

## Article

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

Danilo Dessì <sup>1,2,\*</sup> , Diego Reforgiato Recupero <sup>3,†</sup>  and Harald Sack <sup>1,2,†</sup> <sup>1</sup> FIZ Karlsruhe–Leibniz Institute for Information Infrastructure, Hermann-von-Helmholtz-Platz 1, 76344 Eggenstein-Leopoldshafen, Germany; harald.sack@fiz-karlsruhe.de<sup>2</sup> Karlsruhe Institute of Technology, Institute AIFB, Kaiserstraße 89, 76133 Karlsruhe, Germany<sup>3</sup> Department of Mathematics and Computer Science, University of Cagliari, 09124 Cagliari, Italy; diego.reforgiato@unica.it

\* Correspondence: danilo.dessi@fiz-karlsruhe.de

† The authors equally contributed to the research performed in this paper.

**Abstract:** Today, increasing numbers of people are interacting online and a lot of textual comments are being produced due to the explosion of online communication. 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. Deep learning technologies have recently 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 for identifying 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



**Citation:** Dessì, D.; Reforgiato Recupero, D.; Sack, H. An Assessment of Deep Learning Models and Word Embeddings for Toxicity Detection within Online Textual Comments. *Electronics* **2021**, *10*, 779. <https://doi.org/10.3390/electronics10070779>

Academic Editor: Gwanggil Jeon

Received: 11 February 2021

Accepted: 20 March 2021

Published: 25 March 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 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 to be the digital version of personal attacks (e.g., bullying behaviors) that can cause social problems (e.g., racism), and they are felt to be 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 (<https://medium.com/analytics-vidhya/twitter-toxicity-detector-using-tensorflow-js-1140e5ab57ee>, accessed on 11 February 2021). An attempt to deal with this issue is the introduction of crowdsourcing voting schemes that give the possibility to denounce inappropriate comments in online environments to the users. Among many others, Facebook, for example, allows its users to report a post in terms of violence or hate speech [1]. This scheme allows Facebook to identify fake accounts,

offensive comments, etc. However, these methodologies are often inefficient, as they fail to detect toxic comments in real time [2], becoming a requirement within social network communities. A toxic post might have been published online much earlier than the time it is reported and, during the time it is online, it might cause problems and offenses to several users which might have undesired behaviors (e.g., leaving the underlying social platform). Therefore, detecting toxicity within textual comments through novel technologies has great relevance in the prevention of adverse social effects in a timely and appropriate manner within online environments [3].

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

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

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

Within this scenario, in this paper various deep learning models that are fed by word embeddings are designed and evaluated to recognize toxicity levels within textual comments. In details, four deep learning models built using the Keras (<https://keras.io/>, accessed on 11 February 2021) framework are designed, and four different types of word embeddings are analyzed.

To this aim, the current state-of-the-art toxicity dataset released during the Kaggle challenge on toxic comments (<https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/overview>, accessed on 11 February 2021) is used.

The reader notices that this paper analyzes the performances of deep learning and classical Machine Learning approaches (using tf-idf and word embeddings) when tackling the task of toxicity detection. Basically, we want to assess whether the syntactic and semantic information lying within the text can provide hints on the presence of certain toxicity classes. This is not possible in some domains and tasks: for example, for the problem of identifying empathetic vs. non-empathetic discussion within answers of a therapist during motivational interviews, it has been initially observed that syntactic and

semantic information do not provide any clue for the classification task, leading to very low accuracies [11].

Thus, for a fair analysis, it is important that the dataset does not contain any unbalanceness. Machine Learning classifiers fail to cope with imbalanced training datasets, as they are sensitive to the proportions of the different classes [12]. As a consequence, these algorithms tend to favor the class with the largest proportion of observations, which may lead to misleading accuracies. That is why we preprocessed the mentioned dataset to make it balanced and then applied a 10-fold cross-validation to tackle the proposed task.

Thus, this paper provides the following contributions:

- We analyzed four deep learning models based on Dense, Convolutional Neural Network (CNN), and Long-Short Term Memory (LSTM) layers to detect various levels of toxicity within online textual comments.
- We evaluate the use of four word embedding representations based on *Word2Vec* [9,13] and Bidirectional Encoder Representations from Transformer (*BERT*) [10] algorithms for the task of toxicity detection in online textual comments.
- We provide a comparison between deep learning models against common baselines that are used within classification tasks of textual resources.
- We release contextual word embeddings resource trained on a dataset, including toxic comments.
- We also release mimicked word embeddings of tokens that are missing in the pre-trained Google Word2Vec (<https://code.google.com/archive/p/word2vec/>, accessed on 11 February 2021) word embeddings.

The source code that is used for this study is freely available through a GitHub repository (<https://github.com/danilo-dessi/toxicity>, accessed on 11 February 2021).

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

## 2. Related Work

A few past works have already addressed the challenge of detecting toxicity within textual comments that are left by users within online environments. Generally, they rely on Sentiment Analysis methods [14–19] to detect and extract the subjective information and classify emotions and sentiments to determine whether a toxicity facet is present or not. For doing so, NLP, Machine Learning, Text Mining, and Computational Linguistics are the most prominent technologies employed [20,21]. Sentiment Analysis methods, like many others within the Machine Learning domain, can be mainly split into two categories. i.e., supervised and unsupervised. Supervised techniques require the use of labeled data (training set) to train a model that can be applied to unseen data in order to predict a sentiment or an emotion [22–24]. These methods often are limited by the lack of labeled data, or by the fact that there are not either good or enough examples for certain categories (e.g., in case of dataset imbalance) [25]. On the other hand, unsupervised Sentiment Analysis approaches usually rely on semantic resources, like lexicons, where words are assigned to scores for reflecting words relevance for target categories to infer sentiments and emotions of the input data [26–28]. Supervised and unsupervised approaches are both largely explored in literature for Sentiment Analysis tasks, which include Sentiment Analysis polarity detection (i.e., identifying whether a certain text is either positive or negative) [29], figurative-language uncovering (understanding if the input text is figurative or objective) [21,30], aspect-based polarity detection (e.g., assigning sentiment polarity to features of a certain topic, such as the screen of an iPhone) [31,32], sentiment scores prediction (e.g., identifying a continuous number in  $[-1, 1]$  to a certain topic or text) [4], and so on.

