEF 2019), Co-Located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2019, in: CEUR Workshop Proceedings, CEUR-WS.org, Bilbao, Spain, 2019.
- [59] H.U. Miranda-Belmonte, A.P. López-Monroy, Early fusion of traditional and deep features for irony detection in Twitter, in: Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019), Co-Located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2019, CEUR Workshop Proceedings, CEUR-WS.org, Bilbao, Spain, 2019.
- [60] L. García, D. Moctezuma, V. Muñiz, A contextualized word representation approach for irony detection, in: Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019), Co-Located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2019, in: CEUR Workshop Proceedings, CEUR-WS.org, Bilbao, Spain, 2019.
- [61] H. Calvo, O.J. Gambino, C.V. García, Irony detection using emotion cues, *Comput. Y Sist.* 24 (2020) 1281–1287, <http://dx.doi.org/10.13053/CyS-24-3-3487>.
- [62] M. del Pilar Salas-Zárate, G. Alor-Hernández, J.L. Sánchez-Cervantes, M.A. Paredes-Valverde, J.L. García-Alcaraz, R. Valencia-García, Review of English literature on figurative language applied to social networks, *Knowl. Inf. Syst.* 62 (2020) 2105–2137, <http://dx.doi.org/10.1007/s10115-019-01425-3>.
- [63] M. Abulaish, A. Kamal, M.J. Zaki, A survey of figurative language and its computational detection in online social networks, *ACM Trans. Web* 14 (2020) 1–52, <http://dx.doi.org/10.1145/3375547>.
- [64] J. Karoui, F. Benamara, V. Moriceau, Automatic Detection of Irony, first ed., John Wiley & Sons Inc., 2019, <http://dx.doi.org/10.1002/9781119671183>.
- [65] A. Joshi, P. Bhattacharyya, M.J. Carman, Investigations in computational sarcasm, in: Rüdiger Dillmann, Yoshihiko Nakamura, Stefan Schaal, David Vernon (Eds.), *Springer Nature, Singapore*, 2018.
- [66] P. Carvalho, L. Sarmento, M.J. Silva, E. d. Oliveira, Clues for Detecting Irony in User-generated Contents: Oh!! it's "so easy" :- in: Proceedings of the 1st International Conference on Information Knowledge Management Workshop on Topic-Sentiment Analysis for Mass Opinion, 2009 pp. 53–56.
- [67] D. Davidov, O. Tsur, A. Rappoport, Semi-supervised recognition of sarcastic sentences in Twitter and Amazon, in: Proceedings of the 40th Conference on Computational Natural Language Learning, CoNLL '10, Association for Computational Linguistics, Uppsala, Sweden, 2010, pp. 107–116.
- [68] R. González-Ibáñez, S. Muresan, N. Wacholder, Identifying sarcasm in Twitter: A closer look, in: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, HLT '11, Association for Computational Linguistics, Portland, Oregon, 2011, pp. 581–586.
- [69] F. Kunneman, C. Liebrecht, M. van Mulken, A. van den Bosch, Signaling sarcasm: From hyperbole to hashtag, *Inf. Process. Manage.* 51 (2015) 500–509.
- [70] T. Ptáček, I. Habernal, J. Hong, Sarcasm detection on Czech and English Twitter, in: Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics, Dublin City University and Association for Computational Linguistics, Dublin, Ireland, 2014, pp. 213–223.
- [71] E. Riloff, A. Qadir, P. Surve, L. De Silva, N. Gilbert, R. Huang, Sarcasm as contrast between a positive sentiment and negative situation, in: Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 2013, pp. 704–714.
- [72] F. Barbieri, H. Saggion, Automatic detection of irony and humour in Twitter, in: Proceedings of the 5th International Conference on Computational Creativity, 2014, pp. 155–162.
- [73] F. Barbieri, H. Saggion, Modelling irony in Twitter: Feature analysis and evaluation, in: Proceedings of the 9th International Conference on Language Resources and Evaluation, LREC 2014, 2014, pp. 4258–4264.
- [74] A. Agrawal, A. An, Affective representations for sarcasm detection, in: Proceeding of the 41st International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2018, Association for Computing Machinery, ACM, Ann Arbor, MI, USA, 2018, pp. 1029–1032, <http://dx.doi.org/10.1145/3209978.3210148>.
- [75] D.I. Hernández Farías, V. Patti, P. Rosso, Irony detection in Twitter: The role of affective content, *ACM Trans. Internet Technol.* 16 (2016) 1–24, <http://dx.doi.org/10.1145/2930663>.
