

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/337531908>

A Transformer-based approach to Irony and Sarcasm detection

Preprint · November 2019

DOI: 10.48550/arXiv.1911.10401

CITATIONS

0

READS

733

3 authors, including:



Rolandos Alexandros Potamias
Imperial College London

27 PUBLICATIONS 351 CITATIONS

[SEE PROFILE](#)



Georgios Siolas
National Technical University of Athens

54 PUBLICATIONS 728 CITATIONS

[SEE PROFILE](#)

A Transformer-based approach to Irony and Sarcasm detection

Rolandos Alexandros Potamias · Georgios Siolas · Andreas - Georgios Stafylopatis

Received: 12 November 2019 / Accepted: 4 June 2020

Abstract Figurative Language (FL) seems ubiquitous in all social-media discussion forums and chats, posing extra challenges to sentiment analysis endeavors. Identification of FL schemas in short texts remains largely an unresolved issue in the broader field of Natural Language Processing (NLP), mainly due to their contradictory and metaphorical meaning content. The main FL expression forms are sarcasm, irony and metaphor. In the present paper we employ advanced Deep Learning (DL) methodologies to tackle the problem of identifying the aforementioned FL forms. Significantly extending our previous work [74], we propose a neural network methodology that builds on a recently proposed pre-trained transformer-based network architecture which, is further enhanced with the employment and devise of a recurrent convolutional neural network (RCNN). With this set-up, data preprocessing is kept in minimum. The performance of the devised hybrid neural architecture is tested on four benchmark datasets, and contrasted with other relevant state of the art methodologies and systems. Results demonstrate that the proposed methodology achieves state of the art performance under all benchmark datasets, outperforming,

even by a large margin, all other methodologies and published studies.

Keywords Sentiment Analysis · Natural Language Processing · Figurative Language · Sarcasm · Irony · Deep Learning · Transformer networks

1 Introduction

In the networked-world era the production of (structured or unstructured) data is increasing with most of our knowledge being created and communicated via web-based social channels [95]. Such data explosion raises the need for efficient and reliable solutions for the management, analysis and interpretation of huge data sizes. Analyzing and extracting knowledge from massive data collections is not only a big issue per-se, but also challenges the data analytics state-of-the-art [102], with statistical and machine learning methodologies paving the way, and deep learning (DL) taking over and presenting highly accurate solutions [29]. Relevant applications in the field of social media cover a wide spectrum, from the categorization of major disasters [42] and the identification of suggestions [72] to inducing users appeal to political parties [2].

The raising of computational social science [55], and mainly its social media dimension [66], challenge contemporary computational linguistics and text-analytics endeavors. The challenge concerns the advancement of text analytics methodologies towards the transformation of unstructured excerpts into some kind of structured data via the identification of special passage characteristics, such as its emotional content (e.g., anger, joy, sadness) [48]. In this context, Sentiment Analysis (SA) comes into play, targeting the devise and development of efficient algorithmic processes for the automatic

Rolandos Alexandros Potamias†
Department of Computing,
Imperial College London, United Kingdom
E-mail: r.potamias@imperial.ac.uk
† Work performed while at National Technical University of Athens.

Georgios Siolas
School of Electrical and Computer Engineering,
National Technical University of Athens, Greece
E-mail: gsiolas@islab.ntua.gr

Andreas - Georgios Stafylopatis
School of Electrical and Computer Engineering,
National Technical University of Athens, Greece
E-mail: andreas@cs.ntua.gr

extraction of a writers sentiment or emotion as conveyed in text excerpts. Relevant efforts focus on tracking the sentiment polarity of single utterances, which in most cases is loaded with a lot of subjectivity and a degree of vagueness [57]. Contemporary research in the field utilizes data from social media resources (e.g., Facebook, Twitter) as well as other short text references in blogs, forums etc [73]. However, users of social media tend to violate common grammar and vocabulary rules and even use various figurative language forms to communicate their message. In such situations, the sentiment inclination underlying the literal content of the conveyed concept may significantly differ from its figurative context, making SA tasks even more puzzling. Evidently, single turn text lack in detecting sentiment polarity on sarcastic and ironic expressions, as already signified in the relevant SemEval-2014 Sentiment Analysis task 9 [81]. Moreover, lacking of facial expressions and voice tone require context aware approaches to tackle such a challenging task and overcome its ambiguities [31]. As sentiment is the emotion behind customer engagement, SA finds its realization in automated customer aware services, elaborating over users emotional intensities [13]. Most of the related studies utilize single turn texts from topic specific sources, such as Twitter, Amazon, IMDB etc. Hand crafted and sentiment-oriented features, indicative of emotion polarity, are utilized to represent respective excerpt cases. The formed data are then fed traditional machine learning classifiers (e.g. SVM, Random Forest, multilayer perceptrons) or DL techniques and respective complex neural architectures, in order to induce analytical models that are able to capture the underlying sentiment content and polarity of passages [32, 82, 41].

The linguistic phenomenon of figurative language (FL) refers to the contradiction between the literal and the non-literal meaning of an utterance [17]. Literal written language assigns exact (or real) meaning to the used words (or phrases) without any reference to putative speech figures. In contrast, FL schemas exploit non-literal mentions that deviate from the exact concept presented by the used words and phrases. FL is rich of various linguistic phenomena like metonymy reference to an entity stands for another of the same domain, a more general case of synonymy; and metaphors systematic interchange between entities from different abstract domains [18]. Besides the philosophical considerations, theories and debates about the exact nature of FL, findings from the neuroscience research domain present clear evidence on the presence of differentiating FL processing patterns in the human brain [94, 59, 45, 6, 13], even for woman-man attraction situations! [23]. A fact that makes FL processing even more challeng-

ing and difficult to tackle. Indeed, this is the case of pragmatic FL phenomena like irony and sarcasm that main intention of in most of the cases, are characterized by an oppositeness to the literal language context. It is crucial to distinguish between the literal meaning of an expression considered as a whole from its constituents words and phrases. As literal meaning is assumed to be invariant in all context at least in its classical conceptualization [46], it is exactly this separation of an expression from its context that permits and opens the road to computational approaches in detecting and characterizing FL utterance.

We may identify three common FL expression forms namely, irony, sarcasm and metaphor. In this paper, figurative expressions, and especially ironic or sarcastic ones, are considered as a way of indirect denial. From this point of view, the interpretation and ultimately identification of the indirect meaning involved in a passage does not entail the cancellation of the indirectly rejected message and its replacement with the intentionally implied message (as advocated in [12, 30]). On the contrary ironic/sarcastic expressions presupposes the processing of both the indirectly rejected and the implied message so that the difference between them can be identified. This view differs from the assumption that irony and sarcasm involve only one interpretation [91, 83]. Holding that irony activates both grammatical / explicit as well as ironic / involved notions provides that irony will be more difficult to grasp than a non-ironic use of the same expression.

