

Article

# An Assessment of Deep Learning Models and Word Embeddings for Toxicity Detection within Online Textual Comments

Danilo Dessì <sup>1,†,\*</sup>  0000-0003-3843-3285, Diego Reforgiato Recupero <sup>2,†</sup>  0000-0001-8646-6183 and Harald Sack <sup>1,†</sup>  0000-0001-7069-9804

<sup>1</sup> FIZ Karlsruhe - Leibniz Institute for Information Infrastructure & Karlsruhe Institute of Technology, (Germany); {danilo.dessi, harald.sack}@fiz-karlsruhe.de

<sup>2</sup> Department of Mathematics and Computer Science, University of Cagliari (Italy); diego.reforgiato@unica.it

\* Correspondence: diego.reforgiato@unica.it

† These authors contributed equally to this work.

Version March 23, 2021 submitted to Electronics

**Abstract:** Today, due to the explosion of online communication, more and more people are interacting online and a lot of textual comments are being produced. However, a paramount inconvenience within online environments is that comments that are shared within digital platforms can hide hazards such as fake news, insults, harassment, and, more in general, comments that may hurt someone's feelings. In this scenario, the detection of this kind of toxicity has an important role to moderate online communication. Recently, deep learning technologies delivered impressive performance within Natural Language Processing applications encompassing Sentiment Analysis and emotion detection across numerous datasets. Such models do not need any pre-defined hand-picked features, but they learn sophisticated features from the input datasets by themselves. In such a domain, word embeddings have been widely used as a way of representing words in Sentiment Analysis tasks proving to be very effective. Therefore, in this paper, we investigated the use of deep learning and word embeddings to detect six different types of toxicity within online comments. In doing so, the most suitable deep learning layers and state-of-the-art word embeddings to identify toxicity are evaluated. The results suggest that Long-Short Term Memory layers in combination with mimicked word embeddings are a good choice for this task.

**Keywords:** Deep Learning; Word Embeddings; Toxicity Detection; Binary Classification

## 1. Introduction

In these years, short text information is continuously being created due to the explosion of online communication, social networks, and e-commerce platforms. Through these systems, people can interact with each others, express opinions, engage in discussions, and receive feedback about any topic. However, a paramount inconvenience within online environments is that text spread by digital platforms can hide hazards such as fake news, insults, harassment, and, more in general, comments that may hurt someone's feeling. These comments can be considered as the digital version of personal attacks (e.g., bullying behaviors) that can cause social problems (e.g., racism), and are felt as dangerous and critical by people who are struggling to prevent and avoid them. [The risk of such a phenomenon has increased with the event of social networks and more in general within online communication platforms](#)<sup>1</sup>. An attempt to deal with this issue is the introduction of crowdsourcing voting schemes

<sup>1</sup> <https://medium.com/analytics-vidhya/twitter-toxicity-detector-using-tensorflow-js-1140e5ab57ee>

28 which give the possibility to denounce inappropriate comments in online environments to the users.  
29 Among many others, Facebook for example allows its users to report a post in terms of violence or hate  
30 speech [1]. This scheme allows Facebook to identify fake accounts, offensive comments, etc. However,  
31 these methodologies are often inefficient as they fail to detect toxic comments in real time [2], becoming  
32 a requirement within social network communities. A toxic post might have been published online  
33 much earlier than the time it is reported, and during the time it is online it might cause problems and  
34 offenses to several users which might have undesired behaviors (e.g., leaving the underlying social  
35 platform). Therefore, detecting toxicity within textual comments through novel technologies has great  
36 relevance in the prevention of adverse social effects in a timely and appropriate manner within online  
37 environments [3].

38 In the last years, the use of data for extracting meaningful information to interpret opinions and  
39 sentiments of people about various topics has taken hold. Today, textual online data is parsed to  
40 predict ratings about online courses [4], sentiments associated to companies and stocks within the  
41 financial domain [5] and, recently, healthcare [6], toxicity in online platforms [7]. All these approaches  
42 fall within the Sentiment Analysis research topic, which classifies data into positive or negative classes,  
43 and includes several subtasks such as emotion detection, aspect-based polarity detection [8], etc. To  
44 detect such knowledge, supervised Machine Learning-based systems are designed and provided by  
45 the research community to support and improve online services to mine and use the information. To  
46 employ supervised Machine Learning based tools, training data is required; however, the amount of  
47 labeled data might result insufficient, thus making challenging the design of these tools.

48 This is more stressed with the spread of Neural Networks and deep learning models, which can  
49 reproduce cognitive functions and mimic skills typically performed by the human brain, but need  
50 large amount of data to be trained. With the elapse of time, the interest in these technologies as well as  
51 their use for the identification of various kinds of toxicity within textual documents are grown [1].

52 Word embeddings are one of the cornerstones to represent textual data and feed Machine Learning  
53 tools. They are representations of words mapped to vectors of real numbers. The first word embedding  
54 model (Word2Vec) utilizing Neural Networks was published in 2013 [9] by researchers at Google.  
55 Since then, word embeddings are encountered in almost every Natural Language Processing (NLP)  
56 model used in practice today. The reason for such a mass adoption is their effectiveness. By translating  
57 a word to an embedding it becomes possible to model the semantic importance of a word in a numeric  
58 form and thus perform mathematical operations on it. In 2018, researchers at Google proposed  
59 the Bidirectional Encoder Representations from Transformers (*BERT*) [10], a deeply bidirectional,  
60 unsupervised language representation able to create word embeddings that represent the semantic of  
61 words in the context they are used. On the contrary, context-free models (e.g, Word2Vec) generate a  
62 single word embedding representation for each word in the vocabulary independently from the word  
63 context.

64 Within this scenario, in this paper various deep learning models fed by word embeddings are  
65 designed and evaluated to recognize toxicity levels within textual comments. In details, four deep  
66 learning models built by using the Keras<sup>2</sup> framework are designed, and four different types of word  
67 embeddings are analysed.

68 To this aim, the current state-of-the-art toxicity dataset released during the Kaggle challenge on  
69 toxic comments<sup>3</sup> is used.

70 The reader notices that this paper analyses the performances of deep learning and classical  
71 Machine Learning approaches (using tf-idf and word embeddings) when tackling the task of toxicity  
72 detection. Basically we want to assess whether the syntactic and semantic information lying within the  
73 text can provide hints on the presence of certain toxicity classes. In some domains and tasks this is

---

<sup>2</sup> <https://keras.io/>

<sup>3</sup> <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/overview>

74 not possible: for example, for the problem of identifying empathetic VS non empathetic discussion  
75 within answers of a therapist during motivational interviews it has been initially observed that  
76 syntactic and semantic information do not provide any clue for the classification task leading to very  
77 low accuracies [11]. Thus, for a fair analysis, it is important that the dataset does not contain any  
78 unbalanceness. Machine Learning classifiers fail to cope with imbalanced training datasets as they are  
79 sensitive to the proportions of the different classes [12]. As a consequence, these algorithms tend to  
80 favor the class with the largest proportion of observations, which may lead to misleading accuracies.  
81 That is why we preprocessed the mentioned dataset to make it balanced and then applied a 10-fold  
82 cross-validation to tackle the proposed task.

83 Thus, this paper provides the following contributions:

- 84 • We analysed four deep learning models based on Dense, Convolutional Neural Network (CNN),  
85 and Long-Short Term Memory (LSTM) layers to detect various levels of toxicity within online  
86 textual comments.
- 87 • We evaluate the use of four word embedding representations based on *Word2Vec* [13,14] and  
88 Bidirectional Encoder Representations from Transformer (*BERT*) [15] algorithms for the task of  
89 toxicity detection in online textual comments.
- 90 • We provide a comparison between deep learning models against common baselines used within  
91 classification tasks of textual resources.
- 92 • We release contextual word embeddings resource trained on a dataset including toxic comments.
- 93 • We also release mimicked word embeddings of tokens that are missing in the pre-trained Google  
94 *Word2Vec*<sup>4</sup> word embeddings.

95 The source code used for this study is freely available through a GitHub repository<sup>5</sup>.

96 The remainder of this paper is organized as follows. Section 2 includes a literature review and  
97 discusses current methods for toxicity detection in textual resources. Section 3 formalizes the problem.  
98 Section 4 describes the word embeddings and deep learning models adopted in this research work.  
99 Research results and their discussion are reported in Section 5. Finally, Section 6 concludes the paper  
100 and illustrates future directions to further tackle the detection of toxic comments.