- [76] D.I. Hernández Farías, Benedí, Applying basic features from sentiment analysis for automatic irony detection, in: R. Paredes, J.S. Cardoso, X.M. Pardo (Eds.), Proceedings of the Pattern Recognition and Image Analysis, in: Lecture Notes in Computer Science, vol. 9117, Springer International Publishing, Santiago de Compostela, Spain, 2015, pp. 337–344, [http://dx.doi.org/10.1007/978-3-319-19390-8\\_38](http://dx.doi.org/10.1007/978-3-319-19390-8_38).
- [77] F. Barbieri, H. Saggion, F. Ronzano, Modelling sarcasm in Twitter, a novel approach, in: Proceedings of the 5th Workshop on Computational Approaches To Subjectivity, Sentiment and Social Media Analysis, Association for Computational Linguistics, Baltimore, Maryland, USA, 2014, pp. 50–58.
- [78] A. Reyes, P. Rosso, T. Veale, A multidimensional approach for detecting irony in Twitter, *Lang. Resour. Eval.* 47 (2013) 239–268.
- [79] David Bamman, Noah A. Smith, Contextualized sarcasm detection on Twitter, in: Proceedings of the 9th International Conference on Web and Social Media, ICWSM 2015, AAAI, Oxford, UK, 2015, pp. 574–577.
- [80] A. Khattri, A. Joshi, P. Bhattacharyya, M. Carman, Your sentiment precedes you: Using an author's historical tweets to predict sarcasm, in: Proceedings of the 6th Workshop on Computational Approaches To Subjectivity, Sentiment and Social Media Analysis, Association for Computational Linguistics, Lisboa, Portugal, 2015, pp. 25–30.
- [81] B.C. Wallace, D.K. Choe, E. Charniak, Sparse, contextually informed models for irony detection: Exploiting user communities, entities and sentiment, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Association for Computational Linguistics, Beijing, China, 2015, pp. 1035–1044.
- [82] D. Ghosh, A.R. Fabbri, S. Muresan, Sarcasm analysis using conversation context, *Comput. Linguist.* 44 (2018) 755–792, <http://dx.doi.org/10.1162/coli>.
- [83] D.I. Hernández Farías, R. Prati, F. Herrera, P. Rosso, Irony detection in Twitter with imbalanced class distributions, *J. Intell. Fuzzy Systems* (2020) 1–17, <http://dx.doi.org/10.3233/jifs-179880>.
- [84] C. Burfoot, T. Baldwin, Automatic satire detection: Are you having a laugh? in: Proceedings of the Joint Conference of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL and AFNLP, Suntec, Singapore, 2009, pp. 16–164.
- [85] T. Ahmad, H. Akhtar, A. Chopra, M.W. Akhtar, Satire detection from web documents using machine learning methods, in: Proceeding of the International Conference on Soft Computing & Machine Intelligence, IEEE, 2014, pp. 102–105, <http://dx.doi.org/10.1109/ISCMI.2014.34>.
- [86] F. Barbieri, F. Ronzano, H. Saggion, Do we criticise (and Laugh) in the same way? Automatic detection of multi-lingual satirical news in twitter, in: Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, 2015, pp. 1215–1221.
- [87] C.K. Chung, J.W. Pennebaker, Linguistic inquiry and word count (LIWC): Pronounced "Luke," and other useful facts, in: Applied Natural Language Processing: Identification, Investigation and Resolution, 2011, pp. 206–229, <http://dx.doi.org/10.4018/978-1-60960-741-8.ch012>.
- [88] A.N. Reganti, T. Maheshwari, Modeling satire in English text for automatic detection, in: Proceedings of the IEEE 16th International Conference on Data Mining Workshops, IEEE, 2016, pp. 970–977, <http://dx.doi.org/10.1109/ICDMW.2016.146>.
- [89] P.P. Thu, T.N. Aung, Implementation of emotional features on satire detection, *Int. J. Netw. Distrib. Comput.* 6 (2018) 78–87.
- [90] P.P. Thu, N. Nwe, Impact analysis of emotion in figurative language, in: 16th IEEE/ACIS International Conference on Computer and Information Science, ICIS'17, 2017, pp. 209–214.
- [91] O. Levi, P. Hosseini, M. Diab, D.A. Broniatowski, Identifying nuances in fake news vs. Satire: Using semantic and linguistic cues, 2019, <http://dx.doi.org/10.18653/v1/d19-5004>, arXiv:1910.01160.
