

Multimodal Detection of Bots on X (Twitter) using Transformers

Loukas Ilias, Ioannis Michail Kazelidis, Dimitris Askounis

Abstract—Although not all bots are malicious, the vast majority of them are responsible for spreading misinformation and manipulating the public opinion about several issues, i.e., elections and many more. Therefore, the early detection of bots is crucial. Although there have been proposed methods for detecting bots in social media, there are still substantial limitations. For instance, existing research initiatives still extract a large number of features and train traditional machine learning algorithms or use GloVe embeddings and train LSTMs. However, feature extraction is a tedious procedure demanding domain expertise. Also, language models based on transformers have been proved to be better than LSTMs. Other approaches create large graphs and train graph neural networks requiring in this way many hours for training and access to computational resources. To tackle these limitations, this is the first study employing only the user description field and images of three channels denoting the type and content of tweets posted by the users. Firstly, we create digital DNA sequences, transform them to 3d images, and apply pretrained models of the vision domain, including EfficientNet, AlexNet, VGG16, etc. Next, we propose a multimodal approach, where we use TwHIN-BERT for getting the textual representation of the user description field and employ VGG16 for acquiring the visual representation for the image modality. We propose three different fusion methods, namely concatenation, gated multimodal unit, and crossmodal attention, for fusing the different modalities and compare their performances. Finally, we present a qualitative analysis of the behavior of our best performing model. Extensive experiments conducted on the Cresci’17 and TwiBot-20 datasets demonstrate valuable advantages of our introduced approaches over state-of-the-art ones.

Index Terms—Bot detection, X (Twitter), digital DNA, image classification, transformers, gated multimodal unit, crossmodal attention

I. INTRODUCTION

Social media platforms, such as Twitter (rebranded to X in 2023) and Reddit, constitute a valuable form of information, where the users have the opportunity to express their feelings, communicate with people around the world, and get informed about the news. For instance, social media have been used for the early detection of mental disorders [1]. Especially in Twitter, people can be informed about everything happening at the moment through the Twitter trends. However, social media are often manipulated by bots. Although not all bots are malicious and dangerous, the majority of them constitute the main form of misinformation [2]. Research has showed that bots can influence elections [3], [4], promote phishing attacks [5], and constitute an effective tool for manipulating social media by spreading articles of low-credibility scores [6].

The authors are with the Decision Support Systems Laboratory, School of Electrical and Computer Engineering, National Technical University of Athens, 15780 Athens, Greece (e-mail: liliask@epu.ntua.gr; gkazelid@gmail.com; askous@epu.ntua.gr).

According to [7], bots dominated hate speech during COVID-19. Therefore, the early detection of bots appears to be imperative nowadays.

Existing research initiatives train shallow machine learning classifiers for recognizing bots in Twitter and propose feature extraction approaches [8], [9]. However, feature extraction constitutes a time-consuming process, demands some level of domain expertise, while the optimal feature set for each specific task may not be found. In addition, there have been introduced deep learning approaches. However, these approaches either use as input a large number of features, including user and tweet metadata, or exploit GloVe embeddings and train LSTM neural networks [10], [11]. However, fine-tuning language models based on transformers requires less time for training and yields better evaluation results. Additionally, methods fine-tuning BERT, RoBERTa, etc. models cannot be applied to other languages except the english one. Moreover, graph-based approaches have been introduced for detecting bots. However, feature extraction techniques are also applied for creating feature vectors per users, which represent the nodes in graphs [12]. At the same time, large graphs are often created, which require a lot of time for training and require access to computational resources.

To address these limitations, we present the first study which exploits only the user description and the sequence of actions that an account performs for recognizing bots in Twitter. In the current study, we introduce both unimodal, i.e., neural networks exploiting either text or image, and multimodal approaches, i.e., neural networks utilizing both text and image. First, motivated by [13], [14], we design the digital DNA per user, which indicates the sequence of actions that an account performs. Next, we adopt the methodology proposed by [15] for converting the DNA sequence into an image. After that, we fine-tune several pretrained models, including AlexNet, ResNet, VGG16, etc., and compare their performances. For recognizing bots through user descriptions, this is the first study fine-tuning TwHIN-BERT [16] in the task of bot detection in Twitter. Regarding the multimodal approaches, we employ VGG16, which constitutes our best performing model, and extract the visual representation. In terms of the textual modality, we employ the TwHIN-BERT and extract the textual representation. Then, we propose three methods for fusing the representations of the different modalities. Specifically, first we concatenate the representation vectors of the two modalities. Second, we exploit a gated multimodal unit (GMU), which controls the importance of each modality towards the final classification. Thirdly, we use a cross-attention mechanism for capturing the inter-modal interactions. Experiments conducted on the Cresci’17 and TwiBot-20 datasets show that the introduced

arXiv:2308.14484v2 [cs.CL] 24 Jul 2024

cross-modal model outperforms the competitive multimodal ones.