## 101 2. Related Work

102 A few past works have already addressed the challenge of detecting toxicity within textual  
103 comments left by users within online environments. Generally, they rely on Sentiment Analysis  
104 methods [16–21] to detect and extract the subjective information and classify emotions and sentiments  
105 to determine if a toxicity facet is present or not. For doing so, NLP, Machine Learning, Text Mining, and  
106 Computational Linguistics are the most prominent technologies that are employed [22,23]. Sentiment  
107 Analysis methods, like many others within the Machine Learning domain, can be mainly split into  
108 two categories. i.e., supervised and unsupervised. Supervised techniques require the use of labeled  
109 data (training set) to train a model that can be applied to unseen data to predict a sentiment or an  
110 emotion [24–26]. These methods often are limited by the lack of labeled data, or by the fact that there are  
111 not either good or enough examples for certain categories (e.g., in case of dataset imbalance) [27]. On  
112 the other hand, unsupervised Sentiment Analysis approaches usually rely on semantic resources like  
113 lexicons, where words are assigned to scores for reflecting words relevance for target categories to infer  
114 sentiments and emotions of the input data [28–30]. Both supervised and unsupervised approaches are  
115 largely explored in literature for Sentiment Analysis tasks, which include Sentiment Analysis polarity  
116 detection (i.e., identifying whether a certain text is either positive or negative) [31], figurative-language  
117 uncovering (understanding if the input text if figurative or objective) [23,32], aspect-based polarity  
118 detection (e.g., assigning sentiment polarity to features of a certain topic such as the screen of an

---

4 <https://code.google.com/archive/p/word2vec/>

5 <https://github.com/danilo-dessi/toxicity>

119 iPhone) [33,34], sentiment scores prediction (e.g., identifying a continuous number in [-1,1] to a certain  
120 topic or text) [4], and so on.

121 However, only recently, these methodologies have been explored for toxicity detection [35],  
122 although the need to monitor online communications to identify toxicity and make the communications  
123 safe and respectful is an old and still open issue. Hence, the gap between the current methodologies  
124 and their potential use within toxicity detection remains an open challenge. Therefore, dealing with  
125 toxicity raises new challenges and research opportunities where deep learning-based approaches for  
126 Sentiment Analysis can have a relevant role in making advancements for the identification of toxicity  
127 levels.

128 Also, Semantic Web technologies are being used within Sentiment Analysis tasks. It has been  
129 proved that they bring several benefits leading to higher accuracy [36]. For example, the use of  
130 sentiment-based technologies to detect toxicity is investigated in [37]. However, the use of word  
131 embedding representation is not taken into account. A work worth noting is [23], where authors  
132 analysed the problem of figurative language detection in social media. More in detail, they focused  
133 on the use of semantic features extracted with Framester for identifying irony and sarcasm. Semantic  
134 features have been extracted to enrich the representation of input tweets with event information using  
135 frames and word senses in addition to lexical units. One more example of an unsupervised method  
136 that exploits Semantic Web technologies is represented by Sentilo [38,39]. Given a statement expressing  
137 an opinion, Sentilo recognizes its holder, detects its related topics and subtopics, links them to relevant  
138 situations and events referred to by it, and evaluates the sentiment expressed on each topic/subtopic.  
139 Moreover, Sentilo is domain-independent and relies on a novel lexical resource, which enables a proper  
140 propagation of the sentiment scores from topics to subtopics. Its output is represented as an RDF graph  
141 and, where applicable, it resolves holders' and topics' identity on Linked Data.

142 Recently, authors in [35] discussed the problem of toxicity detection and proved that context  
143 can both amplify or mitigate the perceived toxicity of posts. Besides, they found out no evidence  
144 that context actually improves the performance of toxicity classifiers. In another work [40] authors  
145 presented an interactive tool for auditing toxicity detection models by visualizing explanations for  
146 predictions and providing alternative wordings for detected toxic speech. In particular, they displayed  
147 the attention of toxicity detection models on user input, providing suggestions on how to replace  
148 sensitive text with less toxic words.

149 Others, [41], tackled the problem of identifying disguised offensive language, such as adversarial  
150 attacks that avoid known toxic patterns and lexicons. To do that, they proposed a framework to fortify  
151 existing toxic speech detectors without a large labeled corpus of veiled toxicity. In particular, they  
152 augmented the toxic speech detector's training data with new discovered offensive examples.

153 Deep learning technologies have been leveraged by authors in [42] to tackle the problem of toxic  
154 comments detection. More in details, the authors introduced two state-of-the-art neural network  
155 architectures and demonstrate how to employ a contextual language representation model.

156 One more work that deals with a sentiment toxicity detection problem is [7], where authors adopt  
157 both pre-trained word embeddings and close-domain word embeddings previously trained on a large  
158 dataset of users' comments [43]. However, their approach is based on a Logistic Regression (LR)  
159 classifier and does not use state-of-the-art deep learning technologies. Well established methodologies  
160 (e.g., k-nearest neighbors (kNN), Naive Bayes (NB), Support Vector Machines (SVM), etc.) are today  
161 outperformed for the same tasks by CNN-based models by [44].

162 One more work for toxicity detection is proposed by authors in [45] and it lies within the context  
163 of multiplayer online games. There, social interactions are an essential feature for a growing number  
164 of players worldwide. This interaction might bring undesired and unintended behavior especially if  
165 the game is designed to be highly competitive. They defined toxicity as the use of profane language by  
166 one player to insult or humiliate another player in the same team. Given the specific domain, the use of  
167 bad words is a necessary, but not sufficient condition for toxicity as they can be used to curse without  
168 the intent to offend anyone. Authors looked at the 100 most frequently used n-grams for n=1,2,3,4

169 and manually determined which of them are toxic or not. With such training data they use a SVM  
170 to predict the odds of winning for each team to observers based on their communication, while the  
171 match is still going.

172 Another work that embraces both deep learning and word embeddings for toxicity detection  
173 is reported in [1], where FastText<sup>6</sup> pre-trained embeddings are used to feed four different deep  
174 learning models based on CNN, Long Short Term Memory (LSTM), and Gated Recurrent Unit (GRU)  
175 layers. However, the experiments show weak results probably due to the class imbalance of classes.  
176 Conversely, in this work deep learning models by using a balanced dataset are trained, considering one  
177 toxicity class at a time, and trying to better represent the input texts by using word embeddings tuned  
178 to the target domain. More precisely, in one set of experiments, domain generated word embeddings  
179 are created through mimicking techniques; this allows to face slang, misspellings, or obfuscated  
180 contents not represented within pre-trained word embedding representations [46,47]. Besides the  
181 Word2Vec embeddings, state-of-the-art word embeddings called BERT [15,48,49] are used to tune the  
182 vectors to the context where words are used.

### 183 3. Problem Formulation

The problem faced in this paper is a multi-class multi-label classification problem. We turned it into several binary single-label classification problems. More precisely, given a textual comment  $c$  and a toxicity facet  $t$ , the approach is aimed to build a deep learning model

$$\gamma : (c, t) \rightarrow l$$

184 where  $l$  is a binary label that can only assume values in  $\{0, 1\}$  and indicates if the toxicity  $t$  is present  
185 in  $c$  (i.e.,  $l$  takes the value 1) or not (i.e.,  $l$  takes the value 0). Therefore, with such an approach, an  
186 independent binary classifier for each toxicity label is trained. Given an unseen sample, each binary  
187 classifier predicts whether that underlying toxicity is present or not in the sample. The combined  
188 model then predicts all the labels for this sample for which the respective classifier predicts a positive  
189 result. Although this method of dividing the task into multiple binary tasks may resemble superficially  
190 the one-vs-all and one-vs-rest methods for multi-class classification, it is essentially different from  
191 them because a single binary classifier deals with a single label without any regard to other labels  
192 whatsoever. This means that each binary classification task we formulated does not benefit from  
193 the information of the other labels at training time. However, this mapping is straightforward and  
194 does not change the semantic of the input problem [50]. By building these models for various  $t$ , the  
195 performances of the proposed solutions are evaluated with the goal of finding which combination of  
196 the deep learning layers and word embeddings can better capture the text peculiarities for toxicity  
197 detection.