- [92] G. Guibon, L. Ermakova, H. Seffih, A. Firsov, G.L. Noé-bienvenu, Multi-lingual Fake News Detection with Satire To cite this version : HAL Id : halshs-02391141, International Conference on Computational Linguistics and Intelligent Text Processing (CICLing 2019), La Rochelle, France, 2019.
- [93] A. Ghosh, T. Veale, Fracking sarcasm using neural network, in: Proceedings of the 7th Workshop on Computational Approaches To Subjectivity, Sentiment and Social Media Analysis, Association for Computational Linguistics, San Diego, California, 2016, pp. 161–169, <http://www.aclweb.org/anthology/W16-0425>.
- [94] D. Ghosh, A. Richard Fabbri, S. Muresan, The role of conversation context for sarcasm detection in online interactions, in: Proceedings of the 18th Annual SIGDial Meeting on Discourse and Dialogue, Association for Computational Linguistics, Saarbrücken, Germany, 2017, pp. 186–196, <http://dx.doi.org/10.18653/v1/W17-5523>.
- [95] Y.-H. Huang, H.-H. Huang, H.-H. Chen, Irony detection with attentive recurrent neural networks, in: J.M. Jose, C. Hauff, I.S. Altungovde, D. Song, D. Albakour, S. Watt, J. Tait (Eds.), Proceedings of the Advances in Information Retrieval, Springer International Publishing, Cham, 2017, pp. 534–540.
- [96] A. Joshi, V. Tripathi, K. Patel, P. Bhattacharyya, M.J. Carman, Are Word Embedding-based Features Useful for Sarcasm Detection? in: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November, 2016, 2016, pp. 1006–1011.
- [97] D. Nozza, E. Fersini, E. Messina, Unsupervised irony detection: A probabilistic model with word embeddings, in: Proceedings of the 8th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management, Vol. 1, 2016, pp. 68–76, <http://dx.doi.org/10.5220/0006052000680076>.



- [98] S. Poria, E. Cambria, D. Hazarika, P. Vij, A deeper look into sarcastic tweets using deep convolutional neural networks, in: Proceedings of the COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers Pp, Association for Computational Linguistics, Osaka, Japan, 2016, pp. 1601–1612, URL <https://www.aclweb.org/anthology/C16-1151>.
- [99] D. Hazarika, S. Poria, S. Gorantla, E. Cambria, R. Zimmermann, R. Mihalcea, Cascade: Contextual sarcasm detection in online discussion forums, in: Proceedings of the 27th International Conference on Computational Linguistics, Association for Computational Linguistics (ACL), Santa Fe, New Mexico, USA, 2018, pp. 1837–1848, URL <https://www.aclweb.org/anthology/C18-1156>, arXiv:1805.06413.
- [100] C. Wu, F. Wu, S. Wu, J. Liu, Z. Yuan, Y. Huang, THU\_NGN at SemEval-2018 task 3: Tweet irony detection with densely connected LSTM and multi-task learning, in: Proceeding of the 12th International Workshop on Semantic Evaluation, Association for Computational Linguistics (ACL), New Orleans, Louisiana, 2018, pp. 51–56, <http://dx.doi.org/10.18653/v1/s18-1006>.
- [101] C. Baziotis, A. Nikolaos, P. Papalampidi, A. Kolovou, G. Paraskevopoulos, N. Ellinas, A. Potamianos, NTUA-SLP At SemEval-2018 task 3: Tracking ironic tweets using ensembles of word and character level attentive RNNs, 2018, pp. 613–621, arXiv:1804.06659.
- [102] C. Zhang, M. Abdul-Mageed, Multi-task bidirectional transformer representations for irony detection, in: CEUR Workshop Proceedings, Vol. 2517, 2019, pp. 391–400, arXiv:1909.03526.
- [103] L.S.M. Altin, A. Bravo, H. Saggion, LatUS/TALN at IroSvA: Irony detection in spanish variants, in: Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019), Co-located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2019, in: CEUR Workshop Proceedings, vol. 2421 (2019) 291–296.
- [104] K. Ravi, V. Ravi, Irony detection using neural network language model, psycholinguistic features and text mining, 2018.
- [105] Q. Le, T. Mikolov, Distributed representations of sentences and documents, in: Proceedings of the 31st International Conference on International Conference on Machine Learning, Vol. 32, Beijing, China, 2014, pp. II–1188–II–1196, volume 4, arXiv:1405.4053.
- [106] H. Hotelling, Relations between two sets of variates, *Biometrika* 28 (1936) 321–377, <http://dx.doi.org/10.2307/2333955>.