However, these methodologies have only recently been explored for toxicity detection [33], although the need to monitor online communications to identify toxicity and make the communications safe and respectful is an old and still open issue. Hence, the gap between the current methodologies and their potential use within toxicity detection remains an open challenge. Therefore, dealing with toxicity raises new challenges and research opportunities where deep learning-based approaches for Sentiment Analysis can have a relevant role in making advancements for the identification of toxicity levels.

Additionally, Semantic Web technologies are being used within Sentiment Analysis tasks. It has been proved that they bring several benefits, leading to higher accuracy [34]. For example, the use of sentiment-based technologies to detect toxicity is investigated in [35]. However, the use of word embedding representation is not taken into account. A work worth noting is [21], where the authors analyzed the problem of figurative language detection in social media. More in detail, they focused on the use of semantic features that were extracted with Framester for identifying irony and sarcasm. Semantic features have been extracted to enrich the representation of input tweets with event information using frames and word senses in addition to lexical units. Sentilo [36,37] represents one more example of an unsupervised method that exploits Semantic Web technologies. Given a statement expressing an opinion, Sentilo recognizes its holder, detects its related topics and subtopics, links them to relevant situations and events referred to by it, and evaluates the sentiment expressed on each topic/subtopic. Moreover, Sentilo is domain-independent and it relies on a novel lexical resource, which enables a proper propagation of the sentiment scores from topics to subtopics. Its output is represented as an RDF graph and, where applicable, it resolves holders' and topics' identity on Linked Data.

Recently, the authors in [33] discussed the problem of toxicity detection and proved that context can both amplify or mitigate the perceived toxicity of posts. Besides, they found no evidence that context actually improves the performance of toxicity classifiers. In another work [38] the authors presented an interactive tool for auditing toxicity detection models by visualizing explanations for predictions and providing alternative wordings for detected toxic speech. In particular, they displayed the attention of toxicity detection models on user input, providing suggestions on how to replace sensitive text with less toxic words.

Others, Ref. [39], tackled the problem of identifying disguised offensive language, such as adversarial attacks that avoid known toxic patterns and lexicons. To do that, they proposed a framework to fortify existing toxic speech detectors without a large labeled corpus of veiled toxicity. In particular, they augmented the toxic speech detector's training data with new discovered offensive examples.

Deep learning technologies have been leveraged by the authors in [40] to tackle the problem of toxic comments detection. More in details, the authors introduced two state-of-the-art neural network architectures and demonstrate how to employ a contextual language representation model.

One more work that deals with a sentiment toxicity detection problem is [7], where the authors adopt both pre-trained word embeddings and close-domain word embeddings previously trained on a large dataset of users' comments [41].

However, their approach is based on a Logistic Regression (LR) classifier and it does not use state-of-the-art deep learning technologies. Well established methodologies (e.g., k-nearest neighbors (kNN), Naive Bayes (NB), Support Vector Machines (SVM), etc.) are today outperformed for the same tasks by CNN-based models by [42].

One more work for toxicity detection is proposed by authors in [43] and it lies within the context of multiplayer online games. There, social interactions are an essential feature for a growing number of players worldwide. This interaction might bring undesired and unintended behavior, especially if the game is designed to be highly competitive. They defined toxicity as the use of profane language by one player to insult or humiliate another player in the same team. Given the specific domain, the use of bad words is a necessary, but not sufficient, condition for toxicity, as they can be used to curse without the intent

to offend anyone. The authors looked at the 100 most frequently used n-grams for  $n = 1, 2, 3, 4$  and manually determined which of them are toxic or not. With such training data, they use a SVM to predict the odds of winning for each team to observers based on their communication, while the match is still going.

Another work that embraces both deep learning and word embeddings for toxicity detection is reported in [1], where FastText (<https://fasttext.cc/docs/en/english-vectors.html>, accessed on 11 February 2021) pre-trained embeddings are used to feed four different deep learning models that are based on CNN, Long Short Term Memory (LSTM), and Gated Recurrent Unit (GRU) layers. However, the experiments show weak results probably due to the class imbalance of classes. Conversely, in this work deep learning models by using a balanced dataset are trained, when considering one toxicity class at a time, and trying to better represent the input texts by using word embeddings tuned to the target domain. More precisely, in one set of experiments, domain generated word embeddings are created through mimicking techniques; this allows to face slang, misspellings, or obfuscated contents not represented within pre-trained word embedding representations [24,44]. Besides the Word2Vec embeddings, state-of-the-art word embeddings, called BERT [10,45,46], are used to tune the vectors to the context where words are used.

### 3. Problem Formulation

The problem that is 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

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

where  $l$  is a binary label that can only assume values in  $\{0, 1\}$  and indicates whether the toxicity  $t$  is present in  $c$  (i.e.,  $l$  takes the value 1) or not (i.e.,  $l$  takes the value 0). Therefore, with such an approach, an independent binary classifier for each toxicity label is trained. Given an unseen sample, each binary classifier predicts whether that underlying toxicity is present or not in the sample. The combined model then predicts all of the labels for this sample for which the respective classifier predicts a positive result. Although this method of dividing the task into multiple binary tasks may resemble superficially the one-vs-all and one-vs-rest methods for multi-class classification, it is essentially different from them, because a single binary classifier deals with a single label without any regard to other labels whatsoever. This means that each binary classification task that we formulated does not benefit from the information of the other labels at training time. However, this mapping is straightforward and it does not change the semantic of the input problem [47]. By building these models for various  $t$ , the performances of the proposed solutions are evaluated with the goal of finding which combination of the deep learning layers and word embeddings can better capture the text peculiarities for toxicity detection.