### 198 4. The Proposed Approach

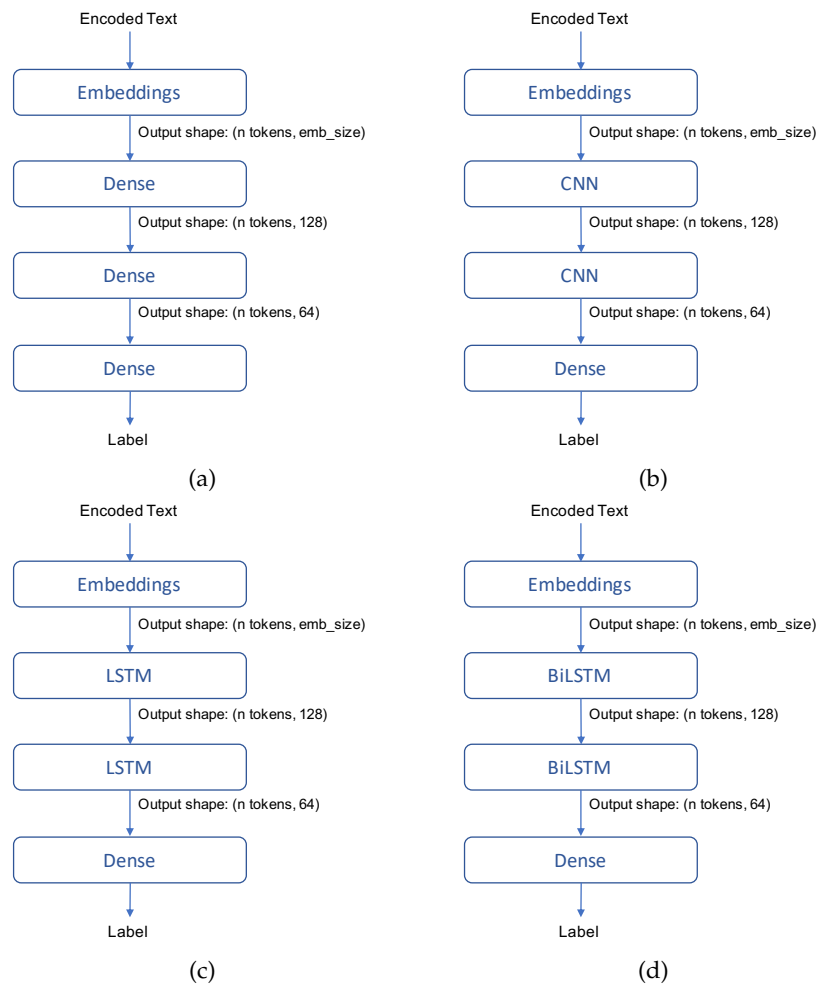
199 In this section we will describe the deep learning models and word embedding representations  
200 for representing the text expressing the various toxicity categories.

#### 201 4.1. Preprocessing

202 Text preprocessing techniques such as stop words and punctuation removal, lemmatization,  
203 stemming, matching words with a dictionary to correct grammar, removing words containing  
204 alpha-numeric characters, and so on, are common practices when Machine Learning algorithms  
205 are applied [51,52], and text representation is generated as a result of different feature engineering  
206 processes. However, with the introduction of deep learning approaches, these techniques have not

---

<sup>6</sup> <https://fasttext.cc/docs/en/english-vectors.html>



**Figure 1.** The deep learning models. (a) Dense (b) CNN (c) LSTM (d) Bidirectional LSTM. The output shape of the employed layers is indicated within the parenthesis.

207 shown promising results. The reason is that neural networks learn from any element found within the  
208 text because each token contributes to the sentence semantics. Therefore, although certain terms might  
209 be included in existing stop word lists, they are maintained because they can enrich the semantics of  
210 text content and improve the performance of the deep learning model [1]. Hence, as suggested by  
211 authors in [1], all the above mentioned preprocessing steps are ignored; only the conversion of texts in  
212 lower case is performed. Afterward, the whole set of input text is ready to feed a deep learning model.  
213 More precisely, imagine to have a toxicity target class  $t$  and a set of pairs  $P = \{(c_0, l_0), \dots, (c_n, l_n)\}$ ,  
214 where  $c_i$  is a textual comment and  $l_0$  is a binary label that can only take either the value 0 if the comment  
215  $c_i$  does not include the toxicity  $t$  or 1 if the comment  $c_i$  expresses some level of toxicity  $t$ . From the set  $P$ ,  
216 the set  $P' = \{(c'_0, l_0), \dots, (c'_n, l_n)\}$  is derived, where each comment  $c'_i$  is an integer-encoded comment  
217 of the original  $c_i$ . In details, let  $W$  be the list of all the words belonging to all textual comments, and  
218  $WS$  the set of all the words in  $W$  without duplicates (i.e.,  $WS$  has only one occurrence for each input  
219 word, whereas  $W$  can contain multiple occurrences of the same element). Then, two functions  $\theta$  and  $\phi$ ,  
220 which map the elements in  $W$  and  $WS$  to unique integer values, respectively, are built. For example,  
221 consider the sentence *you both shut up or you both die* and imagine to have the toy functions  $\theta_{toy}$  and  
222  $\phi_{toy}$ . The function  $\theta_{toy}$  maps “you” to “7”, “both” to “43”, “shut” to “22”, “up” to “76”, “or” to “10”,  
223 “you” to “3”, “both” to “41”, and “die” to “50”. The function  $\phi_{toy}$  maps “you” to “7”, “both” to “43”,  
224 “shut” to “22”, “up” to “76”, “or” to “10”, and “die” to “50”. Then the integer-encoded sentence  
225 is [7, 43, 22, 76, 10, 3, 41, 50] by applying  $\theta_{toy}$ , and [7, 43, 22, 76, 7, 43, 10, 50] by using  $\phi_{toy}$ . The reader  
226 notices that by using  $\theta_{toy}$  the words “you” and “both” are mapped to different integers. Within our  
227 approach, the function  $\theta$  is used for *BERT* word embeddings, whereas the function  $\phi$  is used to encode  
228 the input text when *Word2Vec* word embeddings are employed.

#### 229 4.2. Deep Learning Models

230 The designed deep learning model schemes are shown in Figure 1. In particular, we illustrate four  
231 deep learning models based on *Dense*, *CNN*, and *LSTM* layers available within the Keras framework<sup>7</sup>.  
232 All the models present the same number of layers. It is worth to note that the input and the output  
233 layers among the models are the same to better compare their performances considering only the type  
234 of neural network that they adopt. More precisely, the input layer is an *Embedding* layer, which has the  
235 goal of mapping the words of the input text to the underlying word embeddings. The last layer is a  
236 *Dense* layer that maps the intermediate results of the models in a single label that can only take the  
237 values 0 and 1. For doing so, it uses the *sigmoid* activation function to compute a probability that can  
238 be easily used to obtain the correct label value. In the next paragraphs we will give more details about  
239 the deep learning layers.

240 The literature already showed [53] that deep learning methods trained with word embeddings  
241 outperform those trained with tf-idf features. Therefore, we did not include the latter in our analysis  
242 as we believe that they would not add additional value to the current evaluation.

##### 243 4.2.1. Dense Model

244 The first model is depicted in Figure 1(a). It is composed of two inner *Dense* layers with 128 and  
245 64 neurons. They are densely-connected layers able to reduce the input size of hundred and thousands  
246 of nodes to a few nodes whose weights can be used to predict the final class of the input.

##### 247 4.2.2. CNN Model

248 The CNN model depicted in Figure 1(b) is based on inner CNN layers. These layers perform  
249 filtering operations to detect meaningful features of textual input for the target toxicity facet. Filters

---

<sup>7</sup> <https://keras.io/>

250 can be envisioned as kernels that slide on the vector representation and perform the same operations  
251 on each element until all the vectors have been covered. Two kernels of size 10 for the first layer,  
252 and size 5 for the second layer are used. For these layers, the same number of neurons previously  
253 introduced for the *Dense* layers is used to better compare the model performances.

#### 254 4.2.3. LSTM Model

255 The model depicted in Figure 1(c) exploits the LSTM layers to perform a binary classification of  
256 the input text. LSTMs are an extended version of Recurrent Neural Networks (RNN) and are designed  
257 to work on sequences. They use memory blocks to hold the state of the computation which makes it  
258 possible to learn temporal dependencies of data, binding the chunks of data that are currently being  
259 processed with the chunks of data already processed. This allows to infer semantic patterns that  
260 describe the history of the input data, solving the problem of common RNN whose results mostly  
261 depend on the last seen data fed into the model, smoothing the relevance of data previously processed.

#### 262 4.2.4. Bidirectional LSTM

263 The last model, shown in Figure 1(d), is an evolution of the LSTM model. It uses bidirectional  
264 LSTM layers to find patterns that can be discovered by exploring the history of the input data in both  
265 forward and backward directions. The idea of this kind of network consists of presenting the training  
266 data forwards and backward to the two bidirectional LSTM hidden layers whose results are then  
267 combined by a common output layer.

### 268 4.3. Word Embeddings Representations

269 In this section, the word embedding representations used to model the syntactic and semantic  
270 properties of the words in vectors of real numbers are introduced. Within this work, the employed  
271 word embedding representations are *Word2Vec* [13,14] and BERT [15]. We chose the most common  
272 sizes for the embeddings, i.e., 300 for *Word2Vec* embeddings and 1024 for BERT word embeddings.