- [107] A. Kumar, V.T. Narapareddy, V.A. Srikanth, A. Malapati, L.B.M. Neti, Sarcasm detection using multi-head attention based bidirectional LSTM, *IEEE Access* 8 (2020) 6388–6397, <http://dx.doi.org/10.1109/ACCESS.2019.2963630>.
- [108] N. Majumder, S. Poria, H. Peng, N. Chhaya, E. Cambria, A. Gelbukh, Sentiment and sarcasm classification with multitask learning, *IEEE Intell. Syst.* 34 (2019) 38–43, <http://dx.doi.org/10.1109/MIS.2019.2904691>.
- [109] D.S.R. Chauhan, A. Ekbal, P. Bhattacharyya, Sentiment and emotion help sarcasm? a multi-task learning framework for multi-modal sarcasm, sentiment and emotion analysis, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Online, 2020, pp. 4351–4360, <http://dx.doi.org/10.18653/v1/2020.acl-main.401>, URL <https://www.aclweb.org/anthology/2020.acl-main.401>.
- [110] F. Yang, A. Mukherjee, E. Gragut, Satirical news detection and analysis using attention mechanism and linguistic features, in: Proceeding of the Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Association for Computational Linguistics, Copenhagen, Denmark, 2017a, pp. 1979–1989.
- [111] S.D. Sarkar, F. Yang, A. Mukherjee, Attending sentences to detect satirical fake news, in: 27th International Conference on Computational Linguistics, COLING'18, 2018, pp. 3371–3380.
- [112] S. Dutta, A. Chakraborty, A deep learning-inspired method for social media satire detection, in: Proceeding of the Soft Computing and Signal Processing, Advances in Intelligent Systems and Computing, Springer Singapore, 2019, pp. 243–251, <http://dx.doi.org/10.1007/978-981-13-3393-4>, [http://dx.doi.org/10.1007/978-981-13-3393-4\\_25](http://dx.doi.org/10.1007/978-981-13-3393-4_25).
- [113] Y.-J. Tang, H.-H. Chen, Chinese irony corpus construction and ironic structure analysis, in: Proceedings of the COLING 2014, the 25th International Conference on Computational Linguistics, Association for Computational Linguistics, Dublin, Ireland, 2014, pp. 1269–1278.
- [114] J. Karoui, F. Benamara, V. Moriceau, N. Aussenac-Gilles, L. Hadrich-Belguith, Towards a contextual pragmatic model to detect irony in tweets, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), Association for Computational Linguistics, 2015, pp. 644–650.
- [115] F. Benamara, C. Grouin, J. Karoui, V. Moriceau, I. Robba, Analyse d'Opinion et Langage Figuratif dans des Tweets : Présentation et Résultats du Défi Fouille de Textes DEFT2017, in: Actes de l'atelier DEFT2017 Associé à la Conférence TALN, Orléans, France, 2017.
- [116] C. Bosco, V. Patti, A. Bolioli, Developing corpora for sentiment analysis : The case of irony and senti-TUT, *IEEE Intell. Syst.* 28 (2013) 55–63.
- [117] AC Cignarella, Simona F. V. Basile, C. Bosco, V. Patti, P. Rosso, S. Frenda, V. Basile, C. Bosco, V. Patti, P. Rosso, Overview of the evalita 2018 task on irony detection in Italian tweets (IronITA), in: Proceedings of the 6th Evaluation Campaign of Natural Language Processing and Speech Tools for Italian, EVALITA'18, CEUR.org, Turin, Italy, 2018.
- [118] F. Rangel, D.I. Hernández. Fariás, P. Rosso, Emotions and irony per gender in Facebook, in: Proceeding of the Workshop on Emotion, Social Signals, Sentiment & Linked Open Data (ES3LOD), LREC-2014 pp. 68–73. Reykjavík, Iceland, 2014.
- [119] G. Jasso López, I. Meza Ruiz, Character and word baselines systems for irony detection in Spanish short texts, *Procesamiento Del Lenguaje Nat.* 56 (2016) 41–48, URL <http://journal.sepln.org/sepln/ojs/ojs/index.php/pln/article/view/5285>.
- [120] J. Karouia, F.B. Zitoun, Veronique Moriceau, SOUKHRIA: Towards an irony detection system for arabic in social media, in: Proceedings of the 3rd International Conference on Arabic Computational Linguistics, ACLing 2017, Association for Computational Linguistic (ACL), Dubai, United Arab Emirates, 2017, pp. 116–168.