Despite that all forms of FL are well studied linguistic phenomena [91], computational approaches fail to identify the polarity of them within a text. The influence of FL in sentiment classification emerged both on SemEval-2014 Sentiment Analysis task [81] and [18]. Results show that Natural Language Processing (NLP) systems effective in most other tasks see their performance drop when dealing with figurative forms of language. Thus, methods capable of detecting, separating and classifying forms of FL would be valuable building blocks for a system that could ultimately provide a full-spectrum sentiment analysis of natural language.

In literature we encounter some major drawbacks of previous studies and we aim to resolve with our proposed method:

- Many studies tackle figurative language by utilizing a wide range of engineered features (e.g. lexical and sentiment based features) [21, 28, 74, 76, 77, 85] making classification frameworks not feasible.
- Several approaches search words on large dictionaries which demand large computational times and can be considered as impractical [74, 85]

- Many studies exhaustively preprocess the input texts, including stemming, tagging, emoji processing etc. that tend to be time consuming especially in large datasets [51, 89].
- Many approaches attempt to create datasets using social media APIs to automatically collect data rather than exploiting their system on benchmark datasets, with proven quality. To this end, it is impossible to be compared and evaluated [51, 56, 89].

To tackle the aforementioned problems, we propose an end-to-end methodology containing none hand crafted engineered features or lexicon dictionaries, a preprocessing step that includes only de-capitalization and we evaluate our system on several benchmark dataset. To the best of our knowledge, this is the first time that an unsupervised pre-trained Transformer method is used to capture figurative language in many of its forms.

The rest of the paper is structured as follows, in Section 2 we present the related work on the field of FL detection, in Section 3 we shortly describe the background of recent advances in natural language processing that achieve high performance in a wide range of tasks and will be used to compare performance, in 4 we present our proposed method, the results of our experiments are presented in Section 4, and finally our conclusion is in Section 6.

2 Literature Review

Although the NLP community have researched all aspects of FL independently, none of the proposed systems were evaluated on more than one type. Related work on FL detection and classification tasks could be categorized into two main categories, according to the studied task: (a) irony and sarcasm detection, and (b) sentiment analysis of FL excerpts. Even if sarcasm and irony are not identical phenomena, we will present those types together, as they appear in the literature.

2.1 Irony and Sarcasm Detection

Recently, the detection of ironic and sarcastic meanings from respective literal ones have raised scientific interest due to the intrinsic difficulties to differentiate between them. Apart from English language, irony and sarcasm detection have been widely explored on other languages as well, such as Italian [84], Japanese [35], Spanish [67], Greek [10] etc. In the review analysis that follows we group related approaches according to the their adopted key concepts to handle FL.

Approaches based on unexpectedness and contradictory factors. Reyes et al. [78, 79] were the first

that attempted to capture irony and sarcasm in social media. They introduced the concepts of unexpectedness and contradiction that seems to be frequent in FL expressions. The unexpectedness factor was also adopted as a key concept in other studies as well. In particular, Barbieri et al. [4] compared tweets with sarcastic content with other topics such as, #politics, #education, #humor. The measure of unexpectedness was calculated using the *American National Corpus Frequency Data* source as well as the morphology of tweets, using Random Forests (RF) and Decision Trees (DT) classifiers. In the same direction, Buschmeir et al. [7] considered unexpectedness as an emotional imbalance between words in the text. Ghosh et al. [26] identified sarcasm using Support Vector Machines (SVM) using as features the identified contradictions within each tweet.

Content and context-based approaches. Inspired by the contradictory and unexpectedness concepts, follow-up approaches utilized features that expose information about the content of each passage including: N-gram patterns, acronyms and adverbs [8]; semi-supervised attributes like word frequencies [16]; statistical and semantic features [77]; and *Linguistic Inquiry and Word Count* (LIWC) dictionary along with syntactic and psycholinguistic features [75]. LIWC corpus [68] was also utilized in [28], comparing sarcastic tweets with positive and negative ones using an SVM classifier. Similarly, using several lexical resources [85], and syntactic and sentiment related features [56], the respective researchers explored differences between sarcastic and ironic expressions. Affective and structural features are also employed to predict irony with conventional machine learning classifiers (DT, SVM, Nave Bayes/NB) in [20]. In a follow-up study [21], a knowledge-based k-NN classifier was fed with a feature set that captures a wide range of linguistic phenomena (e.g., structural, emotional). Significant results were achieved in [89], were a combination of lexical, semantic and syntactic features passed through an SVM classifier that outperformed LSTM deep neural network approaches. Apart from local content, several approaches claimed that global context may be essential to capture FL phenomena. In particular, in [92] it is claimed that capturing previous and following comments on Reddit increases classification performance. Users behavioral information seems to be also beneficial as it captures useful contextual information in Twitter post [76]. A novel unsupervised probabilistic modeling approach to detect irony was also introduced in [65].

Deep Learning approaches. Although several DL methodologies, such as recurrent neural networks (RNNs),

are able to capture hidden dependencies between terms within text passages and can be considered as content-based, we grouped all DL studies for readability purposes. Word Embeddings, i.e., learned mappings of words to real valued vectors [61], play a key role in the success of RNNs and other DL neural architectures that utilize pre-trained word embeddings to tackle FL. In fact, the combination of word embeddings with Convolutional Neural Networks (CNN), so called CNN-LSTM units, was introduced by Kumar [52] and Ghosh & Veale [25] achieving state-of-the-art performance. Attentive RNNs exhibit also good performance when matched with pre-trained Word2Vec embeddings [38], and contextual information [101]. Following the same approach an LSTM based intra-attention was introduced in [87] that achieved increased performance. A different approach, founded on the claim that number present significant indicators, was introduced by Dubey et al. [19]. Using an attentive CNN on a dataset with sarcastic tweets that contain numbers, showed notable results. An ensemble of a shallow classifier with lexical, pragmatic and semantic features, utilizing a Bidirectional LSTM model is presented in [50]. In a subsequent study [51], the researchers engineered a soft attention LSTM model coupled with a CNN. Contextual DL approaches are also employed, utilizing pre-trained along with user embeddings structured from previous posts [1] or, personality embeddings passed through CNNs [33]. ELMo embeddings [71] are utilized in [39]. In our previous approach we implemented an ensemble deep learning classifier (DESC) [74], capturing content and semantic information. In particular, we employed an extensive feature set of a total 44 features leveraging syntactic, demonstrative, sentiment and readability information from each text along with Tf-idf features. In addition, an attentive bidirectional LSTM model trained with GloVe pre-trained word embeddings was utilized to structure an ensemble classifier processing different text representations. DESC model performed state-of-the-art results on several FL tasks.

2.2 Sentiment Analysis on Figurative Language

The Semantic Evaluation Workshop-2015 [24] proposed a joint task to evaluate the impact of FL in sentiment analysis on ironic, sarcastic and metaphorical tweets, with a number of submissions achieving highly performance results. The ClaC team [104] exploited four lexicons to extract attributes as well as syntactic features to identify sentiment polarity. The UPF team [3] introduced a regression classification methodology on tweet features extracted with the use of the widely utilized

SentiWordNet and DepecheMood lexicons. The LLT-PolyU team [98] used semi-supervised regression and decision trees on extracted uni-gram and bi-gram features, coupled with features that capture potential contradictions at short distances. An SVM-based classifier on extracted n-gram and Tf-idf features was used by the Elirf team [27] coupled with specific lexicons such as Affin, Patter and Jeffrey 10. Finally, the LT3 team [88] used an ensemble Regression and SVM semi-supervised classifier with lexical features extracted with the use of WordNet and DBpedia11.