Our main contributions can be summarized as follows:

- We create digital DNA sequences per user based on both the type and contents of the tweets and transform these sequences into 3d images.
- We employ and compare many models of the vision domain, including EfficientNet, AlexNet, VGG16, etc. by utilizing the DNA sequence as an image consisting of three channels.
- This is the first study fine-tuning the TwHIN-BERT model for recognizing bots utilizing only the user description field.
- To the best of our knowledge, this is the first study introducing multimodal models employing the user description and the representation of the DNA sequence as an image.
- We present the first study utilizing a gated multimodal unit and cross-attention mechanism, and comparing the fusion methods.

II. RELATED WORK

A. Traditional Machine Learning Algorithms

In [8], the authors proposed two approaches for detecting bots in Twitter both at account and at tweet-level. In terms of the account level classification, the authors extracted a large set of features per user, applied feature selection algorithms, sampling techniques for dealing with the imbalanced datasets, and trained shallow machine learning algorithms.

Similarly, the authors in [10] adopted methods for recognizing bots both at account and tweet level. In terms of the account level, the authors extracted a set of features, including number of statuses (number of tweets and retweets), Followers Count, Friends Count, Favorites Count, Listed Count, Default Profile, etc. Next, the authors combined SMOTE with data enhancement via edited nearest neighbors (SMOTENN) and Tomek Links (SMOTOMEK), and trained traditional machine learning algorithms. In terms of the tweet level classification, the authors introduced a deep neural network consisting of LSTM and dense layers. The authors used as input GloVe embeddings of tweets and tweet metadata. Auxiliary output was also used, whose target was also the classification label.

The authors in [9] extracted a set of style-based features, including the number of punctuation marks, number of hashtags, number of retweets, number of user mentions, number of url links, and many more. They trained the following classification algorithms: Random Forest, Naive Bayes, Logistic Regression, Support Vector Machine with rbf and linear kernel, and Convolution Neural Network.

A wavelet-based approach was introduced by [17]. Specifically, the authors exploited the discrete wavelet transform and extracted a set of features, namely wavelet magnitude, wavelet phase, wavelet domain score, and so on. After this, a feature selection technique, namely a Correlation-based Feature Subset Selection, was proposed for reducing the dimensionality of the feature vector. Finally, a Random Forest Classifier was trained.

The authors in [18] proposed *Botometer*, which is a publicly available bot detection tool via the website, Python, or REST

APIs. The authors use network, user, friends, and temporal features and train a Random Forest Classifier. The only input to this publicly available tool is the screen name.

B. Deep Learning and Transformer-based Approaches

In [19], the authors exploited three types of LSTMs and four types of features, namely tweet content, tweet metadata, account metadata, and user description. In terms of the tweet content and the user description, the authors exploited the GloVe embeddings. Regarding account metadata, the authors extracted 30 features, while 15 tweet metadata features were also extracted. One limitation of this study is the feature extraction procedure and the usage of GloVe embeddings in conjunction with LSTM models instead of pretrained language models based on transformers, which capture better the context and require significantly less time for fine-tuning.

A similar approach was introduced by [11], where the authors utilized GloVe embeddings for representing the user description field as embeddings. Also, the authors extracted several types of features, namely number of followers and friends, number of tweets and retweets posted by the user, length of the name, entropy of the screen name, entropy of the description, etc. They passed the GloVe embeddings of the user description field into LSTM models and concatenated the output of the LSTMs with the rest of the features.

In [8], the authors proposed two approaches for detecting bots in Twitter both at account and at tweet-level. With regards to the tweet-level classification, the authors utilized a deep neural network. Specifically, the authors exploited GloVe embeddings and passed them through BiLSTM layers coupled with an attention mechanism.

Similarly, the authors in [10] adopted methods for recognizing bots both at account and tweet level. In terms of the tweet level classification, the authors introduced a deep neural network consisting of LSTM and dense layers. The authors used as input GloVe embeddings of tweets and tweet metadata. Auxiliary output was also used, whose target was also the classification label.