- [121] B. Ghanem, J. Karoui, F. Benamara, V. Moriceau, P. Rosso, IDAT@FIRE2019: Overview of the track on irony detection in arabic tweets, in: Proceedings of 11th Forum for Information Retrieval Evaluation, CEURS, 2019, pp. 1–11, <http://dx.doi.org/10.1145/3368567.3368585>.
- [122] R.K. Singh, M.K. Sachan, R. Patel, 360 degree view of cross-domain opinion classification: a survey, *Artif. Intell. Rev.* 54 (2021) 1385–1506.
- [123] A. Esuli, A. Moreo, F. Sebastiani, Cross-lingual sentiment quantification, *IEEE Intell. Syst.* 35 (2020) 106–114, <http://dx.doi.org/10.1109/MIS.2020.2979203>.
- [124] S. Galeshchuk, J. Qiu, J. Jourdan, Sentiment analysis for multilingual corpora, in: Proceedings of the 7th Workshop on Balto-Slavic Natural Language Processing, Association for Computational Linguistics, Florence, Italy, 2019, pp. 120–125, <http://dx.doi.org/10.18653/v1/W19-3717>, URL <https://www.aclweb.org/anthology/W19-3717>.
- [125] S.L. Lo, E. Cambria, R. Chiong, D. Cornforth, Multilingual sentiment analysis: from formal to informal and scarce resource languages, *Artif. Intell. Rev.* 48 (2017) 499–527, <http://dx.doi.org/10.1007/s10462-016-9508-4>.
- [126] M. Abdalla, G. Hirst, Cross-lingual sentiment analysis without (good) translation, in: Proceedings of the 8th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 506–515, URL <https://www.aclweb.org/anthology/I17-1051>.
- [127] K. Dashtipour, S. Poria, A. Hussain, E. Cambria, A.Y. Hawalah, A. Gelbukh, Q. Zhou, Multilingual sentiment analysis: state of the art and independent comparison of techniques, *Cognitive Comput.* 8 (2016) 757–771.
- [128] A. Balahur, M. Turchi, Multilingual sentiment analysis using machine translation? In Proceedings of the 3rd Workshop in Computational Approaches to Subjectivity and Sentiment Analysis, WASSA'13 2012, pp. 52–60.
- [129] J. Karoui, F. Benamara, V. Moriceau, V. Patti, C. Bosco, N. Aussenac-Gilles, Exploring the impact of pragmatic phenomena on irony detection in tweets: A multilingual corpus study, in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017, 2017, pp. 262–272, <http://dx.doi.org/10.18653/v1/e17-1025>.
- [130] A.T. Cignarella, V. Basile, M. Sanguinetti, C. Bosco, P. Rosso, F. Benamara, Multilingual irony detection with dependency syntax and neural models, in: Proceedings of the 28th International Conference on Computational Linguistics, Association for Computational Linguistics (ACL), Barcelona, Spain, 2020, pp. 1346–1358, <http://dx.doi.org/10.18653/v1/2020.coling-main.116>, arXiv:2011.05706.
- [131] B. Ghanem, J. Karoui, F. Benamara, P. Rosso, V. Moriceau, Irony detection in a multilingual context, advances in information retrieval, in: Proceedings of 42nd European Conference on IR Research, ECIR 2020, 2020, pp. 114–149, <http://dx.doi.org/10.1007/978-3-030-45442-5>, arXiv:2003.13924.
- [132] P. Zhou, W. Shi, J. Tian, Z. Qi, B. Li, H. Hao, B. Xu, Attention-based bidirectional long short-term memory networks for relation classification, in: Proceeding of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016 - Short Papers, Association for Computational Linguistics (ACL), Berlin, Germany, 2016, pp. 207–212, <http://dx.doi.org/10.18653/v1/p16-2034>.
- [133] J. Perkins, Python 3 Text Processing with NLTK 3 Cookbook, Packt Publishing, 2014, arXiv:arXiv:1011.1669v3.
- [134] L. Padró, E. Stanilovsky, FreeLing 3.0: Towards Wider Multilinguality, in: Proceedings of the LREC 2012, 2012.
- [135] F. Barbieri, L.E. Anke, H. Saggion, Revealing patterns of twitter emoji usage in barcelona and madrid, International Conference of the Catalan Association for Artificial Intelligence, 2016.