#### 273 4.3.1. Word2Vec

274 The *Word2Vec* [13,14] word embedding generator aims to detect the meaning and semantic  
275 relations among the words by investigating the co-occurrence of words in documents within a given  
276 corpus. The idea behind this algorithm is to model the context of words by exploiting Machine  
277 Learning and statistics and come up with a vector representation for each word within the corpus.  
278 The resulting word vector representations allow the recognition of relatedness between words. For  
279 example, the verbs *capture* and *catch*, which are syntactically different but share common meaning and  
280 present analogous co-occurring words, are associated to similar vectors. A *Word2Vec* model can be  
281 trained by using either the Continuous Bag-Of-Words (CBOW) or the Skip-gram algorithm. Within our  
282 work, the Skip-gram algorithm is adopted because from a preliminary evaluation it obtained higher  
283 performances. In details, the following *Word2Vec* word embeddings are used:

- 284 • *Pre-trained*. Pre-trained word embeddings released by Google and available online<sup>8</sup>. They are  
285 trained on the Google news dataset and contain more than 1 billion words. However, their  
286 use can be limited by words that could be misspelled (e.g., words with orthographic errors) or  
287 domain-dependent words within the input data. These words are commonly referred to as Out  
288 Of Vocabulary (OOV) words.
- 289 • *Domain-trained*. Domain-trained word embeddings are trained on [the original unbalanced dataset](#)  
290 [\(we merged the training and the test set\) provided by the Kaggle challenge](#). [The reader notices](#)  
291 [that we computed the domain-trained embeddings on the new training sets only \(at each iteration](#)

---

<sup>8</sup> <https://code.google.com/archive/p/word2vec/>



of the 10-fold cross-validation procedure) of our evaluation strategy. Training the embeddings on the domain data solves the problem of OOV words because for each word it is possible to associate a vector. However, words that are not frequent within our data might have a vector that does not fully and correctly represent words' semantics. The Skip-gram *Word2Vec* algorithm available within the *gensim*<sup>9</sup> library is used. The model is trained using 20 epochs.

- *Mimicked*. Mimicked word embeddings are embeddings of OOV words that are not present within the original model used to represent the text data, but are inferred by exploiting syntactic similarities of words that are in the originally considered vocabulary. More in details, we used the algorithm proposed by [47], which is based on an RNN and works at character level. Words within an original vector model representation are firstly encoded by sequences of characters, and characters are associated with new vector representations. Then, by using a BiLSTM network, an OOV word  $w$  is associated to a new word embedding  $e$ . To create word embeddings for the OOV words we used the default input dataset, the hyperparameters mentioned in [47] and the pre-trained *Word2Vec* Google embeddings.

### 4.3.2. BERT

The BERT word embeddings model was introduced in late 2018 by authors in [15]. It is a novel model of pre-trained language representations that allows the tuning of word vector representations to the meaning that the word has in a given context, overcoming ambiguity issues of words. One of the famous examples is usually reported with the word *bank*. Consider the two sentences "The man was accused of robbing a bank" and "The man went fishing by the bank of the river". The introduced word embedding models describe the word *bank* with the same word embedding, i.e., they express all the possible meanings with the same vector, and, therefore, cannot disambiguate the word senses based on the surrounding context. On the other hand, BERT produces two different word embeddings, coming up with more accurate representations for the two different meanings. For doing so, BERT computes context-tuned word embeddings resulting in more accurate representations which might lead to better model performances. In this work, the *bert\_24\_1024\_16* BERT model trained on *book\_corpus\_wiki\_en\_cased* is employed and fine-tuned by using the *bert\_embedding*<sup>10</sup> library.

### 4.3.3. Word Embeddings Preparation

To load word embeddings into a deep learning model, they have to be organized into a matrix  $M$ . For *Word2Vec* word embeddings, the set  $WS$  of words in the input data is used to build  $M$  as a matrix of size  $(|WS|, 300)$ , where each row with index  $\phi(w) \mid w \in WS$  (i.e.,  $row_{\phi(w)}$ ) contains the word embedding of the word  $w$ . If a word  $w$  is not present in the *Word2Vec* selected resource (e.g., when only pre-trained word embeddings are used), then  $row_{\phi(w)}$  is a row with all its entries set to 0. Similarly, when the *BERT* embeddings are employed, the matrix  $M$  size is  $(|W|, 1024)$ , where each row with index  $\theta(w) \mid w \in W$  (i.e.,  $row_{\theta(w)}$ ) contains the word embedding of the word  $w$ . The generated matrix  $M$  is loaded into the *Embedding* layer of the employed deep learning model to map the encoded textual comments to the correct word embeddings.

## 5. Experimental study

In this section we describe the dataset used to perform our experiments, the obtained results, and the related discussion. All the experiments are run by using a 10-fold cross-validation setup. Each model is trained with batches of size 128. The model is configured to train at most with 20 epochs. However, an early stopping method with patience of 5 epochs and a delta of 0.05 that monitors the accuracy of the model are embedded within the training stage. The loss function used to train the

---

<sup>9</sup> <https://radimrehurek.com/gensim/>

<sup>10</sup> <https://pypi.org/project/bert-embedding/>

**Table 1.** Number of textual comments for each class.

Toxicity class	Number of comments	Percentage	Balanced dataset size
<i>No toxic</i>	201,081	89.95%	-
<i>toxic</i>	21,384	9.57%	42,768
<i>severe toxic</i>	1,962	0.88%	3,924
<i>obscene</i>	12,140	5.43%	24,280
<i>threat</i>	689	0.31%	1,378
<i>insult</i>	11,304	5.06%	22,608
<i>identity hate</i>	2,117	0.95%	4,234

models is the *binary crossentropy* and the used optimizer is *rmsprop* with the default learning rate 0.001 provided by the used library. The domain-trained word embeddings have been computed on the training sets only at each iteration of the 10-fold cross-validation procedure. All the other parameters have been empirically set on the basis of the models performance and previous experiences in past works [4,46]. The experiments have been carried out on a Titan X GPU mounted on a server with 16 GB of RAM memory.

### 5.1. The Dataset

To perform our analysis we employed the dataset released by a Kaggle competition<sup>11</sup>. The dataset is collected from Wikipedia comments that have been manually labeled into 6 different toxicity classes. It consists of training and test files. However, the original split is not kept in order to apply the proposed approach and balance the data. The dataset is composed of more than 200k comments and presents annotations for six different toxicity classes and one more class when no toxicity is present. Table 1 reports the number of comments and the related percentage concerning the original dataset (second and third columns) belonging to each of the seven resulting classes. The first row includes the comments that do not present toxicity, then from the second row on, the number of comments for each toxicity class (*toxic*, *severetoxic*, *obscene*, *threat*, *insult*, *identityhate*) are reported. Besides, from Table 1 it is worth to note that the dataset is strongly unbalanced as nearly 90% of the overall comments do not present toxicity. Therefore, as mentioned early in the paper, the training of a model is biased because the model does not have a sufficient number of examples of the minority class to correctly identify a pattern. A random model that always predicts the majority class can obtain better performances although it is not be able to recognize elements that should belong to the minority class. Hence, having a balanced dataset is a common procedure in several classification tasks [54] and allows understanding better the performances of a model [12]. It follows that for each toxicity class we built a dataset where the number of positive examples (i.e., comments that present the target toxicity class) and the number of negative examples (i.e., comments that do not present that target toxicity class) are the same. The size of the created datasets for each class are reported in Table 1 under the *Balanced dataset size* column. The reader notices that, for a certain toxicity class, the negative examples are chosen among all the other classes including the *No toxic* comments.

### 5.2. Baselines

For evaluation purposes, the deep learning models have been compared to a certain number of baselines. These are classical Machine Learning classifiers that are usually employed with the *tf-idf* to

<sup>11</sup> <https://www.kaggle.com/>

366 represent textual resources [51]. More precisely, the deep learning models are compared against the  
 367 following classifiers:

- 368 • **Decision Tree (DT)**. The Decision Tree algorithm builds a model by learning decision rules that  
 369 when applied to the input features can correctly predict the target class. The model has a root  
 370 node that represents the whole set of input data. This node is subsequently split into its children  
 371 by applying a given rule. The process is then applied to its children recursively as long as there  
 372 are nodes that can be split.
- 373 • **Random Forest (RF)**. This method adopts more DTs applied on different samples of the input  
 374 data and uses a majority voting strategy to predict the output classes. The strength of this  
 375 algorithm is that each DT is individually trained; therefore, overfitting and errors due to biases  
 376 are limited. We adopted a classifier that made use of 100 DTs estimators.
- 377 • **Multi-Layer Perceptron (MLP)**. This is a neural network that is composed of a single layer of  
 378 nodes. In our experiment, we used a layer with 100 nodes.

379 For these classical Machine Learning methods employed as baselines the adoption of just word  
 380 embeddings is not promising and this has already been shown in literature [55]. In particular, when  
 381 employing word embeddings for classical Machine Learning methods, they should be processed by  
 382 operations such as the average or the sum before being fed to a given classifier. This causes loss of  
 383 syntactic and semantic information expressed by the embeddings of each word.

384 To develop the algorithms above we employed the *scikit-learn*<sup>12</sup> library.

385 Additionally, the area under the ROC (Receiver Operating Characteristic) curve (ROC-AUC) is  
 386 also reported in Table 2 in order to understand the performance of our model with respect to the best  
 387 models proposed for the challenge’s task.