A different approach was proposed by [20], where the authors exploited Generative Adversarial Neural networks (GANs) for recognizing bots in Twitter. Specifically, the authors addressed the limitation of original SeqGAN [21] by introducing the GANBOT setting, where a shared LSTM was used between generator and classifier. Similar to [10], the authors used GloVe embeddings as input to LSTMs. GANs were also exploited by [22]. Specifically, the authors extracted first a set of features consisting of user meta-data, sentiment, friends, content, network, and timing. Finally, they proposed a new method based on Conditional GANs by introducing the Wasserstein distance with a gradient penalty and a condition-generation method through the modified density peak clustering algorithm. The authors compared their approaches with Random Oversampling, SMOTE, and ADASYN.

In [23], the authors use language models based on transformers, namely BERT and RoBERTa, for recognizing bots in Twitter. The authors concatenated the text embeddings along with additional metadata and trained a neural network consisting of dense layers.

In [15], the authors proposed a new method for recognizing bots in Twitter based on image recognition. Specifically, they introduced a new approach for converting the digital DNA sequence into an image consisting of three channels. Finally, the authors applied ResNet and stated that this method achieved competitive results to state-of-the-art approaches. However, the authors created a digital DNA sequence based only on the type of the tweets, while they fine-tuned only one pretrained model.

The authors in [24] utilized the concept of digital DNA sequences along with pretrained BERT models. First, the authors trained BERT on sentiment analysis classification tasks. Next, they passed each tweet through a BERT model and predicted the sentiment of each tweet, i.e., positive, negative, and neutral. In this way, they create a digital DNA sequence consisting of characters denoting the sentiment of each tweet.

A deep learning architecture consisting of two branches was proposed by [25]. Specifically, the authors used word2vec and passed the respective embeddings through CNNs. Next, they extracted features, including number of followers, number of tweets, account age/reputation, etc., and passed them through dense layers. After this, they concatenated the two branches and passed the resulting vector through a dense layer for getting the final prediction.

An active learning approach was adopted by [26]. Specifically, the authors extracted a set of features, including metadata-based, interaction-based, content-based, and timing-based. Next, active learning was employed for efficiently expanding the labeled data. Finally, a deep neural network consisting of ResNet, BiGRU, and attention layer was trained.

The study in [27] introduced a CNN and BiLSTM-based deep neural network model coupled with an attention mechanism. Specifically, the authors passed the profile, temporal, and activity information to a two-layers stacked BiLSTM, whereas they passed the content information to a deep CNN. An attention mechanism was exploited at the top of the proposed architecture. Findings showed that the proposed approach outperformed the state-of-the-art approaches.

A CNN-LSTM network was also proposed by [28]. Specifically, the authors model the social behaviour per user, i.e., posting and retweeting, and exploit an LSTM network. Next, the authors consider users' history tweets as temporal text data and use a CNN-LSTM network. Finally, they fuse the representation vectors obtained by the two networks and get the final prediction.

In [29], the authors extract word embeddings, character embeddings, part-of-speech embeddings, and named-entity embeddings. They pass these embeddings through BiLGRU [30] and get the final prediction. Findings stated that the proposed model performed better or comparably to state-of-the-art Twitter bot detection models.

In [31], the authors used GloVe embeddings as input to a 3-layer BiLSTM neural network for distinguishing Twitter bots from human accounts.

C. Graph-based Approaches

Recently, graph convolutional networks have been also used for identifying bots in Twitter. Specifically, the study in [12]

created a heterogeneous graph and employed Relational Graph Convolutional Networks. The authors represented each user as a vector consisting of user description, tweets, numerical and categorical properties. To be more precise, the authors employed RoBERTa to get a representation vector of the user description. Similarly, they used RoBERTa and averaged the representations of all the tweets posted by the user into one single vector. In terms of the numerical properties, the authors created a vector consisting of the number of friends, followers, favourites, etc. Regarding the categorical properties, the resulting vector was composed of binary variables, i.e., if the user has a profile image, if the user is verified, etc. However, this method cannot incorporate the intrinsic heterogeneity of relation. Therefore, the authors in [32] introduced propose a novel Twitter bot detection framework that is graph-based and heterogeneity-aware. Specifically, the authors encoded each user adopting the procedure followed in [12]. Finally, the authors exploited relational graph transformers and semantic attention networks.