- [136] F. Barbieri, G. Kruszewski, F. Ronzano, H. Saggion, How cosmopolitan are emojis? exploring emojis usage and meaning over different languages with distributional semantics, in: Proceedings of the 24th ACM International Conference on Multimedia, MM '16, Association for Computing

- Machinery, New York, NY, USA, 2016, pp. 531–535, <http://dx.doi.org/10.1145/2964284.2967278>.
- [137] F. Barbieri, F. Ronzano, H. Saggion, What does this emoji mean? a vector space skip-gram model for Twitter emojis, in: Proceedings of the 10th International Conference on Language Resources and Evaluation, LREC'16, European Language Resources Association (ELRA), Portorož, Slovenia, 2016d, pp. 3967–3972, URL <https://www.aclweb.org/anthology/L16-1626>.
- [138] F. Barbieri, M. Ballesteros, H. Saggion, Are emojis predictable? in: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, 2017a, pp. 105–111, URL <https://www.aclweb.org/anthology/E17-2017>.
- [139] F. Barbieri, L. Espinosa-Anke, M. Ballesteros, J. Soler-Company, H. Saggion, Towards the understanding of gaming audiences by modeling twitch emotes, in: Proceedings of the 3rd Workshop on Noisy User-Generated Text, Association for Computational Linguistics, Copenhagen, Denmark, 2017b, pp. 11–20, <http://dx.doi.org/10.18653/v1/W17-4402>, URL <https://www.aclweb.org/anthology/W17-4402>.
- [140] M. Pota, M. Ventura, R. Catelli, M. Esposito, An effective bert-based pipeline for twitter sentiment analysis: A case study in italian, Sensors 21 (2021) <http://dx.doi.org/10.3390/s21010133>, URL <https://www.mdpi.com/1424-8220/21/1/133>.
- [141] M. Pota, M. Ventura, H. Fujita, M. Esposito, Multilingual evaluation of pre-processing for bert-based sentiment analysis of tweets, Expert Syst. Appl. 181 (2021b) 115119, <http://dx.doi.org/10.1016/j.eswa.2021.115119>, URL <https://www.sciencedirect.com/science/article/pii/S0957417421005601>.
- [142] D. Cer, Y. Yang, S. y. Kong, N. Hua, N. Limtiaco, R. St. John, N. Constant, M. Guajardo-Céspedes, S. Yuan, C. Tar, Y.H. Sung, B. Strope, R. Kurzweil, Universal sentence encoder for English, in: Proceedings of the Conference on Empirical Methods in Natural Language Processing: System Demonstrations, 2018, pp. 169–174, [arXiv:arXiv:1803.11175v2](https://arxiv.org/abs/1803.11175v2).
- [143] Y. Yang, D. Cer, A. Ahmad, M. Guo, J. Law, N. Constant, G. Hernandez, Abrego, S. Yuan, C. Tar, Y.-h. Sung, B. Strope, R. Kurzweil, Multilingual universal sentence encoder for semantic retrieval, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, Association for Computational Linguistics, Online, 2020, pp. 87–94, <http://dx.doi.org/10.18653/v1/2020.acl-demos.12>, URL <https://www.aclweb.org/anthology/2020.acl-demos.12>.
- [144] R. Ortega-Bueno, C.E. Muñiz, P. Rosso, J.E. Medina-Pagola, Uo\_Upv : Deep linguistic humor detection in spanish social media, in: P. Rosso, J. Gonzalo, R. Martínez, S. Montalvo, J. Carrillo de Albornoz (Eds.), Proceedings of the 3rd Workshop on Evaluation of Human Language Technologies for Iberian Languages (IberEval 2018) Co-Located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2018, CEUR-WS.org, Sevilla, Spain, 2018, pp. 203–213.
- [145] R. Ortega-Bueno, P. Rosso, J.E. Medina Pagola, UO\_UPV2 at HABA 2019: BiGRU neural network informed with linguistic features for humor recognition, in: CEUR Workshop Proceedings, CEUR Workshop Proceedings, CEUR-WS.org, Bilbao, Spain, 2019.
- [146] R. Ortega-Bueno, J.E. Medina Pagola, UO\_IRO: Linguistic informed deep-learning model for irony detection, in: CEUR Workshop Proceedings, Vol. 2263, CEUR Workshop Proceedings, CEUR-WS.org, 2018, pp. 1–6, <http://dx.doi.org/10.4000/books.aaccademia.4638>.
- [147] D. Vilares, H. Peng, R. Satapathy, E. Cambria, Babelsentinet: a commonsense reasoning framework for multilingual sentiment analysis, in: Proceedings of the IEEE Symposium Series on Computational Intelligence, SSCI, IEEE, 2018, pp. 1292–1298.