**Table 2.** ROC-AUC values of our deep learning models on each binary classification and average for each model.

Learning Model	Word Embeddings	Toxic	Severe Toxic	Obscene	Threat	Identity Hate	Insult	Average
Deep Model Dense	pre-trained	0.921	0.968	0.936	0.977	0.944	0.933	0.947
	domain-trained	0.915	0.959	0.928	0.968	0.934	0.924	0.938
	mimicked	0.922	0.969	0.938	0.981	0.941	0.931	0.947
	bert	0.898	0.964	0.904	0.945	0.924	0.906	0.924
Deep Model CNN	pre-trained	0.905	0.964	0.924	0.969	0.934	0.915	0.935
	domain-trained	0.895	0.950	0.857	0.957	0.909	0.903	0.912
	mimicked	0.906	0.961	0.923	0.974	0.935	0.914	0.936
	bert	0.881	0.952	0.894	0.909	0.892	0.895	0.904
Deep Model LSTM	pre-trained	0.970	0.982	0.980	0.983	0.968	0.976	0.977
	domain-trained	0.963	0.980	0.977	0.983	0.968	0.970	0.974
	mimicked	0.971	0.983	0.977	0.985	0.970	0.977	0.977
	bert	0.930	0.974	0.940	0.956	0.950	0.940	0.948
Deep Model Bidirectional LSTM	pre-trained	0.969	0.981	0.973	0.984	0.967	0.975	0.975
	domain-trained	0.963	0.980	0.977	0.984	0.964	0.970	0.973
	mimicked	0.969	0.963	0.980	0.988	0.970	0.976	0.974
	bert	0.930	0.970	0.939	0.951	0.947	0.941	0.946

<sup>12</sup> <https://scikit-learn.org/stable/index.html>

**Table 3.** Precision (p), recall (r), and f-measure (f) related to the binary classification for each toxicity class using the balanced dataset.

Learning Model	Word Embeddings	Toxic			Severe Toxic			Obscene		
		p	r	f	p	r	f	p	r	f
Decision Trees	tf-idf	0.859	0.855	0.857	0.847	<b>0.947</b>	0.894	0.926	<b>0.929</b>	<b>0.928</b>
Random Forests	tf-idf	<b>0.860</b>	0.856	<b>0.858</b>	0.888	0.940	0.913	<b>0.945</b>	0.834	0.913
MLP	tf-idf	0.849	<b>0.857</b>	0.853	<b>0.913</b>	0.918	<b>0.915</b>	0.884	0.895	0.889
Deep Model Dense	pre-trained	0.863	<b>0.856</b>	<b>0.858</b>	0.923	0.910	0.916	<b>0.886</b>	0.867	0.876
	domain-trained	0.855	0.848	0.851	0.893	0.910	0.899	0.874	0.863	0.867
	mimicked	<b>0.868</b>	0.844	0.855	<b>0.926</b>	0.914	<b>0.919</b>	0.880	<b>0.877</b>	<b>0.878</b>
	bert	0.828	0.817	0.822	0.912	<b>0.917</b>	0.913	0.844	0.821	0.832
Deep Model CNN	pre-trained	<b>0.848</b>	0.849	0.848	<b>0.910</b>	0.911	<b>0.909</b>	<b>0.863</b>	0.861	0.861
	domain-trained	0.846	0.841	0.842	0.903	0.875	0.888	0.858	0.849	0.853
	mimicked	0.836	<b>0.865</b>	<b>0.850</b>	0.886	<b>0.919</b>	0.901	0.856	<b>0.870</b>	<b>0.862</b>
	bert	0.801	0.812	0.805	0.899	0.911	0.904	0.819	0.832	0.825
Deep Model LSTM	pre-trained	<b>0.914</b>	0.915	0.914	0.944	0.962	<b>0.953</b>	0.927	<b>0.949</b>	<b>0.938</b>
	domain-trained	0.903	0.916	0.909	<b>0.947</b>	0.948	0.947	<b>0.929</b>	0.944	0.936
	mimicked	0.895	<b>0.938</b>	<b>0.916</b>	0.941	<b>0.966</b>	<b>0.953</b>	0.928	0.938	0.932
	bert	0.866	0.851	0.858	0.927	0.932	0.929	0.889	0.861	0.875
Deep Model Bidirectional LSTM	pre-trained	0.906	<b>0.923</b>	0.914	0.936	0.959	0.947	<b>0.963</b>	0.854	0.905
	domain-trained	0.905	0.915	0.910	<b>0.948</b>	0.962	<b>0.955</b>	0.941	0.933	<b>0.937</b>
	mimicked	<b>0.910</b>	0.921	<b>0.915</b>	0.939	<b>0.963</b>	0.951	0.929	<b>0.945</b>	<b>0.937</b>
	bert	0.875	0.841	0.856	0.933	0.941	0.937	0.892	0.852	0.871

Learning Model	Word Embeddings	Threat			Identity Hate			Insult		
		p	r	f	p	r	f	p	r	f
Decision Trees	tf-idf	0.917	0.891	0.903	0.819	<b>0.927</b>	0.869	0.887	<b>0.891</b>	<b>0.889</b>
Random Forests	tf-idf	<b>0.954</b>	0.897	<b>0.924</b>	0.847	0.911	0.877	<b>0.929</b>	0.851	0.888
MLP	tf-idf	0.914	<b>0.916</b>	0.913	<b>0.889</b>	0.897	<b>0.893</b>	0.871	0.880	0.876
Deep Model Dense	pre-trained	<b>0.934</b>	0.930	<b>0.931</b>	<b>0.897</b>	0.865	0.879	0.872	0.865	<b>0.869</b>
	domain-trained	0.913	0.918	0.914	0.858	0.877	0.866	<b>0.876</b>	0.846	0.860
	mimicked	0.933	<b>0.932</b>	<b>0.931</b>	0.881	<b>0.882</b>	<b>0.880</b>	0.873	<b>0.857</b>	0.863
	bert	0.867	0.891	0.877	0.874	0.865	0.855	0.841	0.827	0.834
Deep Model CNN	pre-trained	<b>0.932</b>	0.870	0.891	<b>0.872</b>	0.863	0.867	0.842	<b>0.862</b>	<b>0.851</b>
	domain-trained	0.898	0.899	0.898	0.823	0.868	0.842	<b>0.874</b>	0.816	0.843
	mimicked	0.927	<b>0.918</b>	<b>0.922</b>	0.860	<b>0.879</b>	<b>0.869</b>	0.847	0.849	0.847
	bert	0.842	0.872	0.849	0.824	0.842	0.832	0.831	0.821	0.826
Deep Model LSTM	pre-trained	0.932	<b>0.967</b>	0.948	0.907	0.909	0.906	0.918	0.939	0.928
	domain-trained	0.949	0.951	0.950	<b>0.913</b>	0.925	<b>0.918</b>	<b>0.919</b>	0.930	0.924
	mimicked	<b>0.953</b>	0.962	<b>0.957</b>	0.887	<b>0.946</b>	0.914	0.916	<b>0.948</b>	<b>0.931</b>
	bert	0.916	0.899	0.907	0.880	0.895	0.886	0.874	0.870	0.872
Deep Model Bidirectional LSTM	pre-trained	0.946	<b>0.961</b>	<b>0.952</b>	<b>0.905</b>	0.921	0.912	0.918	0.931	0.924
	domain-trained	<b>0.949</b>	0.949	0.949	0.904	<b>0.935</b>	<b>0.919</b>	0.918	<b>0.938</b>	<b>0.927</b>
	mimicked	0.941	0.944	0.940	0.902	0.934	0.916	<b>0.920</b>	0.935	<b>0.927</b>
	bert	0.913	0.900	0.905	0.900	0.857	0.874	0.889	0.866	0.877

### 388 5.3. Results and Discussion

389 In this section, we discuss the results of the experiments we have carried. They are reported  
 390 in Tables 2 and 3 in terms of ROC-AUC, precision, recall, and f-measure scores (for computing the  
 391 ROC-AUC, the true positive rates and false positive rates are computed accordingly to Equations (1)  
 392 and (2); precision, recall and f-measure are computed according to Equations (3), (4), and (5)). In  
 393 the equations,  $TP$  (true positives) is the number of comments with the target toxicity class correctly  
 394 guessed by the model,  $FP$  (false positives) is the number of comments erroneously associated to a  
 395 target toxicity class,  $TN$  (true negatives) is the number of comments that the classifier correctly does not  
 396 classify for a target class, and  $FN$  (false negatives) is the number of comments erroneously classified  
 397 with a class different than the target class.

$$\text{True positive rate} = \frac{TN}{TN + FP} \quad (1)$$

$$\text{False positive rate} = \frac{FP}{FP + TN} \quad (2)$$

$$\text{Precision } (p) = \frac{TP}{TP + FP} \quad (3)$$

$$\text{Recall } (r) = \frac{TP}{TP + FN} \quad (4)$$

$$F - \text{measure } (f) = 2 \cdot \frac{P \cdot R}{P + R} \quad (5)$$

398 Results depicted in Table 3 show how the deep learning models perform against the baselines  
 399 (classical Machine Learning approaches). For each deep learning model, the performance of the model  
 400 in combination with the embedding representations is illustrated as well.