### 4. The Proposed Approach

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

#### 4.1. Preprocessing

Text preprocessing techniques, such as stop words and punctuation removal, lemmatization, stemming, matching words with a dictionary to correct grammar, removing words containing alpha-numeric characters, and so on, are common practices when Machine Learning algorithms are applied [48,49], and text representation is generated as a result of different feature engineering processes. However, with the introduction of deep learning approaches, these techniques have not shown promising results. The reason is that neural networks learn from any element found within the text, because each token contributes to the sentence semantics. Therefore, although certain terms might be included in existing stop word lists, they are maintained because they can enrich the semantics of text con-



tent and improve the performance of the deep learning model [1]. Hence, as suggested by authors in [1], all of the above mentioned preprocessing steps are ignored; only the conversion of texts in lower case is performed. Afterward, the whole set of input text is ready to feed a deep learning model. More precisely, imagine to have a toxicity target class  $t$  and a set of pairs  $P = \{(c_0, l_0), \dots, (c_n, l_n)\}$ , 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  $c_i$  does not include the toxicity  $t$  or 1 if the comment  $c_i$  expresses some level of toxicity  $t$ . From the set  $P$ , the set  $P' = \{(c'_0, l_0), \dots, (c'_n, l_n)\}$  is derived, where each comment  $c'_i$  is an integer-encoded comment of the original  $c_i$ . In details, let  $W$  be the list of all the words that belong to all textual comments, and  $WS$  the set of all the words in  $W$  without duplicates (i.e.,  $WS$  has only one occurrence for each input word, whereas  $W$  can contain multiple occurrences of the same element). Subsequently, two functions  $\theta$  and  $\phi$ , which map the elements in  $W$  and  $WS$  to unique integer values, respectively, are built. For example, consider the sentence *you both shut up or you both die* and imagine to have the toy functions  $\theta_{toy}$  and  $\phi_{toy}$ . The function  $\theta_{toy}$  maps “you” to “7”, “both” to “43”, “shut” to “22”, “up” to “76”, “or” to “10”, “you” to “3”, “both” to “41”, and “die” to “50”. The function  $\phi_{toy}$  maps “you” to “7”, “both” to “43”, “shut” to “22”, “up” to “76”, “or” to “10”, and “die” to “50”. Subsequently, the integer-encoded sentence 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 notices that, by using  $\theta_{toy}$ , the words “you” and “both” are mapped to different integers. Within our approach, the function  $\theta$  is used for BERT word embeddings, whereas the function  $\phi$  is used to encode the input text when Word2Vec word embeddings are employed.

#### 4.2. Deep Learning Models

Figure 1 shows the designed deep learning model schemes. In particular, we illustrate four deep learning models based on *Dense*, *CNN*, and *LSTM* layers that are available within the Keras framework (<https://keras.io/>, accessed on 11 February 2021). All the models present the same number of layers. It is worth noting that the input and output layers among the models are the same to better compare their performances when only considering the type of neural network that they adopt. More precisely, the input layer is an *Embedding* layer, which has the goal of mapping the words of the input text to the underlying word embeddings. The last layer is a *Dense* layer that maps the intermediate results of the models in a single label that can only take the values 0 and 1. For doing so, it uses the *sigmoid* activation function to compute a probability that can be easily used to obtain the correct label value. In the next paragraphs, we will give more details about the deep learning layers.

The literature already showed [50] that deep learning methods that are trained with word embeddings outperform those trained with tf-idf features. Therefore, we did not include the latter in our analysis, as we believe that they would not add additional value to the current evaluation.