- [148] E. Cambria, Y. Li, F.Z. Xing, S. Poria, K. Kwok, Senticnet 6: Ensemble application of symbolic and subsymbolic ai for sentiment analysis, in: Proceedings of the 29th ACM International Conference on Information & Knowledge Management, Association for Computing Machinery, New York, NY, USA, 2020, pp. 105–114.
- [149] M. Iyyer, V. Manjunatha, J. Boyd-Graber, H. Daumé, Deep unordered composition rivals syntactic methods for text classification, in: Proceedings of the 53rd Annual Meeting Of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, Association for Computational Linguistics, Beijing, China, 2015, pp. 1681–1691.
- [150] Q. Le, T. Mikolov, Distributed representations of sentences and documents, in: Proceedings of the 31st International Conference on International Conference on Machine Learning, Vol. 32, ICML'14, JMLR.org, 2014, II–1188–II–1196.
- [151] M. Pagliardini, P. Gupta, M. Jaggi, Unsupervised learning of sentence embeddings using compositional n-gram features, in: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, New Orleans, Louisiana, 2018, pp. 528–540, <http://dx.doi.org/10.18653/v1/N18-1049>, URL <https://www.aclweb.org/anthology/N18-1049>.
- [152] P. Bojanowski, E. Grave, A. Joulin, T. Mikolov, Enriching word vectors with subword information, Trans. ACL 5 (2017) 135–146.
- [153] A. Conneau, D. Kiela, H. Schwenk, L. Barrault, A. Bordes, Supervised learning of universal sentence representations from natural language inference data, in: Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Copenhagen, Denmark, 2017, pp. 670–680, <http://dx.doi.org/10.18653/v1/D17-1070>, URL <https://www.aclweb.org/anthology/D17-1070>.
- [154] T. Mikolov, I. Sutskever, K. Chen, G.S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, Adv. Neural Inf. Process. Syst. (2013) 3111–3119.
- [155] J. Pennington, R. Socher, C.D. Manning, GloVe: Global vectors for word representation, in: Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP, Association for Computational Linguistics, Doha, Qatar, 2014, pp. 1532–1543.
- [156] M.E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, L. Zettlemoyer, Deep contextualized word representations, in: Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Vol. 1, 2018, pp. 2227–2237, [arXiv:1802.05365](https://arxiv.org/abs/1802.05365).
- [157] J. Howard, S. Ruder, Universal language model fine-tuning for text classification, in: Proceedings of the 56th Annual Meeting Of the Association for Computational Linguistics, Association for Computational Linguistics, Melbourne, Australia, 2018, pp. 328–339, <http://dx.doi.org/10.3760/cma.j.issn.04124081.2010.02.006>.
- [158] D.P. Kingma, J. Ba, Adam: A method for stochastic optimization, in: Y. Bengio, Y. LeCun (Eds.), Proceedings of the 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings, 2015, URL <http://arxiv.org/abs/1412.6980>.
- [159] D. Bahdanau, K.H. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate, in: Proceedings of the 3rd International Conference on Learning Representations, ICLR 2015, 2015, pp. 1–15, [arXiv:1409.0473](https://arxiv.org/abs/1409.0473).
- [160] O. Iroay, C. Cardie, Deep recursive neural networks for compositionality in language, Adv. Neural Inf. Process. Syst. 3 (2014) 2096–2104, URL <http://www.scopus.com/inward/record.url?eid=s-2-s2.0-84937828128&partnerID=tZOTx3y1>.
- [161] T. Rocktäschel, E. Grefenstette, K.M. Hermann, T. Kočiský, P. Blunsom, Reasoning about entailment with neural attention, in: 4th International Conference on Learning Representations, ICLR 2016, 2016, pp. 1–9, [arXiv:1509.06664](https://arxiv.org/abs/1509.06664).
- [162] L. Chiruzzo, S. Castro, A. Rosá, HABA 2019 dataset: A corpus for humor analysis in spanish, in: Proceedings of the 12th Conference on Language Resources and Evaluation, LREC 2020, European Language Resources Association (ELRA), Marseille, France, 2020, pp. 5106–5112.
- [163] L. Chiruzzo, S. Castro, M. Etcheverry, D. Garat, J.J. Prada, A. Rosá, Overview of HABA at IberLEF 2019: Humor analysis based on human annotation, in: Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019), Co-Located with 34th Conference of the Spanish Society for Natural Language Processing, SEPLN 2019, CEUR Workshop Proceedings, CEUR-WS.org, Bilbao, Spain, 2019, pp. 132–144.