The authors in [33] proposed a text-graph interaction module along with a semantic consistency module. Specifically, the authors exploited the description and tweets and passed them through RoBERTa for getting the textual representation (embeddings). For graph-based approach, they adopted the method proposed by [12]. Next, they introduced a method for capturing the interactions between the text and graph modality. After that, a semantic consistency detection module was proposed, which exploits the attention weights obtained by the RoBERTa model. Finally, the respective vectors were concatenated and were passed through a dense layer for the final classification.

Multi-view graph attention networks were exploited in [34]. The authors used different datasets for training and evaluating their proposed approaches under a transfer learning scenario. Profile features were also utilized, including number of followers/friends, age in days, and many more.

Motivated by the fact that existing methods almost ignore the differences in bot behaviors in multiple domains, the study in [35] introduces a domain-aware approach for recognizing bots. Specifically, multi-relational graphs were exploited coupled with a user representation learning module, consisting of a series of graph embedding layers and semantic attention layers. Finally, domain-aware classifiers were exploited for detecting bots. Additionally, according to the authors this is the first study employing a federated learning framework.

The authors in [36] introduced a Graph Convolutional Neural Network for exploiting the characteristics of the accounts' neighbours and identifying bots. Each user is represented by a node consisting of the following feature set: age, favourites_count, statuses_count, account length name, followers_count, and friends_count. Results showed that the proposed approach outperformed the state-of-the-art ones.

D. Unsupervised Learning

Cresci et al. [14] introduced digital DNA for recognizing bots in Twitter. Specifically, the authors set the tweets of each user into chronological user and created a digital DNA sequence according to the type and content of each tweet. In terms of the

type of the tweets, the resulting DNA sequence was composed of A, C, and T, where a "C" indicates a reply, a "T" denotes a retweet, and an "A" denotes a simple tweet. With regards to the content of the tweet, the authors examined whether the tweet contained hashtags, URLs, mentions, media, etc. Finally, the authors calculated similarities between these sequences using the longest common subsequence (LCS) similarity and clustered users based on their similarity scores.

In [37], the authors introduce MulBot, which constitutes an unsupervised bot detector based on multivariate time series. Specifically, the proposed approach consists of the following steps: multidimensional temporal features extracted from user timelines, dimensionality reduction using an autoencoder, extraction of statistical global features (optional), concatenation of global features with vectorial features in output from the encoder (optional), and clustering algorithm. Results showed that the proposed approach yielded satisfactory results in both the binary task and the multiclass classification one.

The bot detection task was considered as an anomaly detection task in [38]. The authors utilized 95 one-gram features from tweet text along with user features and modified two stream clustering algorithms, namely StreamKM++ and DenStream.

E. Reinforcement Learning

A reinforcement learning approach was proposed by [39] for searching the GNN architecture. In this way, the most suitable multi-hop neighborhood and the number of layers in the GNN architecture are found. An Heterogeneous Information Network was also exploited for modelling the entities and relationships in the social networks. Finally, the authors exploited self-supervised learning approaches. Also, the authors in [40] introduced a deep Q-network architecture by incorporating a Deep Q-Learning (DQL) model. The authors extracted tweet-based, user profile-based, and social graph-based features. In terms of the tweet-based, the authors utilized syntax, semantic, and temporal behaviour features. Regarding the user profile features, the authors employed features pertaining to user behaviour (Posting tweets in several languages, URL/hashtag/mention ratio) and user interactions (Number of active days, Number of retweeted tweets). Finally, they defined social graph-based attributes, such as clustering coefficient, closeness centrality, betweenness centrality and pagerank centrality for each user.

The authors in [41] extracted URL-based features, including URL redirection, frequency of shared URLs, and spam content in URL, and designed a learning automata based model. Specifically, the authors designed a trust computation model, which contains two parameters, namely direct trust and indirect trust. Findings suggested that the proposed approach improved the existing ones.

F. Related Work Review Findings

In spite of the rise of deep learning algorithms, existing research initiatives extract still a large number of features and train traditional machine learning algorithms. However, feature extraction constitutes a tedious procedure demanding domain expertise. Thus, the optimal set of features may not be found.

Additionally, existing research initiatives still exploit GloVe embeddings and train LSTMs instead of employing language models based on transformers, which achieve state-of-the-art results across a large number of domains. In addition, methods employing BERT cannot be employed for languages other than the english one. Recently, there have been introduced graph-based approaches. However, these methods still extract a large number of features per user for creating a feature vector per node-user. At the same time, large graphs are created, which require access to computational resources and need a lot of time for training.