3 The background: Recent advances in Natural Language Processing

Due to the limitations of annotated datasets and the high cost of data collection, unsupervised learning approaches tend to be an easier way towards training networks. Recently, *transfer learning* approaches, i.e., the transfer of already acquired knowledge to new conditions, are gaining attention in several domain adaptation problems [22]. In fact, pre-trained embeddings representations, such as GloVe, ELMo and USE, coupled with transfer learning architectures were introduced and managed to achieve state-of-the-art results on various NLP tasks [36]. In the current section we summarize those methods in order to introduce our proposed transfer learning system in Section 5. Model specifications used for the state-of-the-art models can be found in Appendix A.

3.1 Contextual Embeddings

Pre-trained word embeddings proved to increase classification performances in many NLP tasks. In particular, Global Vectors (GloVe) [69] and Word2Vec [62] became popular in various tasks due to their ability to capture representative semantic representations of words, trained on large amount of data. However, in various studies (e.g., [70,71,60]) it is argued that the actual meaning of words along with their semantics representations varies according to their context. Following this assumption, researchers in [71] present an approach that is based on the creation of pre-trained word embeddings through building a bidirectional Language model, i.e. predicting next word within a sequence. The ELMo model was exhaustingly trained on 30 million sentences corpus [11], with a two layered bidirectional LSTM architecture, aiming to predict both next and previous words, introducing the concept of contextual embeddings. The final embeddings vector is

produced by a task specific weighted sum of the two directional hidden layers of LSTM models. Another contextual approach for creating embedding vector representations is proposed in [9] where, complete sentences, instead of words, are mapped to a latent vector space. The approach provides two variations of Universal Sentence Encoder (USE) with some trade-offs in computation and accuracy. The first approach consists of a computationally intensive transformer that resembles a transformer network [90], proved to achieve higher performance figures. In contrast, the second approach provides a light-weight model that averages input embedding weights for words and bi-grams by utilizing of a Deep Average Network (DAN) [40]. The output of the DAN is passed through a feedforward neural network in order to produce the sentence embeddings. Both approaches take as input lowercased PTB tokenized¹ strings, and output a 512-dimensional sentence embedding vectors.

3.2 Transformer Methods

Sequence-to-sequence (seq2seq) methods using encoder-decoder schemes are a popular choice for several tasks such as Machine Translation, Text Summarization, Question Answering etc. [86]. However, encoders contextual representations are uncertain when dealing with long-range dependencies. To address these drawbacks, Vaswani et al. in [90] introduced a novel network architecture, called Transformer, relying entirely on self-attention units to map input sequences to output sequences without the use of RNNs. The Transformers decoder unit architecture contains a masked multi-head attention layer followed by a multi-head attention unit and a feed forward network whereas the encoder unit is almost identical without the masked attention unit. Multi-head self-attention layers are calculated in parallel facing the computational costs of regular attention layers used by previous seq2seq network architectures. In [17] the authors presented a model that is founded on findings from various previous studies (e.g., [14, 37, 71, 75, 90]), which achieved state-of-the-art results on eleven NLP tasks, called BERT - Bidirectional Encoder Representations from Transformers. The BERT training process is split in two phases, the unsupervised pre-training phase and the fine-tuning phase using labelled data for down-streaming tasks. In contrast with previous proposed models (e.g., [71, 75]), BERT uses masked language models (MLMs) to enable pre-trained deep bidirectional representations. In the pre-training phase the model is trained with a large amount of unlabeled

data from Wikipedia, BookCorpus [103] and WordPiece [97] embeddings. In this training part, the model was tested on two tasks; on the first task, the model randomly masks 15% of the input tokens aiming to capture conceptual representations of word sequences by predicting masked words inside the corpus, whereas in the second task the model is given two sentences and tries to predict whether the second sentence is the next sentence of the first. In the second phase, BERT is extended with a task-related classifier model that is trained on a supervised manner. During this supervised phase, the pre-trained BERT model receives minimal changes, with the classifiers parameters trained in order to minimize the loss function. Two models presented in [17], a Base Bert model with 12 encoder layers (i.e. transformer blocks), feed-forward networks with 768 hidden units and 12 attention heads, and a Large Bert model with 24 encoder layers 1024 feed-the pre-trained Bert model, an architecture almost identical with the aforementioned Transformer network. A [CLS] token is supplied in the input as the first token, the final hidden state of which is aggregated for classification tasks. Despite the achieved breakthroughs, the BERT model suffers from several drawbacks. Firstly, BERT, as all language models using Transformers, assumes (and pre-supposes) independence between the masked words from the input sequence, and neglects all the positional and dependency information between words. In other words, for the prediction of a masked token both word and position embeddings are masked out, even if positional information is a key-aspect of NLP [15]. In addition, the [MASK] token which, is substituted with masked words, is mostly absent in fine-tuning phase for down-streaming tasks, leading to a pre-training fine-tuning discrepancy. To address the cons of BERT, a permutation language model was introduced, so-called XLnet, trained to predict masked tokens in a non-sequential random order, factorizing likelihood in an autoregressive manner without the independence assumption and without relying on any input corruption [99]. In particular, a query stream is used that extends embedding representations to incorporate positional information about the masked words. The original representation set (content stream), including both token and positional embeddings, is then used as input to the query stream following a scheme called Two-Stream SelfAttention. To overcome the problem of slow convergence the authors propose the prediction of the last token in the permutation phase, instead of predicting the entire sequence. Finally, XLnet uses also a special token for the classification and separation of the input sequence, [CLS] and [SEP] respectively, however it also learns an embedding that denotes whether the two

¹ <https://nlp.stanford.edu/software/tokenizer.html>

words are from the same segment. This is similar to relative positional encodings introduced in TransformerXL [15], and extends the ability of XLnet to cope with tasks that encompass arbitrary input segments. Recently, a replication study, [58], suggested several modifications in the training procedure of BERT which, outperforms the original XLNet architecture on several NLP tasks. The optimized model, called Robustly Optimized BERT Approach (RoBERTa), used 10 times more data (160GB compared with the 16GB originally exploited), and is trained with far more epochs than the BERT model (500K vs 100K), using also 8-times larger batch sizes, and a byte-level BPE vocabulary instead of the character-level vocabulary that was previously utilized. Another significant modification, was the dynamic masking technique instead of the single static mask used in BERT. In addition, RoBERTa model removes the next sentence prediction objective used in BERT, following advises by several other studies that question the NSP loss term [54,100,43].