- [164] W. Haynes, Wilcoxon rank sum test, in: W. Dubitzky, O. Wolkenhauer, K.-H. Cho, H. Yokota (Eds.), Encyclopedia of Systems Biology, Springer New York, New York, NY, 2013, pp. 2354–2355, [http://dx.doi.org/10.1007/978-1-4419-9863-7\\_1185](http://dx.doi.org/10.1007/978-1-4419-9863-7_1185).
- [165] G. Sidorov, S. Miranda-Jiménez, F. Viveros-Jiménez, A. Gelbukh, N. Castro-Sánchez, F. Velásquez, I. Dí az Rangel, S. Suárez-Guerra, A. Treviño, J. Gordon, Empirical study of machine learning based approach for opinion mining in tweets, in: Proceedings of the 11th Mexican International Conference on Advances in Artificial Intelligence - Volume Part I, MICAI'12, Springer-Verlag, Berlin, Heidelberg, 2013, pp. 1–14, [http://dx.doi.org/10.1007/978-3-642-37807-2\\_1](http://dx.doi.org/10.1007/978-3-642-37807-2_1).
- [166] L. Hernández, A. López-Lopez, J.E. Medina-Pagola, Classification of attitude words for opinions mining, Int. J. Comput. Linguistics Appl. 2 (2011) 267–283.
- [167] L.R. Gurillo, M.B.A. Ortega, S.R. Rosique, S. Attardo, E.M.-G. Paredes, X.A. Padilla-García, J. Muñoz Basols, P. Adrjan, M. David, A. Viana, K. Feysaerts, F. Yus, Irony and Humor: From Pragmatics To Discourse Volume 231, John Benjamins Publishing Company, Amsterdam / Philadelphia, 2013.
- [168] J. Garmendia, Irony, first ed., Cambridge University Press, New York, USA, 2018, <http://dx.doi.org/10.1017/9781316136218>.
- [169] G.A. Miller, Wordnet: a lexical database for English, Commun. ACM 38 (1995) 39–41.
- [170] A. Gonzalez-Agirre, E. Laparra, G. Rigau, Multilingual central repository version 3.0, LREC (2012) 2525–2529.
- [171] M.D.M. González, E.M. Cámara, M.T.M. Valdivia, CRISOL: Base de conocimiento de opiniones para el español, Procesamiento Del Lenguaje Natural (2015) 143–150.
- [172] X. Saralegi, I.S. Vicente, Elhuyar at TASS 2013, in: XXIX Congreso de la Sociedad Española de Procesamiento de Lenguaje Natural, Workshop on Sentiment Analysis at SEPLN, TASS2013, 2013, pp. 143–150.

- [173] A. Hogenboom, D. Bal, F. Frasincar, Exploiting emoticons in sentiment analysis, in: Proceedings of the 28th Annual ACM Symposium on Applied Computing, SAC '13, ACM, New York, NY, USA, ISBN: 978-1-4503-1656-9, 2013, pp. 703–710, <http://dx.doi.org/10.1145/2480362.2480498>, URL <http://doi.acm.org/10.1145/2480362.2480498>.
- [174] R. Ortega-Bueno, J.E. Medina-Pagola, C.E. Muñiz Cuza, P. Rosso, Improving attitude words classification for opinion mining using word embedding, in: R. Vera-Rodríguez, J. Fierrez, A. Morales (Eds.), Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications - 23rd Iberoamerican Congress, CIARP 2018, Madrid, Spain, November 19–22, 2018, Proceedings, in: Lecture Notes in Computer Science, vol. 11401, Springer, 2018, pp. 971–982, [http://dx.doi.org/10.1007/978-3-030-13469-3\\_112](http://dx.doi.org/10.1007/978-3-030-13469-3_112).
- [175] Y. Susanto, A.G. Livingstone, B.C. Ng, E. Cambria, The hourglass model revisited, IEEE Intell. Syst. 35 (2020) 96–102, <http://dx.doi.org/10.1109/MIS.2020.2992799>.
- [176] J. W. Pennebaker, James W. Volume 1890, first ed., Bloomsbury Press, 2011.
- [177] A.S. Peña, L.A. García, A.R. Dosina, Detección de ironía en textos cortos enfocada a la minería de opinión, in: IV Conferencia Internacional en Ciencias Computacionales e Informáticas (CICCI' 2018) 1–10. Havana, Cuba, 2018.