Therefore, our proposed work differs significantly from the aforementioned research works, since we (1) exploit only the user description field and create images according to the type and content of the tweets posted by each user, (2) create images consisting of three channels and train several pretrained models comparing their performances, (3) employ TwHIN-BERT which is a new multi-lingual Tweet language model that is trained on 7 billion Tweets from over 100 distinct languages, (4) introduce multimodal models, which take as input the user description field and the 3d images.

III. DATASETS

A. Cresci'17 Dataset

We use the dataset introduced in [42] for conducting our experiments. Specifically, this dataset consists of genuine accounts and some sets of social and traditional spambots. In the present study, we exploit the set of genuine accounts and the set of social spambots #1. For collecting genuine users, the authors adopted the methodology proposed in [43]. Specifically, the authors contacted with some accounts and asked them a question. Then, the authors examined the users, who replied to the question and verified 3,474 accounts. In terms of the social spambots users, the authors used a group of social bots that was discovered on Twitter during the last Mayoral election in Rome, in 2014. Specifically, this set of social spambots corresponds to retweeters of an Italian political candidate. As mentioned in Section II-D, the authors presented a figure illustrating the LCS curves of humans and bots. Findings showed that the LCS of bots were long, even when the number of account increases. A sudden drop in LCS length of bots is observed when the number of accounts gets close to the group size. On the contrary, the LCS curve of humans indicates little to no similarity. Also, in [42], the authors illustrate the cumulative distribution function (CDF) of join date and number of followers. Results indicate that social spambots have anomalous distributions. In our experiments, we keep users having available the user description field in their profile. To address the issue of the imbalanced dataset, we create a dataset consisting of 943 real users and 943 social spambots. This technique, which addresses the issue of imbalanced dataset by downsampling the set of real users, has been adopted by previous studies [15], [37].

B. TwiBot-20 Dataset

We use the TwiBot-20 dataset to conduct our experiments [44]. Contrary to other datasets which include a specific type of users, this dataset consists of diversified bots that co-exist on the

real-world Twittersphere. Users of this dataset have an interest in four domains, including politics, business, entertainment, and sports. Thus, this dataset constitutes a challenging task, since it includes multiple types of bots. For proving the high quality of annotations, the authors illustrate the CDF of account reputation (number of friends and followers), user tweet counts, and screen name likelihood. CDF plots show that genuine accounts present higher reputation scores than bots in the TwiBot-20 dataset. In terms of user tweet counts, CDF plots indicate that bots generate fewer tweets than humans, in order to avoid the traditional detection methods. Finally, findings of CDF plot regarding the screen name likelihood show that bot users in TwiBot-20 do have slightly lower screen name likelihood. Also, the authors in [44] conduct a user diversity analysis and study the distribution of profile locations and user interests. Findings of a geographic analysis study demonstrate that India, United States, Europe, and Africa are represented in the TwiBot-20 dataset. Finally, the authors examine the most frequent hashtags, present their analysis through a bar plot, and show that twitter users in TwiBot-20 are proved to be diversified in interest domains. TwiBot-20 dataset includes 229,573 users. However, the authors have labelled 11,826 users. This dataset includes the 200 recent tweets per user. The authors in [44] have divided TwiBot-20 dataset into a train, validation, and test set. We remove users who have not posted any tweets and have not added a description field in their profile. Table I presents the distribution of genuine accounts and bots into each set.

TABLE I: TwiBot-20 dataset statistics.

	genuine users	bots
train set	3,262	3,911
development set	959	1,101
test set	494	533

IV. UNIMODAL MODELS UTILIZING ONLY IMAGES FOR DETECTING BOTS IN TWITTER

We adopt the methodology introduced by [13], [14] for constructing a DNA sequence. Inspired by the biological DNA sequence consisting of A (adenine), C (cytosine), G (guanine) and T (thymine), the authors in [13], [14] introduce the digital DNA for creating a DNA sequence based either on type of the tweet or the content of the tweet. In this study, we create two digital DNA sequences based on both the type and content of the tweet. First, we set all tweets into a chronological order. Regarding the first approach, we create a sequence consisting of A, T, and C. A "T" denotes a retweet, a "C" indicates a reply, while an "A" denotes a tweet. In terms of the second approach, which is based on the content of the tweet, we create a digital DNA sequence consisting of N, U, H, M, and X. Specifically, a "N" indicates that the tweet contains no entities (plain text), a "U" denotes that the tweet contains one or more ULRs, a "H" means that the tweet contains one or more hashtags, a "M" indicates that the tweet contains one or more mentions, and a "X" indicates that the tweet contains entities of mixed types.