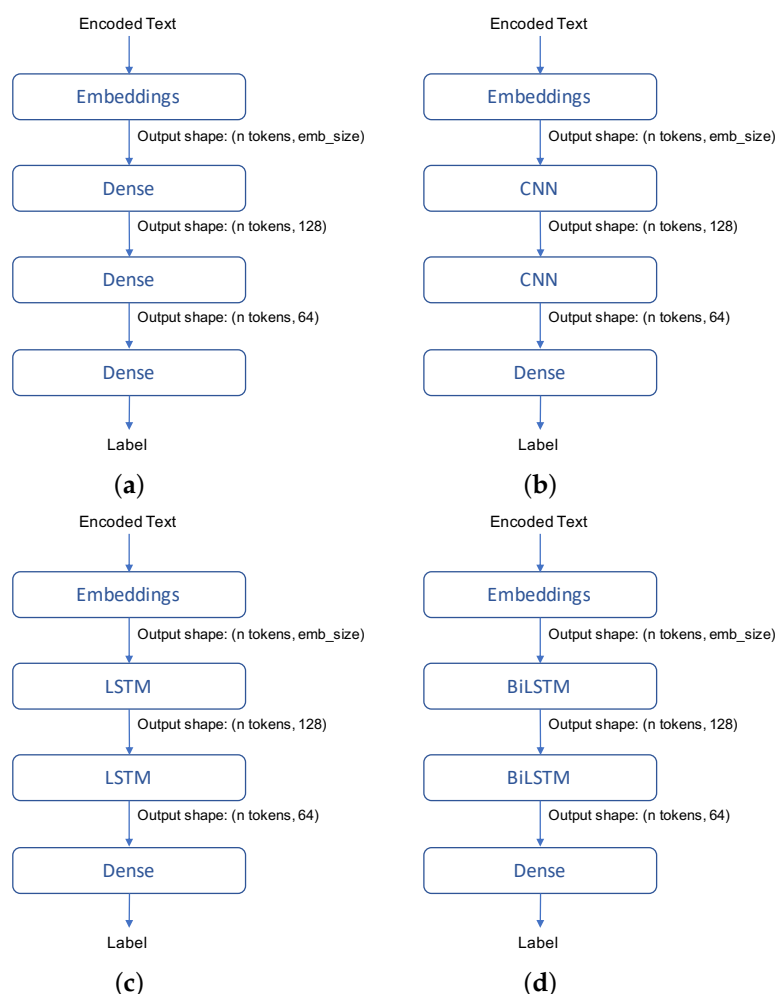
##### 4.2.1. Dense Model

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

##### 4.2.2. CNN Model

The CNN model that is depicted in Figure 1b is based on inner CNN layers. These layers perform filtering operations to detect meaningful features of textual input for the target toxicity facet. Filters can be envisioned as kernels that slide on the vector representation and perform the same operations on each element until all of the vectors have been covered. Two kernels of size 10 for the first layer and size 5 for the second layer are used.

For these layers, the same number of neurons previously introduced for the *Dense* layers is used to better compare the model performances.



**Figure 1.** The deep learning models. (a) Dense (b) Convolutional Neural Network (CNN) (c) Long-Short Term Memory (LSTM) (d) Bidirectional LSTM. The output shape of the employed layers is indicated within the parenthesis.

#### 4.2.3. LSTM Model

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

#### 4.2.4. Bidirectional LSTM

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

### 4.3. Word Embeddings Representations

In this section, the word embedding representations that are used to model the syntactic and semantic properties of the words in vectors of real numbers are introduced. Within this work, the employed word embedding representations are *Word2Vec* [9,13] and BERT [10]. We chose the most common sizes for the embeddings, i.e., 300 for *Word2Vec* embeddings, and 1024 for BERT word embeddings.

#### 4.3.1. Word2Vec

The *Word2Vec* [9,13] word embedding generator aims to detect the meaning and semantic relations among the words by investigating the co-occurrence of words in documents within a given corpus. The idea behind this algorithm is to model the context of words by exploiting Machine Learning and statistics and come up with a vector representation for each word within the corpus. The resulting word vector representations allow the recognition of relatedness between words. For example, the verbs *capture* and *catch*, which are syntactically different but share common meaning and present analogous co-occurring words, are associated to similar vectors. A *Word2Vec* model can be trained by using either the Continuous Bag-Of-Words (CBOW) or the Skip-gram algorithm. Within our work, the Skip-gram algorithm is adopted, because, from a preliminary evaluation, it obtained higher performances. In details, the following *Word2Vec* word embeddings are used:

- *Pre-trained*. Pre-trained word embeddings that are released by Google and available online (<https://code.google.com/archive/p/word2vec/>, accessed on 11 February 2021). They are trained on the Google news dataset and they contain more than one billion words. However, their use can be limited by words that could be misspelled (e.g., words with orthographic errors) or domain-dependent words within the input data. These words are commonly referred to as Out Of Vocabulary (OOV) words.
- *Domain-trained*. Domain-trained word embeddings are trained on the original unbalanced dataset (we merged the training and the test set) provided by the Kaggle challenge. The reader notices that we computed the domain-trained embeddings on the new training sets only (at each iteration 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* (<https://radimrehurek.com/gensim/>, accessed on 11 February 2021) 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 they are inferred by exploiting syntactic similarities of words that are in the originally considered vocabulary. More in details, we used the algorithm that was proposed by [44], 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. Subsequently, 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 [44], and the pre-trained Word2Vec Google embeddings.

#### 4.3.2. BERT

The BERT word embeddings model was introduced in late 2018 by the authors in [10]. It is a novel model of pre-trained language representations that allows for 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 that 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* (<https://pypi.org/project/bert-embedding/>, accessed on 11 February 2021) 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 of 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 of 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 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 of the other parameters have been empirically set on the basis of the models performance and previous experiences in past works [4,24]. The experiments have been carried out on a Titan X GPU that was mounted on a server with 16 GB of RAM memory.

#### 5.1. The Dataset

To perform our analysis, we employed the dataset that was released by a Kaggle competition (<https://www.kaggle.com/>, accessed on 11 February 2021). The dataset is collected from Wikipedia comments, which 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 200 k comments and it 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 noting that the dataset is strongly unbalanced, as nearly 90% of the overall comments do not present toxicity. Therefore, 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, as mentioned earlier in the paper. A random model that always predicts the majority class can obtain better performances although it is not able to recognize elements that should belong to the minority class. Hence, having a balanced

dataset is a common procedure in several classification tasks [51] and allows for better understanding 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. Table 1 reports the size of the created datasets for each class 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.

**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>	1962	0.88%	3924
<i>obscene</i>	12,140	5.43%	24,280
<i>threat</i>	689	0.31%	1378
<i>insult</i>	11,304	5.06%	22,608
<i>identity hate</i>	2117	0.95%	4234

## 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 represent textual resources [48]. More precisely, the deep learning models are compared against the following classifiers:

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

For these classical Machine Learning methods employed as baselines, the adoption of just word embeddings is not promising and this has already been shown in literature [52]. In particular, when employing word embeddings for classical Machine Learning methods, they should be processed by operations such as the average or the sum before being fed to a given classifier. This causes a loss of syntactic and semantic information expressed by the embeddings of each word.

To develop the algorithms above, we employed the *scikit-learn* (<https://scikit-learn.org/stable/index.html>, accessed on 11 February 2021) library.