### 401 5.4. Comparison with the Kaggle Challenge

402 The results indicated in Table 2 report the ROC-AUC values of our deep learning approaches for  
 403 each toxicity class and the average over all the classes. The reader notices that it is not the purpose of  
 404 this paper to compete with the other participants of the Kaggle challenge where the data have been  
 405 extracted and the evaluation has been reported using the ROC-AUC. The best three approaches of the  
 406 challenge were *Toxic Crusaders*, *neongen & Computer says no*, and *Adversarial Autoencoder*, which reported  
 407 a ROC-AUC value of 0.989, 0.988, and 0.988, respectively. The challenge task was to test any proposed  
 408 approach on a highly unbalanced dataset. In this paper we wanted to study how deep learning  
 409 methods and classical Machine Learning approaches (using tf-idf and word embeddings) perform on  
 410 the toxicity problem without any bias (unbalanceness of the data). Moreover, it has been proved that  
 411 optimizing a method for the ROC-AUC does not guarantee the optimization on the precision-recall  
 412 curve [56]. This is why we included Table 3 with precision, recall and f-measure metrics computed on  
 413 the preprocessed balanced dataset. There are several heuristics and tuning that can be done in presence  
 414 of unbalanced datasets to help achieving high values of ROC-AUC. Those could not be performed by  
 415 us since we used a balanced version of the original dataset.

#### 416 5.4.1. Baseline Comparison

417 The results indicate that *Dense*- and *CNN*-based models are not much better than the baseline  
 418 methods. Actually, in some cases, they are outperformed. For example, considering the toxicity  
 419 classes *obscene* and *insult*, it is possible to observe that the f-measure computed on the baseline  
 420 predictions is higher than the one obtained by *Dense*- and *CNN*-based models. On the other hand,  
 421 *LSTM*-based models are able to outperform the baseline methods with a minimum improvement in

422 terms of f-measure of 0.01, i.e., in percentage 1% (see *obscene* class), and a maximum of 0.058, i.e.,  
423 in percentage 5.8% (see *toxic* class). These results are similar, and sometimes still more noticeable  
424 when the *Bidirectional LSTM* layers are employed. Moreover, considering that by using the balanced  
425 dataset every classifier is able to obtain a f-measure always higher than 0.8, the improvements can be  
426 considered remarkable. The only drawback is related to the computational time needed to train the  
427 deep learning model. Nevertheless, the training time is not reported since i) it is out of the scope of  
428 this study ii) with modern GPUs it is feasible to train complex deep learning models iii) the training  
429 step must be executed only once, and iv) the computational time needed for the prediction step does  
430 not depend on the underlying model used for the training step.

#### 431 5.4.2. Dense-based Model

432 For the task of toxicity detection the *Dense*-based model never obtains the best performances. In  
433 most of the cases, the best results with this model are obtained with the *mimicked* word embeddings  
434 where for four out of six classes the achieved f-measure score is the highest. The *pre-trained* word  
435 embeddings obtain high performances too, especially for classes such as *Toxic*, *Threat* (in this case  
436 the f-measure is very close to the case when using *mimicked*), and *Insult*. The use of *domain-trained*  
437 word embeddings never meets high scores, except when the *precision* is considered for the *Insult* class.  
438 Similarly, *BERT* word embeddings performances are the worst.

#### 439 5.4.3. CNN-based Model

440 Using the *CNN*-based model the results do not improve further with respect to the *Dense*-based  
441 model. In some cases, the performances of the model are even lower. With this model, the best results  
442 are obtained by employing the *mimicked* word embeddings for the toxicity classes *Toxic*, *Obscene*, *Threat*,  
443 and *Identity Hate*. For the other toxicity classes, the best results are obtained using the *pre-trained* word  
444 embeddings. *Domain-trained* and *BERT* embeddings are not able to properly represent the domain  
445 knowledge for the *CNN* model, thus the results are poor.

#### 446 5.4.4. LSTM-based Model

447 The *LSTM* model outperforms both *Dense* and *CNN*-based models, proving its suitability to detect  
448 patterns for toxic detection. As previously mentioned, *mimicked* word embeddings are employed  
449 for the deep learning model to learn and uncover toxicity from the text comments. *Pre-trained* and  
450 *Domain-trained* word embeddings obtain good performances, and their results are not far from the  
451 model using the *mimicked* word embeddings. On the other hand, once again *BERT* is not a good  
452 representation for the *LSTM* model. Except for *BERT*, the three other word embeddings adopted with  
453 the *LSTM* model outperform the baseline methods for almost each toxicity level.

#### 454 5.4.5. BiLSTM-based Model

455 Although the higher complexity of the employed layers, the results of the *BiLSTM* (Bidirectional  
456 *LSTM*) model are similar to those obtained by the *LSTM* model. In some cases, the *BiLSTM* is able to  
457 outperform the *LSTM*, in others it is not. Moreover, it differs from the other models because its best  
458 performances for many classes are obtained using the *domain-trained* word embeddings. The *pre-trained*  
459 and *mimicked* word embeddings continued to show good ability to represent domain knowledge, and  
460 *BERT* embeddings confirm to be the last choice for the task of toxicity detection. Similarly to the *LSTM*  
461 model, except using *BERT*, the model outperforms the baselines in almost each toxicity class.

#### 462 5.4.6. Overall evaluation of the deep learning models

463 The use of deep learning for the task of toxicity detection has shown good performances in all the  
464 toxicity classes. Also, it turns out that although the small size of datasets employed for certain classes,  
465 they are able to detect patterns that allow to correctly perform the classification. More in details, the

466 results suggest that the *Dense* and *CNN* models perform well since their f-measure is always higher  
467 than 0.8 but, for the toxicity detection task, they are outperformed by the *LSTM* and *BiLSTM* models,  
468 which obtain a f-measure higher than 0.9 in most of the cases. Results are comparable among the  
469 *LSTM* and *BiLSTM* models. However, because *BiLSTM*-based models need higher computational  
470 time to be trained than *LSTM* models, the latter are slightly preferred. It is worth to mention that the  
471 current models are trained without the context that was surrounding the comments in the Wikipedia  
472 pages (where the dataset has been originally collected) and, therefore, they might lack the necessary  
473 information to predict the correct class. One more obstacle might be also due to the presence of  
474 figurative language within the comments, which might change the meaning of the sentences, thus  
475 misleading the models. For example, a frequent sentence like *I am going to kill you* pronounced after a  
476 mistake or an undesired change in the Wikipedia pages does not necessarily convey a threat or hate  
477 emotion but it may be simply a joke.

#### 478 5.4.7. Overall evaluation of word embeddings

479 From the results, it is noticeable that the *Word2Vec* algorithm is a good choice to represent textual  
480 resources to be parsed with deep learning models. Results suggest that *mimicked* word embeddings are  
481 the best choice because they enclose the knowledge of *pre-trained* word embeddings that have been  
482 built on a large dataset and do not suffer from the OOV words problem [46]. *Domain-trained* word  
483 embeddings obtain good results but, for most of the cases, they are outperformed. This may depend  
484 on the fact that the resources employed to train these embeddings are not very large and, besides, there  
485 are not a sufficient number of examples of toxicity due to the unbalanced number of toxic comments  
486 in the input dataset (i.e., more than 200k comments do not present toxicity, the reader can see Table 1).

487 Surprisingly, *BERT* embeddings perform badly for the task of toxicity detection although they  
488 are currently the state-of-the-art word embedding representations. A possible motivation behind  
489 this finding is that assigning a different embedding to the same word is somehow misleading to the  
490 training of the deep learning models. More precisely, the tuning step performed to generate the *BERT*  
491 embeddings on our data is not able to capture the context of the words due to the length of some input  
492 textual comments and to the typos and incorrect grammar often present within them, thus transferring  
493 possible erroneous information to our deep learning models. One more reason might be due to the  
494 lack of the surrounding context of the comments; it might have limited the fine-tuning of the model,  
495 therefore leading the semantics of words to be captured badly. This fact is worth to be investigated,  
496 and a close analysis to this problem is required.

## 497 6. Conclusion and Future Work

498 In this paper, we presented an assessment of various deep learning models fed by various word  
499 embedding representations to detect toxicity within textual comments. From the obtained results  
500 we can definitely state that toxicity can be identified by machine and deep learning approaches fed  
501 with syntactic and semantic information extracted from the text. We show how *LSTM*-based model is  
502 the first choice among the experimented models to detect toxicity. We also show how various word  
503 embeddings may represent the domain knowledge in a variety of ways, and an unique model for all  
504 cases might be insufficient. In particular, the results are encouraging when using mimicking techniques  
505 to deal with OOV words where there are not many examples to build significant domain-dependent  
506 word embeddings. As future works, we plan to perform a deeper assessment of deep learning models  
507 by using and combining different layers, to better detect patterns and on real scenarios where classes  
508 may be unbalanced as well. Moreover, we would like to investigate other contextualized word  
509 embedding representations such as *ELMO* [57] for the toxicity detection task. An analysis of the  
510 proposed approaches on which configuration, parameter settings and heuristic may be added to tackle  
511 the same problem but in presence of highly unbalanced datasets is definitely a research direction we  
512 would like to investigate as well. Finally, we would like to investigate the impact of using different  
513 embeddings for the same word since it might be the cause of failure of *BERT* embeddings in our

514 experiments. We also think that an ensemble strategy of the proposed approaches should result in  
515 better overall performances and are then investigating this direction as well.