4 Proposed Method: Recurrent CNN RoBERTa (RCNN-RoBERTa)

The intuition behind our proposed RCNN-RoBERTa approach is founded on the following observation: as pre-trained networks are beneficial for several downstream tasks, their outputs could be further enhanced if processed properly by other networks. Towards this end, we devised an end-to-end model that utilizes pre-trained RoBERTa [58] weights combined with a RCNN in order to capture contextual information. The RoBERTa network architecture is utilized in order to efficiently map words onto a rich embedding space. To improve RoBERTa's performance and identify FL within a sentence, it is essential to capture the dependencies within RoBERTa's pre-trained word-embeddings. This task can be tackled with an RNN layer suited to capture temporal reliant information, in contrast, to fully-connected and 1D convolution layers that are not able to delineate with such dependencies. In addition, aiming to enhance the proposed network architecture, the RNN layer is followed with a fully connected layer that simulates 1D convolution with a large kernel (see below), which is capable to capture spatio-temporal dependencies in RoBERTa's projected latent space. Actually, the proposed leaning model is based on a hybrid DL neural architecture that utilizes pre-trained transformer models and feed the hidden representations of the transformer into a Recurrent Convolutional Neural Network (RCNN), similar to [53]. In particular, we employed the RoBERTa base model with 12 hidden states and 12 attention heads, and used its output hidden states as an

Table 1 Selected hyperparameters used in our proposed method RCNN-RoBERTa. The hyperparameters were settled following a grid search based on a 5-fold cross-validation process; the finally selected parameters are the ones that exhibit the best performance.

<i>Hyperparameter</i>	Value
RoBERTa Layers	12
RoBERTa Attention Heads	12
LSTM units	64
LSTM dropout	0.1
Batch size	10
Adam epsilon	1e-6
Epochs	5
Learning rate	2e-5
Weight decay	1e-5

embedding layer to a RCNN. As already stated, contradictions and long-time dependencies within a sentence may be used as strong identifiers of FL expressions. RNNs are often used to capture temporal relationships between words. However they are strongly biased, i.e. later words are tending to be more dominant than previous ones. This problem can be alleviated with CNNs, which, as unbiased models, can determine semantic relationships between words with max-pooling [53,64]. Nevertheless, contextual information in CNNs is depended totally on kernel sizes. Thus, we appropriately modified the RCNN model presented in [53] in order to capture unbiased recurrent informative relationships within text. In particular, we implemented a Bidirectional LSTM (BiLSTM) layer, which is fed with RoBERTa's final hidden layer weights. The output of LSTM is concatenated with the embedded weights, and passed through a feedforward network, acting as a 1D convolution layer with large kernel, and a max-pooling layer. Finally, softmax function is used for the output layer. Table 1 shows the parameters used in training and Figure 1 illustrates the proposed deep network architecture.

5 Experimental Results

To assess the performance of the proposed method we performed an exhaustive comparison with several advanced state-of-the-art methodologies along with published results. Nowadays trends in NLP community tend to explicitly utilize deep learning methodologies as the most convenient way to approach various semantic analysis tasks. In the past decade, RNNs such as LSTM and GRUs were the most popular choice, whereas the last years the impact of attention-based models such as Transformers seems to outperform all previous methods, even by a large margin [90,17]. On the contrary, classical machine learning algorithms such as SVM, k-

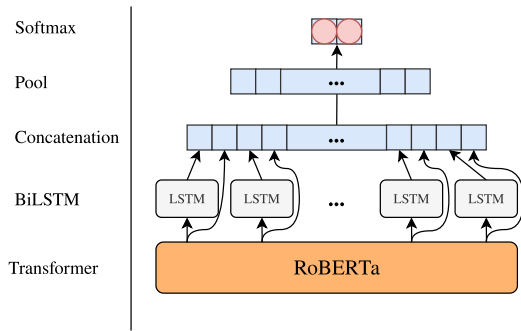


Fig. 1 The proposed RCNN-RoBERTa methodology, consisting of a RoBERTa pre-trained transformer followed by a Bidirectional LSTM layer (BiLSTM). Pooling is applied to the representation vector of concatenated RoBERTa and LSTM outputs and passed through a fully connected softmax-activated layer. We refer the reader to [58,90] for RoBERTa Transformer-based architecture.

Nearest Neighbors (kNN) and tree-based models (Decision Trees, Random Forest) have been considered inappropriate for real world applications, due to their demand on hand-crafted feature extraction and exhaustive preprocessing strategies. In order to have a reasonable kNN or SVM algorithm, there should be a lot of effort to embed sentences on word level to a higher space that a classifier may recognize patterns. In support of the arguments made, in our previous study [74], classical machine learning algorithms supported with rich and informative features failed to compete deep learning methodologies and proved non-feasible to FL detection. To this end, in this study we acquired several state-of-the-art models to compare our proposed method. The used methodologies were appropriately implemented using the available codes and guidelines, and include: ELMo [71], USE [9], NBSVM [93], FastText [44], XLnet base cased model (XLnet) [99], BERT [17] in two setups: BERT base cased (BERT-Cased) and BERT base uncased (BERT-Uncased) models, and RoBERTa base model [58]. The settings and the hyperparameters used for training the aforementioned models can be found in Appendix A. The published results were acquired from the respective original publication (the reference publication is indicated in the respective tables). For the comparison we utilized benchmark datasets that include ironic, sarcastic and metaphoric expressions. Namely, we used the dataset provided in Semantic Evaluation Workshop Task 3 (SemEval-2018) that contains ironic tweets [34]; Riloffs high quality sarcastic unbalanced dataset [80]; a large dataset containing political comments from Reddit [47]; and a SA dataset that contains tweets with various FL forms from SemEval-2015 Task 11 [24]. All datasets are used in a binary classification manner (i.e., irony/sarcasm vs. lit-

Table 2 Comparison of RCNN-RoBERTa with state-of-the-art neural network classifiers and published results on SemEval-2018 dataset; bold figures indicate superior performance.

Irony/SemVal-2018-Task 3.A [34]					
System	Acc	Pre	Rec	F1	AUC
ELMo	0.66	0.66	0.67	0.66	0.72
USE	0.69	0.67	0.67	0.67	0.74
NBSVM	0.69	0.70	0.69	0.69	0.73
FastText	0.69	0.71	0.69	0.69	0.73
XLnet	0.71	0.71	0.72	0.70	0.80
BERT-Cased	0.70	0.69	0.70	0.69	0.77
BERT-Uncased	0.69	0.68	0.69	0.68	0.77
RoBERTa	0.79	0.78	0.79	0.78	0.89
Wu et al. [96]	0.74	0.63	0.80	0.71	-
Ili et al. [39]	0.71	0.70	0.70	0.70	-
THU_NGN [96]	0.73	0.63	0.80	0.71	-
NTUA-SLP [5]	0.73	0.65	0.69	0.67	-
Zhang et al. [101]	-	-	-	0.71	-
DESC [74]	0.74	0.73	0.73	0.73	0.78
Proposed	0.82	0.81	0.80	0.80	0.89

eral), except from the SemEval-2015 Task 11 dataset where the task is to predict a sentiment integer score (from -5 to 5) for each tweet (refer to [74] for more details). For a fair comparison, we splitted the datasets on train/test stets as proposed by the authors providing the datasets or by following the settings of the respective published studies. The evaluation was made across standard five metrics namely, Accuracy (Acc), Precision (Pre), Recall (Rec), F1-score (F1), and Area Under the Receiver Operating Characteristics Curve (AUC). For the SA task the cosine similarity metric (Cos) and mean squared error (MSE) metrics are used, as proposed in the original study [24].