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

**Table 2.** Receiver Operating Characteristic Area Under the Curve (ROC-AUC) values of our deep learning models on each binary classification and average for each model.

Learning Model	Feature	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
		0.915	0.959	0.928	0.968	0.934	0.924	0.938
	domain-trained mimicked bert	0.922	0.969	0.938	0.981	0.941	0.931	0.947
		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
		0.895	0.950	0.857	0.957	0.909	0.903	0.912
	domain-trained mimicked bert	0.906	0.961	0.923	0.974	0.935	0.914	0.936
		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
		0.963	0.980	0.977	0.983	0.968	0.970	0.974
	domain-trained mimicked bert	0.971	0.983	0.977	0.985	0.970	0.977	0.977
		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
		0.963	0.980	0.977	0.984	0.964	0.970	0.973
	domain-trained mimicked bert	0.969	0.963	0.980	0.988	0.970	0.976	0.974
		0.930	0.970	0.939	0.951	0.947	0.941	0.946

### 5.3. Results and Discussion

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

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

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

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

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

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

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

**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	Feature	Toxic			Severe Toxic			Obscene		
		<i>p</i>	<i>r</i>	<i>f</i>	<i>p</i>	<i>r</i>	<i>f</i>	<i>p</i>	<i>r</i>	<i>f</i>
Decision Trees Random Forests MLP	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>
	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
	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	Feature	Threat			Identity Hate			Insult		
		<i>p</i>	<i>r</i>	<i>f</i>	<i>p</i>	<i>r</i>	<i>f</i>	<i>p</i>	<i>r</i>	<i>f</i>
Decision Trees Random Forests MLP	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>
	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
	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

#### 5.4. Comparison with the Kaggle Challenge

The results that are indicated in Table 2 report the ROC-AUC values of our deep learning approaches for each toxicity class and the average over all the classes. The reader notices that it is not the purpose of this paper to compete with the other participants of the Kaggle challenge where the data have been extracted and the evaluation has been reported while using the ROC-AUC. The best three approaches of the challenge were *Toxic Crusaders*, *neongen* & *Computer says no*, and *Adversarial Autoencoder*, which reported

a ROC-AUC value of 0.989, 0.988, and 0.988, respectively. The challenge task was to test any proposed approach on a highly unbalanced dataset. In this paper, we wanted to study how deep learning methods and classical Machine Learning approaches (using tf-idf and word embeddings) perform on the toxicity problem without any bias (unbalanceness of the data). Moreover, it has been proved that optimizing a method for the ROC-AUC does not guarantee the optimization on the precision–recall curve [53]. This is why we included Table 3 with precision, recall, and f-measure metrics computed on the preprocessed balanced dataset. There are several heuristics and tuning that can be done in the presence of unbalanced datasets to help achieve high values of ROC-AUC. Those could not be performed by us, since we used a balanced version of the original dataset.

#### 5.4.1. Baseline Comparison

The results indicate that *Dense*- and *CNN*-based models are not much better than the baseline methods. Actually, in some cases, they are outperformed. For example, considering the toxicity classes *obscene* and *insult*, it is possible to observe that the f-measure that is computed on the baseline predictions is higher than the one obtained by *Dense*- and *CNN*-based models. On the other hand, *LSTM*-based models are able to outperform the baseline methods with a minimum improvement in terms of f-measure of 0.01, i.e., in percentage 1% (see *obscene* class), and a maximum of 0.058, i.e., in percentage 5.8% (see *toxic* class). These results are similar, and sometimes still more noticeable when the *Bidirectional LSTM* layers are employed. Moreover, when considering that, by using the balanced dataset, every classifier is able to obtain a f-measure that is always higher than 0.8, the improvements can be considered to be remarkable. The only drawback is related to the computational time needed to train the deep learning model. Nevertheless, the training time is not reported, since: (i) it is out of the scope of this study; (ii) with modern GPUs, it is feasible to train complex deep learning models; (iii) the training step must be executed only once; and, (iv) the computational time that is needed for the prediction step does not depend on the underlying model used for the training step.

#### 5.4.2. Dense-Based Model

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

#### 5.4.3. CNN-Based Model

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

#### 5.4.4. LSTM-Based Model

The *LSTM* model outperforms both *Dense* and *CNN*-based models, proving its suitability to detect patterns for toxic detection. As previously mentioned, *mimicked* word embeddings are employed for the deep learning model to learn and uncover toxicity from the text comments. *Pre-trained* and *Domain-trained* word embeddings obtain good performances, and their results are not far from the model using the *mimicked* word embeddings.



On the other hand, once again *BERT* is not a good representation for the *LSTM* model. Except for *BERT*, the three other word embeddings that are adopted with the *LSTM* model outperform the baseline methods for almost each toxicity level.

#### 5.4.5. BiLSTM-Based Model

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

#### 5.4.6. Overall Evaluation of the Deep Learning Models

The use of deep learning for the task of toxicity detection has shown good performances in all the toxicity classes. Additionally, it turns out that, despite the small size of datasets employed for certain classes, they are able to detect patterns that allow for correctly performing the classification. More in details, the results suggest that the *Dense* and *CNN* models perform well, since their f-measure is always higher than 0.8, but, for the toxicity detection task, they are outperformed by the *LSTM* and *BiLSTM* models, which obtain a f-measure higher than 0.9 in most of the cases. The results are comparable among the *LSTM* and *BiLSTM* models. However, because *BiLSTM*-based models need a higher computational time to be trained than *LSTM* models, the latter are slightly preferred. It is worth mentioning that the current models are trained without the context that was surrounding the comments in the Wikipedia pages (where the dataset has been originally collected) and, therefore, they might lack the necessary information to predict the correct class. One more obstacle might be also due to the presence of figurative language within the comments, which might change the meaning of the sentences, thus misleading the models. For example, a frequent sentence like *I am going to kill you* pronounced after a mistake or an undesired change in the Wikipedia pages does not necessarily convey a threat or hate emotion, but it may be simply a joke.