516 **Acknowledgments:** We gratefully acknowledge the support of NVIDIA Corporation with the donation of the  
517 Titan X GPU used for this research.

518 **Conflicts of Interest:** The authors declare no conflict of interest.

## 519 Abbreviations

520 The following abbreviations are used in this manuscript:

521	AUC	Area under the curve
	BERT	Bidirectional Encoder Representations from Transformers
	BiLSTM	Bidirectional Long-Short Term Memory
	CBOW	Continuous Bag-Of-Words
	CNN	Convolutional Neural Network
	DT	Decision Tree
	ELMO	Embeddings from Language Models
	GPU	Graphics Processing Unit
	GRU	Gated Recurrent Unit
	kNN	k-Nearest Neighbors
522	LR	Logistic Regression
	LSTM	Long-Short Term Memory
	MLP	Multi-Layer Perceptron
	NB	Naive Bayes
	NLP	Natural Language Processing
	OOV	Out Of Vocabulary
	RF	Random Forest
	ROC	Receiver Operating Characteristic
	RNN	Recurrent Neural Network
	TF-IDF	Term Frequency–Inverse Document Frequency
	SVM	Support Vector Machine

## 523 References

- 524 1. Saeed, H.H.; Shahzad, K.; Kamiran, F. Overlapping Toxic Sentiment Classification Using Deep Neural  
525 Architectures. 2018 IEEE International Conference on Data Mining Workshops (ICDMW). IEEE, 2018, pp.  
526 1361–1366.
- 527 2. Hosseini, H.; Kannan, S.; Zhang, B.; Poovendran, R. Deceiving google’s perspective api built for detecting  
528 toxic comments. *arXiv preprint arXiv:1702.08138* **2017**.
- 529 3. Srivastava, S.; Khurana, P.; Tewari, V. Identifying aggression and toxicity in comments using capsule  
530 network. Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018),  
531 2018, pp. 98–105.
- 532 4. Dessì, D.; Dragoni, M.; Fenu, G.; Marras, M.; Recupero, D.R. Evaluating neural word embeddings created  
533 from online course reviews for sentiment analysis. Proceedings of the 34th ACM/SIGAPP Symposium on  
534 Applied Computing. ACM, 2019, pp. 2124–2127.
- 535 5. Dridi, A.; Atzeni, M.; Recupero, D.R. FineNews: fine-grained semantic sentiment analysis on financial  
536 microblogs and news. *International Journal of Machine Learning and Cybernetics* **2019**, *10*, 2199–2207.
- 537 6. Consoli, S.; Dessì, D.; Fenu, G.; Marras, M. Deep Attention-based Model for Helpfulness Prediction of  
538 Healthcare Online Reviews. Proceedings of the First Workshop on Smart Personal Health Interfaces  
539 co-located with 25th International Conference on Intelligent User Interfaces, SmartPhil@IUI 2020, Cagliari,  
540 Italy, March 17, 2020, 2020, pp. 33–49.
- 541 7. Carta, S.; Corrigan, A.; Mulas, R.; Recupero, D.R.; Saia, R. A Supervised Multi-class Multi-label Word  
542 Embeddings Approach for Toxic Comment Classification. Proceedings of the 11th International Joint  
543 Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management, Vienna,  
544 Austria, 2019, pp. 17–19.



- 545 8. Schouten, K.; Frasinca, F. Survey on aspect-level sentiment analysis. *IEEE Transactions on Knowledge and*  
546 *Data Engineering* **2015**, *28*, 813–830.
- 547 9. Mikolov, T.; Corrado, G.; Chen, K.; Dean, J. Efficient Estimation of Word Representations in Vector Space.  
548 Proceedings of the International Conference on Learning Representations (ICLR 2013), 2013, pp. 1–12.
- 549 10. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers  
550 for Language Understanding. Proceedings of the 2019 Conference of the North American Chapter of  
551 the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and  
552 Short Papers); Association for Computational Linguistics: Minneapolis, Minnesota, 2019; pp. 4171–4186.  
553 doi:10.18653/v1/N19-1423.
- 554 11. Wu, Z.; Helaoui, R.; Kumar, V.; Recupero, D.R.; Riboni, D. Towards Detecting Need for Empathetic  
555 Response in Motivational Interviewing. Companion Publication of the 2020 International Conference on  
556 Multimodal Interaction, ICMI Companion 2020, Virtual Event, The Netherlands, October, 2020; Truong,  
557 K.P.; Heylen, D.; Czerwinski, M.; Berthouze, N.; Chetouani, M.; Nakano, M., Eds. ACM, 2020, pp. 497–502.  
558 doi:10.1145/3395035.3425228.
- 559 12. Dragoni, M.; Petrucci, G. A neural word embeddings approach for multi-domain sentiment analysis. *IEEE*  
560 *Transactions on Affective Computing* **2017**, *8*, 457–470.
- 561 13. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient estimation of word representations in vector space.  
562 *arXiv preprint arXiv:1301.3781* **2013**.
- 563 14. Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G.S.; Dean, J. Distributed representations of words and phrases  
564 and their compositionality. *Advances in neural information processing systems*, 2013, pp. 3111–3119.
- 565 15. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. Bert: Pre-training of deep bidirectional transformers for  
566 language understanding. *arXiv preprint arXiv:1810.04805* **2018**.
- 567 16. Reforgiato Recupero, D.; Cambria, E. ESWC 14 challenge on Concept-Level Sentiment Analysis.  
568 *Communications in Computer and Information Science* **2014**, *475*, 3–20.
- 569 17. Recupero, D.R.; Dragoni, M.; Presutti, V. ESWC 15 Challenge on Concept-Level Sentiment Analysis.  
570 Semantic Web Evaluation Challenges - Second SemWebEval Challenge at ESWC 2015, Portorož, Slovenia,  
571 May 31 - June 4, 2015, Revised Selected Papers, 2015, pp. 211–222. doi:10.1007/978-3-319-25518-7\_18.
- 572 18. Recupero, D.R.; Consoli, S.; Gangemi, A.; Nuzzolese, A.G.; Spampinato, D. A Semantic Web Based Core  
573 Engine to Efficiently Perform Sentiment Analysis. *The Semantic Web: ESWC 2014 Satellite Events*; Presutti,  
574 V.; Blomqvist, E.; Troncy, R.; Sack, H.; Papadakis, I.; Tordai, A., Eds.; Springer International Publishing:  
575 Cham, 2014; pp. 245–248.
- 576 19. Dragoni, M.; Reforgiato Recupero, D. Challenge on fine-grained sentiment analysis within ESWC2016.  
577 *Communications in Computer and Information Science* **2016**, *641*, 79–94.
- 578 20. Reforgiato Recupero, D.; Cambria, E.; Di Rosa, E. Semantic sentiment analysis challenge at ESWC2017.  
579 *Communications in Computer and Information Science* **2017**, *769*, 109–123.
- 580 21. Kumar, V.; Recupero, D.R.; Riboni, D.; Helaoui, R. Ensembling Classical Machine Learning and Deep  
581 Learning Approaches for Morbidity Identification From Clinical Notes. *IEEE Access* **2021**, *9*, 7107–7126.  
582 doi:10.1109/ACCESS.2020.3043221.
- 583 22. Dridi, A.; Recupero, D.R. Leveraging semantics for sentiment polarity detection in social media.  
584 *International Journal of Machine Learning and Cybernetics* **2019**, *10*, 2045–2055.
- 585 23. Recupero, D.R.; Alam, M.; Buscaldi, D.; Grezka, A.; Tavazoe, F. Frame-Based Detection of Figurative  
586 Language in Tweets [Application Notes]. *IEEE Computational Intelligence Magazine* **2019**, *14*, 77–88.
- 587 24. Poria, S.; Cambria, E.; Gelbukh, A. Deep convolutional neural network textual features and multiple  
588 kernel learning for utterance-level multimodal sentiment analysis. Proceedings of the 2015 conference on  
589 empirical methods in natural language processing, 2015, pp. 2539–2544.
- 590 25. Tang, D.; Qin, B.; Liu, T. Document modeling with gated recurrent neural network for sentiment  
591 classification. Proceedings of the 2015 conference on empirical methods in natural language processing,  
592 2015, pp. 1422–1432.
- 593 26. Atzeni, M.; Recupero, D.R. Multi-domain sentiment analysis with mimicked and polarized word  
594 embeddings for human-robot interaction. *Future Gener. Comput. Syst.* **2020**, *110*, 984–999.  
595 doi:10.1016/j.future.2019.10.012.
- 596 27. Yin, H.; Gai, K. An empirical study on preprocessing high-dimensional class-imbalanced data  
597 for classification. 2015 IEEE 17th International Conference on High Performance Computing and