The results are summarized in the tables 2-5; each table refers to the respective comparison study. All tables present the performance results of our proposed method (Proposed) and contrast them to eight state-of-the-art baseline methodologies along with published results using the same dataset. Specifically, Table 2 presents the results obtained using the ironic dataset used in SemEval-2018 Task 3.A, compared with recently published studies and two high performing teams from the respective SemEval shared task [5,96]. Tables 3,4 summarize results obtained using Sarcastic datasets (Reddit SARC politics [47] and Riloff Twitter [80]). Finally, Table 5 compares the results from baseline models, from top two ranked task participants [3,104], from our previous study with the DESC methodology [74] with the proposed RCNN-RoBERTa framework on a Sentiment Analysis task with figurative language, using the SemEval 2015 Task 11 dataset.

As it can be easily observed, the proposed RCNN-RoBERTa approach outperforms all approaches as well

Table 3 Comparison of RCNN-RoBERTa with state-of-the-art neural network classifiers and published results on Reddit Politics dataset.

Reddit SARC2.0 politics [47]					
System	Acc	Pre	Rec	F1	AUC
ELMo	0.70	0.70	0.70	0.70	0.77
USE	0.75	0.75	0.75	0.75	0.82
NBSVM	0.65	0.65	0.65	0.65	0.68
FastText	0.63	0.65	0.61	0.63	0.64
XLnet	0.76	0.77	0.74	0.76	0.83
BERT-Cased	0.76	0.76	0.75	0.76	0.84
BERT-Uncased	0.77	0.77	0.77	0.77	0.84
RoBERTa	0.77	0.77	0.77	0.77	0.85
CASCADE [33]	0.74	-	-	0.75	-
Ili et al. [39]	0.79	-	-	-	-
Khodak et al. [47]	0.77	-	-	-	-
Proposed	0.79	0.78	0.78	0.78	0.85

Table 4 Comparison of RCNN-RoBERTa with state-of-the-art neural network classifiers and published results on on Sarcastic Rillofs dataset.

Riloff Sarcastic Dataset[80]					
System	Acc	Pre	Rec	F1	AUC
ELMo	0.85	0.85	0.86	0.85	0.89
USE	0.87	0.81	0.76	0.78	0.89
NBSVM	0.75	0.59	0.57	0.58	0.60
FastText	0.83	0.83	0.61	0.64	0.85
XLnet	0.86	0.88	0.86	0.86	0.92
BERT-Cased	0.86	0.87	0.85	0.86	0.91
BERT-Uncased	0.87	0.88	0.87	0.87	0.91
RoBERTa	0.89	0.85	0.84	0.85	0.91
Farrias et al. [20]	-	-	-	0.75	-
Ili et al. [39]	0.86	0.78	0.77	0.75	-
Tay el at. [87]	0.82	0.74	0.73	0.73	-
DESC [74]	0.87	0.86	0.87	0.87	0.86
Ghosh [25]	-	0.88	0.88	0.88	-
Proposed	0.91	0.90	0.90	0.90	0.94

Table 5 Comparison of RCNN-RoBERTa with state-of-the-art neural network classifiers and published results on Task11 - SemEval-2015 dataset (sentiment analysis of figurative language expression).

SemEval-2015 Task 11 [24]		
System	COS	MSE
ELMo	0.710	3.610
USE	0.71	3.17
NBSVM	0.69	3.23
FastText	0.72	2.99
XLnet	0.76	1.84
BERT-Cased	0.72	1.97
BERT-Uncased	0.79	1.54
RoBERTa	0.78	1.55
UPF [3]	0.71	2.46
Clac [104]	0.76	2.12
DESC [74]	0.82	2.48
Proposed	0.81	1.45

as all methods with published results, for the respective binary classification tasks (Tables 2, 3, and 4). In particular, the RCNN architecture seems to reinforce RoBERTa model by 2-5% F1 score, increasing also the classification confidence, in terms of AUC performance. Note also that RoBERTa-RCNN show better behaviour, compared to RoBERTa, on imbalanced datasets (Riloff [80], SemEval-2015 [24]). Also, one-way ANOVA Tukey test [63] revealed that RoBERTa-RCNN model outperforms by a statistical significant margin the maximum values of all metrics of previously published approaches, i.e. $p = 0.015; p < 0.05$ for Irony tweets and $p = 0.003; p < 0.01$ for Riloff Sarcastic tweets. Furthermore, the proposed method increased the state-of-the-art performance even by a large margin in terms of Accuracy, F1 and AUC score. Our previous approach, DESC (introduced in [74]), performs slightly better in terms of cosine similarity for the sentiment scoring task (Table 5, 0,820 vs. 0,810), with the RCNN-RoBERTa approach to perform better and managing to significantly improve the MSE measure by almost 33.5% (2,480 vs. 1,450).

6 Conclusion

In this study, we propose the first transformer based methodology, leveraging the pre-trained RoBERTa model combined with a recurrent convolutional neural network, to tackle figurative language in social media. Our network is compared with all, to the best of our knowledge, published approaches under four different benchmark dataset. In addition, we aim to minimize preprocessing and engineered feature extraction steps which are, as we claim, unnecessary when using overly trained deep learning methods such as transformers. In fact, hand crafted features along with preprocessing techniques such as stemming and tagging on huge datasets containing thousands of samples are almost prohibited in terms of their computation cost. Our proposed model, RCNN-RoBERTa, achieve state-of-the-art performance under six metrics over four benchmark dataset, denoting that transfer learning non-literal forms of language. Moreover, RCNN-RoBERTa model outperforms all other state-of-the-art approaches tested including BERT, XLnet, ELMo, and USE under all metric, some by a large factor.