#### 5.4.7. Overall Evaluation of Word Embeddings

From the results, it is noticeable that the *Word2Vec* algorithm is a good choice for representing textual resources to be parsed with deep learning models. The results suggest that *mimicked* word embeddings are the best choice, because they enclose the knowledge of *pre-trained* word embeddings that have been built on a large dataset and do not suffer from the OOV words problem [24]. *Domain-trained* word embeddings obtain good results, but, for most of the cases, they are outperformed. This may depend on the fact that the resources that are employed to train these embeddings are not very large and, besides, there is not a sufficient number of examples of toxicity due to the unbalanced number of toxic comments in the input dataset (i.e., more than 200 k comments do not present toxicity; the reader can see Table 1).

Surprisingly, *BERT* embeddings perform badly for the task of toxicity detection, although they are currently the state-of-the-art word embedding representations. A possible motivation behind this finding is that assigning a different embedding to the same word is somehow misleading to the training of the deep learning models. More precisely, the tuning step performed to generate the *BERT* embeddings on our data is not able to capture the context of the words due to the length of some input textual comments and to the typos and incorrect grammar often present within them, thus transferring possible erroneous information to our deep learning models. One more reason might be due to the lack of the surrounding context of the comments; it might have limited the fine-tuning of the

model, therefore leading the semantics of words to be captured badly. This fact is worth investigating, and a close analysis to this problem is required.

## 6. Conclusions and Future Work

In this paper, we presented an assessment of various deep learning models fed by various word embedding representations to detect toxicity within textual comments. From the obtained results, we can definitely state that toxicity can be identified by machine and deep learning approaches fed with syntactic and semantic information extracted from the text. We show how *LSTM*-based model is the first choice among the experimented models to detect toxicity. We also show how various word embeddings may represent the domain knowledge in a variety of ways, and an unique model for all cases might be insufficient. In particular, the results are encouraging when using mimicking techniques to deal with OOV words where there are not many examples to build significant domain-dependent word embeddings. As future works, we plan to perform a deeper assessment of deep learning models by using and combining different layers, to better detect patterns and on real scenarios where classes may be unbalanced as well. Moreover, we would like to investigate other contextualized word embedding representations, such as *ELMO* [54] for the toxicity detection task. An analysis of the proposed approaches on which configuration, parameter settings and heuristic may be added to tackle the same problem but in presence of highly unbalanced datasets is definitely a research direction we would like to investigate as well. Finally, we would like to investigate the impact of using different embeddings for the same word, since it might be the cause of failure of *BERT* embeddings in our experiments. We also think that an ensemble strategy of the proposed approaches should result in better overall performances and are also investigating this direction.

**Author Contributions:** All authors equally contributed to the research performed in this paper. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

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

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

## Abbreviations

The following abbreviations are used in this manuscript:

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
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