- 598 Communications, 2015 IEEE 7th International Symposium on Cyberspace Safety and Security, and 2015  
599 IEEE 12th International Conference on Embedded Software and Systems. IEEE, 2015, pp. 1314–1319.
- 600 28. Momtazi, S. Fine-grained German Sentiment Analysis on Social Media. Proceedings of the Eighth  
601 International Conference on Language Resources and Evaluation, LREC 2012, Istanbul, Turkey, May 23–25,  
602 2012, 2012, pp. 1215–1220.
- 603 29. Rothfels, J.; Tibshirani, J. Unsupervised sentiment classification of English movie reviews using automatic  
604 selection of positive and negative sentiment items. Technical Report, Stanford University, 2010.
- 605 30. Cheng, K.; Li, J.; Tang, J.; Liu, H. Unsupervised sentiment analysis with signed social networks. Thirty-First  
606 AAAI Conference on Artificial Intelligence, 2017.
- 607 31. Tripathy, A.; Agrawal, A.; Rath, S.K. Classification of sentiment reviews using n-gram machine learning  
608 approach. *Expert Systems with Applications* **2016**, *57*, 117–126.
- 609 32. Reyes, A.; Rosso, P.; Buscaldi, D. From humor recognition to irony detection: The figurative language  
610 of social media. *Data & Knowledge Engineering* **2012**, *74*, 1 – 12. Applications of Natural Language to  
611 Information Systems.
- 612 33. Hamdan, H.; Bellot, P.; Bechet, F. Lsif: Crf and logistic regression for opinion target extraction and  
613 sentiment polarity analysis. Proceedings of the 9th international workshop on semantic evaluation  
614 (SemEval 2015), 2015, pp. 753–758.
- 615 34. Dragoni, M.; da Costa Pereira, C.; Tettamanzi, A.G.; Villata, S. Combining argumentation and aspect-based  
616 opinion mining: the smack system. *AI Communications* **2018**, *31*, 75–95.
- 617 35. Pavlopoulos, J.; Sorensen, J.; Dixon, L.; Thain, N.; Androutsopoulos, I. Toxicity Detection:  
618 Does Context Really Matter? Proceedings of the 58th Annual Meeting of the Association for  
619 Computational Linguistics; Association for Computational Linguistics: Online, 2020; pp. 4296–4305.  
620 doi:10.18653/v1/2020.acl-main.396.
- 621 36. Saif, H.; He, Y.; Alani, H. Semantic Sentiment Analysis of Twitter. Proceedings of the 11th International  
622 Conference on The Semantic Web - Volume Part I; Springer-Verlag: Berlin, Heidelberg, 2012; ISWC12, pp.  
623 508–524.
- 624 37. Brassard-Gourdeau, E.; Khoury, R. Subversive toxicity detection using sentiment information. Proceedings  
625 of the Third Workshop on Abusive Language Online, 2019, pp. 1–10.
- 626 38. Gangemi, A.; Presutti, V.; Recupero, D.R. Frame-Based Detection of Opinion Holders and Topics: A Model  
627 and a Tool. *IEEE Comp. Int. Mag.* **2014**, *9*, 20–30. doi:10.1109/MCI.2013.2291688.
- 628 39. Recupero, D.R.; Presutti, V.; Consoli, S.; Gangemi, A.; Nuzzolese, A.G. Sentilo: Frame-Based Sentiment  
629 Analysis. *Cognitive Computation* **2015**, *7*, 211–225. doi:10.1007/s12559-014-9302-z.
- 630 40. Wright, A.; Shaikh, O.; Park, H.; Epperson, W.; Ahmed, M.; Pinel, S.; Yang, D.; Chau, D.H. RECAST:  
631 Interactive Auditing of Automatic Toxicity Detection Models. Conference: Chinese CHI 2020: The eighth  
632 International Workshop of Chinese CHI, 2020, pp. 80–82. doi:10.1145/3403676.3403691.
- 633 41. Han, X.; Tsvetkov, Y. Fortifying Toxic Speech Detectors Against Veiled Toxicity. Proceedings of the  
634 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP); Association for  
635 Computational Linguistics: Online, 2020; pp. 7732–7739. doi:10.18653/v1/2020.emnlp-main.622.
- 636 42. Morzhov, S. Avoiding Unintended Bias in Toxicity Classification with Neural Networks. 2020 26th  
637 Conference of Open Innovations Association (FRUCT), 2020, pp. 314–320.  
638 doi:10.23919/FRUCT48808.2020.9087368.
- 639 43. Dessi, D.; Fenu, G.; Marras, M.; Recupero, D.R. COCO: Semantic-Enriched Collection of Online Courses at  
640 Scale with Experimental Use Cases. World Conference on Information Systems and Technologies. Springer,  
641 2018, pp. 1386–1396.
- 642 44. Georgakopoulos, S.V.; Tasoulis, S.K.; Vrahatis, A.G.; Plagianakos, V.P. Convolutional neural networks for  
643 toxic comment classification. Proceedings of the 10th Hellenic Conference on Artificial Intelligence. ACM,  
644 2018, p. 35.
- 645 45. Martens, M.; Shen, S.; Iosup, A.; Kuipers, F. Toxicity Detection in Multiplayer Online Games. The 14th  
646 International Workshop on Network and Systems Support for Games (NetGames), At Zagreb, Croatia,  
647 2015. doi:10.1109/NetGames.2015.7382991.
- 648 46. Atzeni, M.; Recupero, D.R. Multi-domain sentiment analysis with mimicked and polarized word  
649 embeddings for human–robot interaction. *Future Generation Computer Systems* **2019**.

- 650 47. Pinter, Y.; Guthrie, R.; Eisenstein, J. Mimicking word embeddings using subword rnns. *arXiv preprint*  
651 *arXiv:1707.06961* **2017**.
- 652 48. Si, Y.; Wang, J.; Xu, H.; Roberts, K. Enhancing Clinical Concept Extraction with Contextual Embedding.  
653 *arXiv preprint arXiv:1902.08691* **2019**.
- 654 49. Reimers, N.; Schiller, B.; Beck, T.; Daxenberger, J.; Stab, C.; Gurevych, I. Classification and Clustering of  
655 Arguments with Contextualized Word Embeddings. *arXiv preprint arXiv:1906.09821* **2019**.
- 656 50. Read, J.; Pfahringer, B.; Holmes, G.; Frank, E. Classifier chains for multi-label classification. *Machine*  
657 *Learning* **2011**, *85*, 333. doi:10.1007/s10994-011-5256-5.
- 658 51. Dessì, D.; Fenu, G.; Marras, M.; Recupero, D.R. Bridging learning analytics and Cognitive Computing for  
659 Big Data classification in micro-learning video collections. *Computers in Human Behavior* **2019**, *92*, 468–477.
- 660 52. Dessì, D.; Recupero, D.R.; Fenu, G.; Consoli, S. A recommender system of medical reports leveraging  
661 cognitive computing and frame semantics. In *Machine Learning Paradigms*; Springer, 2019; pp. 7–30.
- 662 53. Dang, N.C.; Moreno-García, M.N.; De la Prieta, F. Sentiment Analysis Based on Deep Learning: A  
663 Comparative Study. *Electronics* **2020**, *9*, 483. doi:10.3390/electronics9030483.
- 664 54. Lemaître, G.; Nogueira, F.; Aridas, C.K. Imbalanced-learn: A python toolbox to tackle the curse of  
665 imbalanced datasets in machine learning. *The Journal of Machine Learning Research* **2017**, *18*, 559–563.
- 666 55. Lilleberg, J.; Zhu, Y.; Zhang, Y. Support vector machines and Word2vec for text classification with semantic  
667 features. 2015 IEEE 14th International Conference on Cognitive Informatics Cognitive Computing  
668 (ICCI\*CC), 2015, pp. 136–140. doi:10.1109/ICCI-CC.2015.7259377.
- 669 56. Davis, J.; Goadrich, M. The Relationship between Precision-Recall and ROC Curves. Proceedings of the  
670 23rd International Conference on Machine Learning; Association for Computing Machinery: New York,  
671 NY, USA, 2006; ICML '06, p. 233–240. doi:10.1145/1143844.1143874.
- 672 57. Peters, M.E.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; Zettlemoyer, L. Deep contextualized  
673 word representations. Proc. of NAACL, 2018.

674 © 2021 by the authors. Submitted to *Electronics* for possible open access publication  
675 under the terms and conditions of the Creative Commons Attribution (CC BY) license  
676 (<http://creativecommons.org/licenses/by/4.0/>).