References

1. Amir, S., Wallace, B.C., Lyu, H., Silva, P.C.M.J.: Modelling context with user embeddings for sarcasm detection in social media. arXiv preprint arXiv:1607.00976 (2016)

2. Antonakaki, D., Spiliotopoulos, D., V. Samaras, C., Pratikakis, P., Ioannidis, S., Fragopoulou, P.: Social media analysis during political turbulence. *PLOS ONE* **12**(10) (2017)
3. Barbieri, F., Ronzano, F., Saggion, H.: UPF-taln: SemEval 2015 Tasks 10 and 11. Sentiment Analysis of Literal and Figurative Language in Twitter. In: Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), pp. 704–708. Association for Computational Linguistics, Denver, Colorado (2015)
4. Barbieri, F., Saggion, H.: Modelling Irony in Twitter. In: EACL (2014)
5. Baziotis, C., Nikolaos, A., Papalampidi, P., Kolovou, A., Paraskevopoulos, G., Ellinas, N., Potamianos, A.: NTUA-SLP at SemEval-2018 Task 3: Tracking Ironic Tweets using Ensembles of Word and Character Level Attentive RNNs. In: Proceedings of The 12th International Workshop on Semantic Evaluation, pp. 613–621. Association for Computational Linguistics, New Orleans, Louisiana (2018)
6. Benedek, M., Beaty, R., Jauk, E., Koschutnig, K., Fink, A., Silvia, P.J., Dunst, B., Neubauer, A.C.: Creating metaphors: The neural basis of figurative language production. *NeuroImage* **90**, 99–106 (2014)
7. Buschmeier, K., Cimiano, P., Klinger, R.: An impact analysis of features in a classification approach to irony detection in product reviews. In: Proceedings of the 5th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, pp. 42–49. Association for Computational Linguistics, Baltimore, Maryland (2014)
8. Carvalho, P.: Clues for detecting irony in user-generated contents: Oh...!! its so easy. In: International CIKM Workshop on Topic-Sentiment Analysis for Mass Opinion Measurement, Hong Kong (2009)
9. Cer, D., Yang, Y., Kong, S.y., Hua, N., Limtiaco, N., John, R.S., Constant, N., Guajardo-Cespedes, M., Yuan, S., Tar, C., et al.: Universal sentence encoder. arXiv preprint arXiv:1803.11175 (2018)
10. Charalampakis, B., Spathis, D., Kouslis, E., Kermanidis, K.: A comparison between semi-supervised and supervised text mining techniques on detecting irony in greek political tweets. *Engineering Applications of Artificial Intelligence* **51**, 50–57 (2016)
11. Chelba, C., Mikolov, T., Schuster, M., Ge, Q., Brants, T., Koehn, P., Robinson, T.: One billion word benchmark for measuring progress in statistical language modeling. arXiv preprint arXiv:1312.3005 (2013)
12. Clark, H.H., Gerrig, R.J.: On the pretense theory of irony. (1984)
13. Cuccio, V., Ambrosecchia, M., Ferri, F., Carapezza, M., Piparo, F.L., Fogassi, L., Gallese, V.: How the context matters. Literal and figurative meaning in the embodied language paradigm. *PLoS ONE* **9**(12) (2014)
14. Dai, A.M., Le, Q.V.: Semi-supervised sequence learning. In: Advances in neural information processing systems, pp. 3079–3087 (2015)
15. Dai, Z., Yang, Z., Yang, Y., Cohen, W.W., Carbonell, J., Le, Q.V., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context. arXiv preprint arXiv:1901.02860 (2019)
16. Davidov, D., Tsur, O., Rappoport, A.: Semi-supervised Recognition of Sarcastic Sentences in Twitter and Amazon. In: Proceedings of the Fourteenth Conference on Computational Natural Language Learning, CoNLL ’10, pp. 107–116. Association for Computational Linguistics, Stroudsburg, PA, USA (2010)
17. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (2019)
18. Dridi, A., Recupero, D.R.: Leveraging semantics for sentiment polarity detection in social media. *International Journal of Machine Learning and Cybernetics* **10**(8), 2045–2055 (2019)
19. Dubey, A., Kumar, L., Somani, A., Joshi, A., Bhattacharyya, P.: when numbers matter!: Detecting sarcasm in numerical portions of text. In: Proceedings of the Tenth Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, pp. 72–80 (2019)
20. Farias, D.I.H., Montes-y Gómez, M., Escalante, H.J., Rosso, P., Patti, V.: A knowledge-based weighted knn for detecting irony in twitter. In: Mexican International Conference on Artificial Intelligence, pp. 194–206. Springer (2018)
21. Farias, D.I.H., Patti, V., Rosso, P.: Irony detection in twitter: The role of affective content. *ACM Transactions on Internet Technology (TOIT)* **16**(3), 19 (2016)
22. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., Lempitsky, V.: Domain-adversarial training of neural networks. *The Journal of Machine Learning Research* **17**(1), 2096–2030 (2016)
23. Gao, Z., Gao, S., Xu, L., Zheng, X., Ma, X., Luo, L., Kendrick, K.M.: Women prefer men who use metaphorical language when paying compliments in a romantic context. *Scientific Reports* **7**, 40871 (2017)
24. Ghosh, A., Li, G., Veale, T., Rosso, P., Shutova, E., Barnden, J., Reyes, A.: SemEval-2015 task 11: Sentiment analysis of figurative language in twitter. In: Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), pp. 470–478. Association for Computational Linguistics, Denver, Colorado (2015)
25. Ghosh, A., Veale, T.: Fracking sarcasm using neural network. In: Proceedings of the 7th workshop on computational approaches to subjectivity, sentiment and social media analysis, pp. 161–169 (2016)
26. Ghosh, D., Guo, W., Muresan, S.: Sarcastic or Not: Word Embeddings to Predict the Literal or Sarcastic Meaning of Words. In: EMNLP (2015)
27. Gimnez, M., Pla, F., Hurtado, L.F.: ELiRF: A SVM Approach for SA tasks in Twitter at SemEval-2015. In: Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), pp. 574–581. Association for Computational Linguistics, Denver, Colorado (2015)
28. González-Ibez, R.I., Muresan, S., Wacholder, N.: Identifying Sarcasm in Twitter: A Closer Look. In: ACL (2011)
29. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press (2016)
30. Grice, H.P.: Further Notes on Logic and Conversation. In: J.E. Adler, L.J. Rips (eds.) Reasoning: Studies of Human Inference and its Foundations, pp. 765–773. Cambridge University Press, Cambridge (2008)
31. Gupta, U., Chatterjee, A., Srikanth, R., Agrawal, P.: A Sentiment-and-Semantics-Based Approach for Emotion Detection in Textual Conversations. (2017)