## References

1. Saeed, H.H.; Shahzad, K.; Kamiran, F. Overlapping Toxic Sentiment Classification Using Deep Neural Architectures. In Proceedings of the 2018 IEEE International Conference on Data Mining Workshops (ICDMW), Singapore, 17–20 November 2018; pp. 1361–1366.
2. Hosseini, H.; Kannan, S.; Zhang, B.; Poovendran, R. Deceiving google’s perspective api built for detecting toxic comments. *arXiv* **2017**, arXiv:1702.08138.
3. Srivastava, S.; Khurana, P.; Tewari, V. Identifying aggression and toxicity in comments using capsule network. In Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018), Santa Fe, NM, USA, 25 August 2018; pp. 98–105.
4. Dessì, D.; Dragoni, M.; Fenu, G.; Marras, M.; Recupero, D.R. Evaluating neural word embeddings created from online course reviews for sentiment analysis. In Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing, Limassol, Cyprus, 8–12 April 2019; pp. 2124–2127.
5. Dridi, A.; Atzeni, M.; Recupero, D.R. FineNews: Fine-grained semantic sentiment analysis on financial microblogs and news. *Int. J. Mach. Learn. Cybern.* **2019**, *10*, 2199–2207. [[CrossRef](#)]
6. Consoli, S.; Dessì, D.; Fenu, G.; Marras, M. Deep Attention-based Model for Helpfulness Prediction of Healthcare Online Reviews. In Proceedings of the First Workshop on Smart Personal Health Interfaces Co-Located with 25th International Conference on Intelligent User Interfaces (SmartPhil@IUI 2020), Cagliari, Italy, 17 March 2020; pp. 33–49.
7. Carta, S.; Corrigan, A.; Mulas, R.; Recupero, D.R.; Saia, R. A Supervised Multi-class Multi-label Word Embeddings Approach for Toxic Comment Classification. In Proceedings of the 11th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management, Vienna, Austria, 17–19 September 2019; pp. 17–19.
8. Schouten, K.; Frasnica, F. Survey on aspect-level sentiment analysis. *IEEE Trans. Knowl. Data Eng.* **2015**, *28*, 813–830. [[CrossRef](#)]
9. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient estimation of word representations in vector space. *arXiv* **2013**, arXiv:1301.3781.
10. 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); Association for Computational Linguistics: Minneapolis, MN, USA, 2019; pp. 4171–4186. [[CrossRef](#)]
11. Wu, Z.; Helaoui, R.; Kumar, V.; Recupero, D.R.; Riboni, D. Towards Detecting Need for Empathetic Response in Motivational Interviewing. In Proceedings of the Companion Publication of the 2020 International Conference on Multimodal Interaction, ICMI Companion 2020, Virtual Event, Utrecht, The Netherlands, 25–29 October 2020; pp. 497–502. [[CrossRef](#)]
12. Dragoni, M.; Petrucci, G. A neural word embeddings approach for multi-domain sentiment analysis. *IEEE Trans. Affect. Comput.* **2017**, *8*, 457–470. [[CrossRef](#)]
13. Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G.S.; Dean, J. Distributed representations of words and phrases and their compositionality. In Proceedings of the Advances in Neural Information Processing Systems, Lake Tahoe, NV, USA, 5–10 December 2013; pp. 3111–3119.
14. Reforgiato Recupero, D.; Cambria, E. ESWC 14 challenge on Concept-Level Sentiment Analysis. *Commun. Comput. Inf. Sci.* **2014**, *475*, 3–20.
15. Recupero, D.R.; Dragoni, M.; Presutti, V. ESWC 15 Challenge on Concept-Level Sentiment Analysis. Semantic Web Evaluation Challenges. In Proceedings of the Second SemWebEval Challenge at ESWC 2015, Portorož, Slovenia, 31 May–4 June 2015; pp. 211–222. [[CrossRef](#)]
16. Recupero, D.R.; Consoli, S.; Gangemi, A.; Nuzzolese, A.G.; Spampinato, D. A Semantic Web Based Core Engine to Efficiently Perform Sentiment Analysis. In *The Semantic Web: ESWC 2014 Satellite Events*; Presutti, V., Blomqvist, E., Troncy, R., Sack, H., Papadakis, I., Tordai, A., Eds.; Springer: Cham, The Netherlands, 2014; pp. 245–248.
17. Dragoni, M.; Reforgiato Recupero, D. Challenge on fine-grained sentiment analysis within ESWC2016. *Commun. Comput. Inf. Sci.* **2016**, *641*, 79–94.
18. Reforgiato Recupero, D.; Cambria, E.; Di Rosa, E. Semantic sentiment analysis challenge at ESWC2017. *Commun. Comput. Inf. Sci.* **2017**, *769*, 109–123.
19. Kumar, V.; Recupero, D.R.; Riboni, D.; Helaoui, R. Ensembling Classical Machine Learning and Deep Learning Approaches for Morbidity Identification From Clinical Notes. *IEEE Access* **2021**, *9*, 7107–7126. [[CrossRef](#)]
20. Dridi, A.; Recupero, D.R. Leveraging semantics for sentiment polarity detection in social media. *Int. J. Mach. Learn. Cybern.* **2019**, *10*, 2045–2055. [[CrossRef](#)]
21. Recupero, D.R.; Alam, M.; Buscaldi, D.; Grezka, A.; Tavazoei, F. Frame-Based Detection of Figurative Language in Tweets. *IEEE Comput. Intell. Mag.* **2019**, *14*, 77–88. [[CrossRef](#)]
22. Poria, S.; Cambria, E.; Gelbukh, A. Deep convolutional neural network textual features and multiple kernel learning for utterance-level multimodal sentiment analysis. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Lisbon, Portugal, 17–21 September 2015; pp. 2539–2544.

23. Tang, D.; Qin, B.; Liu, T. Document modeling with gated recurrent neural network for sentiment classification. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Lisbon, Portugal, 17–21 September 2015; pp. 1422–1432.
24. Atzeni, M.; Recupero, D.R. Multi-domain sentiment analysis with mimicked and polarized word embeddings for human-robot interaction. *Future Gener. Comput. Syst.* **2020**, *110*, 984–999. [\[CrossRef\]](#)
25. Yin, H.; Gai, K. An empirical study on preprocessing high-dimensional class-imbalanced data for classification. In Proceedings of the 2015 IEEE 17th International Conference on High Performance Computing and Communications, 2015 IEEE 7th International Symposium on Cyberspace Safety and Security, and 2015 IEEE 12th International Conference on Embedded Software and Systems, Lisbon, Portugal, 17–21 September 2015; pp. 1314–1319.
26. Momtazi, S. Fine-grained German Sentiment Analysis on Social Media. In Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC 2012), Istanbul, Turkey, 23–25 May 2012; pp. 1215–1220.
27. Rothfels, J.; Tibshirani, J. *Unsupervised Sentiment Classification of English Movie Reviews Using Automatic Selection of Positive and Negative Sentiment Items*; Technical Report; Stanford University: Stanford, CA, USA, 2010.
28. Cheng, K.; Li, J.; Tang, J.; Liu, H. Unsupervised sentiment analysis with signed social networks. In Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, San Francisco, CA, USA, 4–9 February 2017.
29. Tripathy, A.; Agrawal, A.; Rath, S.K. Classification of sentiment reviews using n-gram machine learning approach. *Expert Syst. Appl.* **2016**, *57*, 117–126. [\[CrossRef\]](#)
30. Reyes, A.; Rosso, P.; Buscaldi, D. From humor recognition to irony detection: The figurative language of social media. *Data Knowl. Eng.* **2012**, *74*, 1–12. [\[CrossRef\]](#)
31. Hamdan, H.; Bellot, P.; Bechet, F. Lsislif: Crf and logistic regression for opinion target extraction and sentiment polarity analysis. In Proceedings of the 9th international workshop on semantic evaluation (SemEval 2015), Denver, CO, USA, 4–5 June 2015; pp. 753–758.
32. Dragoni, M.; da Costa Pereira, C.; Tettamanzi, A.G.; Villata, S. Combining argumentation and aspect-based opinion mining: The smack system. *AI Commun.* **2018**, *31*, 75–95. [\[CrossRef\]](#)
33. Pavlopoulos, J.; Sorensen, J.; Dixon, L.; Thain, N.; Androutsopoulos, I. Toxicity Detection: Does Context Really Matter? In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Online, 5–10 July 2020; pp. 4296–4305. [\[CrossRef\]](#)
34. Saif, H.; He, Y.; Alani, H. Semantic Sentiment Analysis of Twitter. In Proceedings of the 11th International Conference on The Semantic Web (ISWC12), Boston, MA, USA, 11–15 November 2012; pp. 508–524.
35. Brassard-Gourdeau, E.; Khoury, R. Subversive toxicity detection using sentiment information. In Proceedings of the Third Workshop on Abusive Language, Florence, Italy, 1 August 2019; pp. 1–10.
36. Gangemi, A.; Presutti, V.; Recupero, D.R. Frame-Based Detection of Opinion Holders and Topics: A Model and a Tool. *IEEE Comp. Int. Mag.* **2014**, *9*, 20–30. [\[CrossRef\]](#)
37. Recupero, D.R.; Presutti, V.; Consoli, S.; Gangemi, A.; Nuzzolese, A.G. Sentilo: Frame-Based Sentiment Analysis. *Cogn. Comput.* **2015**, *7*, 211–225. [\[CrossRef\]](#)
38. Wright, A.; Shaikh, O.; Park, H.; Epperson, W.; Ahmed, M.; Pinel, S.; Yang, D.; Chau, D.H. RECAST: Interactive Auditing of Automatic Toxicity Detection Models. In Proceedings of the Chinese CHI 2020: The Eighth International Workshop of Chinese CHI, Honolulu, HI, USA, 25–30 April 2020; pp. 80–82. [\[CrossRef\]](#)
39. Han, X.; Tsvetkov, Y. Fortifying Toxic Speech Detectors Against Veiled Toxicity. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), Online, 16–20 November 2020; pp. 7732–7739. [\[CrossRef\]](#)
40. Morzhov, S. Avoiding Unintended Bias in Toxicity Classification with Neural Networks. In Proceedings of the 2020 26th Conference of Open Innovations Association (FRUCT), Yaroslavl, Russia, 20–24 April 2020; pp. 314–320. [\[CrossRef\]](#)
41. Dessì, D.; Fenu, G.; Marras, M.; Recupero, D.R. COCO: Semantic-Enriched Collection of Online Courses at Scale with Experimental Use Cases. In Proceedings of the World Conference on Information Systems and Technologies, Naples, Italy, 27–29 March 2018; pp. 1386–1396.
42. Georgakopoulos, S.V.; Tasoulis, S.K.; Vrahatis, A.G.; Plagianakos, V.P. Convolutional neural networks for toxic comment classification. In Proceedings of the 10th Hellenic Conference on Artificial Intelligence, Patras, Greece, 9–12 July 2018; p. 35.
43. Martens, M.; Shen, S.; Iosup, A.; Kuipers, F. Toxicity Detection in Multiplayer Online Games. In Proceedings of the 14th International Workshop on Network and Systems Support for Games (NetGames), Zagreb, Croatia, 3–4 December 2015. [\[CrossRef\]](#)
44. Pinter, Y.; Guthrie, R.; Eisenstein, J. Mimicking word embeddings using subword rnns. *arXiv* **2017**, arXiv:1707.06961.
45. Si, Y.; Wang, J.; Xu, H.; Roberts, K. Enhancing Clinical Concept Extraction with Contextual Embedding. *arXiv* **2019**, arXiv:1902.08691.
46. Reimers, N.; Schiller, B.; Beck, T.; Daxenberger, J.; Stab, C.; Gurevych, I. Classification and Clustering of Arguments with Contextualized Word Embeddings. *arXiv* **2019**, arXiv:1906.09821.
47. Read, J.; Pfahringer, B.; Holmes, G.; Frank, E. Classifier chains for multi-label classification. *Mach. Learn.* **2011**, *85*, 333. [\[CrossRef\]](#)
48. Dessì, D.; Fenu, G.; Marras, M.; Recupero, D.R. Bridging learning analytics and Cognitive Computing for Big Data classification in micro-learning video collections. *Comput. Hum. Behav.* **2019**, *92*, 468–477. [\[CrossRef\]](#)
49. Dessì, D.; Recupero, D.R.; Fenu, G.; Consoli, S. A recommender system of medical reports leveraging cognitive computing and frame semantics. In *Machine Learning Paradigms*; Springer: Berlin, Germany, 2019; pp. 7–30.

- 
50. Dang, N.C.; Moreno-García, M.N.; De la Prieta, F. Sentiment Analysis Based on Deep Learning: A Comparative Study. *Electronics* **2020**, *9*, 483. [\[CrossRef\]](#)
  51. Lemaître, G.; Nogueira, F.; Aridas, C.K. Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. *J. Mach. Learn. Res.* **2017**, *18*, 559–563.
  52. Lilleberg, J.; Zhu, Y.; Zhang, Y. Support vector machines and Word2vec for text classification with semantic features. In Proceedings of the 2015 IEEE 14th International Conference on Cognitive Informatics Cognitive Computing (ICCI\*CC), Beijing, China, 6–8 July 2015; pp. 136–140. [\[CrossRef\]](#)
  53. Davis, J.; Goadrich, M. The Relationship between Precision-Recall and ROC Curves. In Proceedings of the 23rd International Conference on Machine Learning (ICML '06), Pittsburgh, PA, USA, 25–29 June 2006; pp. 233–240. [\[CrossRef\]](#)
  54. Peters, M.E.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; Zettlemoyer, L. Deep contextualized word representations. In Proceedings of the NAACL, North American Chapter of the Association for Computational Linguistics: Human Language Technologies, New Orleans, LA, USA, 1–6 June 2018.