32. Hangya, V., Farkas, R.: A comparative empirical study on social media sentiment analysis over various genres and languages. *Artificial Intelligence Review* **47**(4), 485–505 (2017)
33. Hazarika, D., Poria, S., Gorantla, S., Cambria, E., Zimmermann, R., Mihalcea, R.: Cascade: Contextual sarcasm detection in online discussion forums. *arXiv preprint arXiv:1805.06413* (2018)
34. Hee, C.V., Lefever, E., Hoste, V.: SemEval-2018 Task 3: Irony Detection in English Tweets. In: *SemEval@NAACL-HLT* (2018)
35. Hiai, S., Shimada, K.: Sarcasm detection using features based on indicator and roles. In: *International Conference on Soft Computing and Data Mining*, pp. 418–428. Springer (2018)
36. Howard, J., Ruder, S.: Universal language model fine-tuning for text classification. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 328–339. Association for Computational Linguistics, Melbourne, Australia (2018)
37. Howard, J., Ruder, S.: Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146* (2018)
38. Huang, Y.H., Huang, H.H., Chen, H.H.: Irony Detection with Attentive Recurrent Neural Networks. In: *ECIR* (2017)
39. Ilić, S., Marrese-Taylor, E., Balazs, J.A., Matsuo, Y.: Deep contextualized word representations for detecting sarcasm and irony. *arXiv preprint arXiv:1809.09795* (2018)
40. Iyyer, M., Manjunatha, V., Boyd-Graber, J., Daumé III, H.: 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 (Volume 1: Long Papers)*, pp. 1681–1691. Association for Computational Linguistics, Beijing, China (2015)
41. Jianqiang, Z., Xiaolin, G., Xuejun, Z.: Deep Convolution Neural Networks for Twitter Sentiment Analysis. *IEEE Access* (2018)
42. Joseph, J.K., Dev, K.A., Pradeepkumar, A.P., Mohan, M.: Chapter 16 - Big Data Analytics and Social Media in Disaster Management. In: P. Samui, D. Kim, C.B.T.I.D.S. Ghosh, Management (eds.) *Integrating Disaster Science and Management*, pp. 287–294. Elsevier (2018)
43. Joshi, M., Chen, D., Liu, Y., Weld, D.S., Zettlemoyer, L., Levy, O.: Spanbert: Improving pre-training by representing and predicting spans. *arXiv preprint arXiv:1907.10529* (2019)
44. Joulin, A., Grave, E., Bojanowski, P., Douze, M., Jégou, H., Mikolov, T.: Fasttext. zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651* (2016)
45. Kasparian, K.: Hemispheric differences in figurative language processing: Contributions of neuroimaging methods and challenges in reconciling current empirical findings (2013)
46. Katz, J.J.: Propositional structure and illocutionary force : a study of the contribution of sentence meaning to speech acts / Jerrold J. Katz. *The Language and thought series*. Crowell, New York (1977)
47. Khodak, M., Saunshi, N., Vodrahalli, K.: A Large Self-Annotated Corpus for Sarcasm. *ArXiv e-prints* (2017)
48. Kim, E., Klinger, R.: A Survey on Sentiment and Emotion Analysis for Computational Literary Studies. (2018)
49. Kingma, D.P., Ba, J.: Adam: A Method for Stochastic Optimization. *ArXiv e-prints* (2014)
50. Kumar, A., Garg, G.: Empirical study of shallow and deep learning models for sarcasm detection using context in benchmark datasets. *Journal of Ambient Intelligence and Humanized Computing* pp. 1–16 (2019)
51. Kumar, A., Sangwan, S.R., Arora, A., Nayyar, A., Abdel-Basset, M., et al.: Sarcasm detection using soft attention-based bidirectional long short-term memory model with convolution network. *IEEE Access* **7**, 23319–23328 (2019)
52. Kumar, L., Somani, A., Bhattacharyya, P.: Having 2 hours to write a paper is fun!: Detecting Sarcasm in Numerical Portions of Text. *ArXiv e-prints* (2017)
53. Lai, S., Xu, L., Liu, K., Zhao, J.: Recurrent convolutional neural networks for text classification. In: *Twenty-ninth AAAI conference on artificial intelligence* (2015)
54. Lample, G., Conneau, A.: Cross-lingual language model pretraining. *arXiv preprint arXiv:1901.07291* (2019)
55. Lazer, D., Pentland, A., Adamic, L., Aral, S., Barabasi, A.L., Brewer, D., Christakis, N., Contractor, N., Fowler, J., Gutmann, M., Jebara, T., King, G., Macy, M., Roy, D., Van Alstyne, M.: Life in the network: the coming age of computational social science. *Science (New York, N.Y.)* **323**(5915), 721–723 (2009)
56. Ling, J., Klinger, R.: An empirical, quantitative analysis of the differences between sarcasm and irony. In: *European Semantic Web Conference*, pp. 203–216. Springer (2016)
57. Liu, B.: *Sentiment Analysis - Mining Opinions, Sentiments, and Emotions*. Cambridge University Press (2015)
58. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019)
59. Loenneker-Rodman, B., Narayanan, S.: Computational approaches to figurative language. *Cambridge Encyclopedia of Psycholinguistics*. (2010)
60. McCann, B., Bradbury, J., Xiong, C., Socher, R.: Learned in translation: Contextualized word vectors. In: *Advances in Neural Information Processing Systems*, pp. 6294–6305 (2017)
61. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient Estimation of Word Representations in Vector Space. *ArXiv e-prints* (2013)
62. Mikolov, T., Sutskever, I., Chen, K., Corrado, G., Dean, J.: Distributed Representations of Words and Phrases and their Compositionality. *ArXiv e-prints* (2013)
63. Montgomery, D.C.: *Design and Analysis of Experiments*, ninth edition edn. John Wiley & Sons (2017)
64. Nguyen, T.H., Grishman, R.: Relation extraction: Perspective from convolutional neural networks. In: *Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing*, pp. 39–48. Association for Computational Linguistics, Denver, Colorado (2015)
65. Nozza, D., Fersini, E., Messina, E.: Unsupervised irony detection: A probabilistic model with word embeddings. In: *KDIR*, pp. 68–76 (2016)
66. Oboler, A., Welsh, K., Cruz, L.: The danger of big data: Social media as computational social science. *First Monday* **17**(7) (2012)
67. Ortega-Bueno, R., Rangel, F., Hernández Farias, D., Rosso, P., Montes-y Gómez, M., Medina Pagola, J.E.:

- 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-WS. org (2019)
68. Pennebaker, J., Francis, M.: *Linguistic Inquiry and Word Count*. Lawrence Erlbaum Associates, Incorporated (1999)
 69. Pennington, J., Socher, R., Manning, C.D.: Glove: Global Vectors for Word Representation. In: EMNLP, vol. 14, pp. 1532–1543 (2014)
 70. Peters, M.E., Ammar, W., Bhagavatula, C., Power, R.: Semi-supervised sequence tagging with bidirectional language models. arXiv preprint arXiv:1705.00108 (2017)
 71. Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. arXiv preprint arXiv:1802.05365 (2018)
 72. Potamias, R.A., Neofytou, A., Siolas, G.: NTUA-ISLab at SemEval-2019 task 9: Mining suggestions in the wild. In: Proceedings of the 13th International Workshop on Semantic Evaluation, pp. 1224–1230. Association for Computational Linguistics, Minneapolis, Minnesota, USA (2019)
 73. Potamias, R.A., Siolas, G.: NTUA-ISLab at SemEval-2019 Task 3: Determining emotions in contextual conversations with deep learning. In: Proceedings of the 13th International Workshop on Semantic Evaluation, pp. 277–281. Association for Computational Linguistics, Minneapolis, Minnesota, USA (2019)
 74. Potamias, R.A., Siolas, G., Stafylopatis, A.: A robust deep ensemble classifier for figurative language detection. In: International Conference on Engineering Applications of Neural Networks, pp. 164–175. Springer (2019)
 75. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.: Improving language understanding by generative pre-training (2018)
 76. Rajadesingan, A., Zafarani, R., Liu, H.: Sarcasm Detection on Twitter: A Behavioral Modeling Approach. In: WSDM (2015)
 77. Ravi, K., Ravi, V.: A novel automatic satire and irony detection using ensemble feature selection and data mining. *Knowledge-Based Systems* **120**, 15–33 (2017)
 78. Reyes, A., Rosso, P., Buscaldi, D.: From humor recognition to irony detection: The figurative language of social media. *Data & Knowledge Engineering* **74**, 1–12 (2012)
 79. Reyes, A., Rosso, P., Veale, T.: A Multidimensional Approach for Detecting Irony in Twitter. *Lang. Resour. Eval.* **47**(1), 239–268 (2013)
 80. Riloff, E., Qadir, A., Surve, P., De Silva, L., Gilbert, N., Huang, R.: Sarcasm as contrast between a positive sentiment and negative situation. In: EMNLP 2013 - 2013 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference, pp. 704–714. Association for Computational Linguistics (ACL) (2013)
 81. Rosenthal, S., Ritter, A., Nakov, P., Stoyanov, V.: SemEval-2014 task 9: Sentiment analysis in twitter. In: Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014), pp. 73–80. Association for Computational Linguistics, Dublin, Ireland (2014)
 82. Singh, N.K., Tomar, D.S., Sangaiah, A.K.: Sentiment analysis: a review and comparative analysis over social media. *Journal of Ambient Intelligence and Humanized Computing* (2018)
 83. Sperber, D., Wilson, D.: Irony and the use-mention distinction. *Radical Pragmatics* (1981)
 84. Stranisci, M., Bosco, C., FARIAS, H., IRAZU, D., Patti, V.: Annotating sentiment and irony in the online italian political debate on# labuonascuola. In: Tenth International Conference on Language Resources and Evaluation LREC 2016., pp. 2892–2899. elra (2016)
 85. Sulis, E., Farías, D.I.H., Rosso, P., Patti, V., Ruffo, G.: Figurative messages and affect in twitter: Differences between# irony,# sarcasm and# not. *Knowledge-Based Systems* **108**, 132–143 (2016)
 86. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: Advances in neural information processing systems, pp. 3104–3112 (2014)
 87. Tay, Y., Luu, A.T., Hui, S.C., Su, J.: Reasoning with sarcasm by reading in-between. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1010–1020. Association for Computational Linguistics, Melbourne, Australia (2018)
 88. Van Hee, C., Lefever, E., Hoste, V.: LT3: Sentiment Analysis of Figurative Tweets: piece of cake #NotReally. In: Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), pp. 684–688. Association for Computational Linguistics, Denver, Colorado (2015)
 89. Van Hee, C., Lefever, E., Hoste, V.: Exploring the fine-grained analysis and automatic detection of irony on twitter. *Language Resources and Evaluation* **52**(3), 707–731 (2018)
 90. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems, pp. 5998–6008 (2017)
 91. W. Gibbs, R.: On the psycholinguistics of sarcasm. *Journal of Experimental Psychology: General* **115** (1986)
 92. Wallace, B.C., Choe, D.K., Charniak, E.: Sparse, contextually informed models for irony detection: Exploiting user communities, entities and sentiment. In: ACL-IJCNLP 2015 - 53rd Annual Meeting of the Association for Computational Linguistics (ACL), Proceedings of the Conference, vol. 1 (2015)
 93. Wang, S., Manning, C.D.: Baselines and bigrams: Simple, good sentiment and topic classification. In: Proceedings of the 50th annual meeting of the association for computational linguistics: Short papers-volume 2, pp. 90–94. Association for Computational Linguistics (2012)
 94. Weiland, H., Bambini, V., Schumacher, P.B.: The role of literal meaning in figurative language comprehension: evidence from masked priming ERP. *Frontiers in Human Neuroscience* **8** (2014)
 95. Winbey, J.P.: *The social fact*. The MIT Press, Cambridge, Massachusetts (2019)
 96. Wu, C., Wu, F., Wu, S., Liu, J., Yuan, Z., Huang, Y.: THU_ngn at SemEval-2018 Task 3: Tweet Irony Detection with Densely connected LSTM and Multi-task Learning. In: SemEval@NAACL-HLT (2018)
 97. Wu, Y., Schuster, M., Chen, Z., Le, Q.V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al.: Google’s neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:1609.08144 (2016)
 98. Xu, H., Santus, E., Laszlo, A., Huang, C.R.: LLT-PolyU: Identifying Sentiment Intensity in Ironic Tweets. In:

- Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), pp. 673–678. Association for Computational Linguistics, Denver, Colorado (2015)
99. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., Le, Q.V.: Xlnet: Generalized autoregressive pretraining for language understanding. arXiv preprint arXiv:1906.08237 (2019)
 100. You, Y., Li, J., Hseu, J., Song, X., Demmel, J., Hsieh, C.J.: Reducing bert pre-training time from 3 days to 76 minutes. arXiv preprint arXiv:1904.00962 (2019)
 101. Zhang, S., Zhang, X., Chan, J., Rosso, P.: Irony detection via sentiment-based transfer learning. *Information Processing & Management* **56**(5), 1633–1644 (2019)
 102. Zhou, L., Pan, S., Wang, J., Vasilakos, A.V.: Machine learning on big data: Opportunities and challenges. *Neurocomputing* **237**, 350–361 (2017)
 103. Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., Fidler, S.: Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In: Proceedings of the IEEE international conference on computer vision, pp. 19–27 (2015)
 104. zdemir, C., Bergler, S.: CLaC-SentiPipe: SemEval2015 Subtasks 10 B,E, and Task 11. In: Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), pp. 479–485. Association for Computational Linguistics, Denver, Colorado (2015)

A Appendix

In our experiments we compared our model with several seven different classifiers under different settings. For the ELMo system we used the mean-pooling of all contextualized word representations, i.e. character-based embedding representations and the output of the two layer LSTM resulting with a 1024 dimensional vector, and passed it through two deep dense ReLu activated layers with 256 and 64 units. Similarly, USE embeddings are trained with a Transformer encoder and output 512 dimensional vector for each sample, which is also passed through two deep dense ReLu activated layers with 256 and 64 units. Both ELMo and USE embeddings retrieved from TensorFlow Hub². NBSVM system was modified according to [93] and trained with a 10^{-3} learning rate for 5 epochs with Adam optimizer [49]. FastText system was implemented by utilizing pre-trained embeddings [44] passed through a global max-pooling and a 64 unit fully connected layer. System was trained with Adam optimizer with learning rate 0.1 for 3 epochs. XLnet model implemented using the base-cased model with 12 layers, 768 hidden units and 12 attention heads. Model trained with learning rate 4×10^{-5} using 10^{-5} weight decay for 3 epochs. We exploited both cased and uncased BERT-base models containing 12 layers, 768 hidden units and 12 attention heads. We trained models for 3 epochs with learning rate 2×10^{-5} using 10^{-5} weight decay. We trained RoBERTa model following the setting of BERT model. RoBERTa, XLnet and BERT models implemented using pytorch-transformers library³ and were topped with two dense fully connected layers used as the output classifier.

² <https://tfhub.dev/s?module-type=text-embedding>

³ <https://huggingface.co/transformers/>