After having created a digital DNA sequence per user, we adopt the method introduced in [15] for transforming the DNA

sequence into an image consisting of three channels. First, we check whether the length of the longest DNA sequence is a perfect square. If it is a perfect square, then we define the image size as the square root of the length of the longest DNA sequence. Otherwise, we consider the perfect square closest to and larger than the maximum length. In this way, all strings can be converted to images of equal sizes. Next, we assign to each symbol of the DNA sequence, i.e., A, C, and T, a "color" for creating the image. The image is colored pixel by pixel based on the coors assigned to the correspondent symbol. The same methodology is adopted in terms of the digital DNA sequence based on the content of the tweets, which consists of N, U, H, M, and X. The resulting image is a grayscale image in a (1, H, W) format. Similar to [15], we convert the images into a (3, H, W) format utilizing the grayscale transformation¹.

Each image is resized to 256×256 pixels. Next, we fine-tune the following pretrained models: **GoogLeNet (Inception v1)** [45], **ResNet50** [46], **WideResNet-50-2** [47], **AlexNet** [48], **SqueezeNet1_0** [49], **DenseNet-201** [50], **MobileNetV2** [51], **ResNeXt-50 32 \times 4d** [52], **VGG16** [53], and **EfficientNet-B4**² [54].

A. Experiments

1) *Experimental Setup*: In terms of the Cresci'17 dataset, we divide the dataset into a train, validation, and test set (80%-10%-10%). Regarding the TwiBot-20 dataset, the train, validation, and test sets are provided by the authors [44]. We use *EarlyStopping* and stop training if the validation loss has stopped decreasing for 6 consecutive epochs. We use Adam optimizer with a learning rate of $1e-5$ and exploit *ReduceLROnPlateau*, where we reduce the learning rate by a factor of 0.1 if the validation loss has not presented an improvement after 3 consecutive epochs. We minimize the cross-entropy loss function. In terms of the TwiBot-20 dataset, we apply class-weights to the loss function to deal with the data imbalance. We train the proposed approaches for a maximum of 30 epochs. All models have been created using the PyTorch library [55]. All experiments are conducted on a single Tesla T4 GPU.

2) *Evaluation Metrics*: We use Precision, Recall, F1-score, Accuracy, and Specificity for evaluating the results of the proposed approaches. These metrics have been computed by considering the class of bots as the positive one (label=1). Results are obtained over five runs using different random seeds reporting the average and the standard deviation.

3) *Results*: The results of our proposed approaches are reported in Tables II-V. Specifically, Tables II and III report the results on the Cresci'17 dataset, while Tables IV and V report the results on the TwiBot-20 dataset.

a) *Cresci'17 Dataset*: The results of our proposed approaches are reported in Tables II and III.

As one can easily observe in Table II, VGG16 constitutes the best performing model outperforming the other pretrained

¹<https://pytorch.org/vision/main/generated/torchvision.transforms.Grayscale.html>

²We experimented with EfficientNet-B0 to B7, but EfficientNet-B4 was the best performing model.

models in Recall, F1-score, and Accuracy. It is worth noting that there are models surpassing VGG16 or obtaining equal performance in Precision and Specificity. However, VGG16 outperforms these models in F1-score, which constitutes the weighted average of Precision and Recall. Additionally, F1-score is a more important metric than Specificity, since high Specificity and low F1-score indicates that some social spambots are falsely detected as real accounts. VGG16 outperforms the other models in Recall by 0.44-3.48%, in F1-score by 0.22-2.06%, and in Accuracy by 0.19-2.12%. WideResNet-50-2 constitutes the second best performing model attaining an Accuracy of 99.58% and a F1-score of 99.57%. Specifically, it surpasses the other models, except for VGG16, in Accuracy by 0.64-1.91% and in F1-score by 0.59-1.84%. F1-score ranging from 98.17% to 98.98% is obtained by the rest of the models, except for EfficientNet-B4. Similarly, Accuracy ranging from 98.20% to 98.94% is attained by the rest of the models, except for EfficientNet-B4. The highest Precision and Specificity scores are obtained by GoogLeNet and ResNet50 and are equal to 1.00. EfficientNet-B4 obtains the worst evaluation results reaching Accuracy and F1-score up to 97.67% and 97.73% respectively.

As one can easily observe in Table III, VGG16 constitutes our best performing model obtaining an Accuracy of 99.89%, a Precision of 1.00%, a Recall of 99.78%, a F1-score of 99.89%, and a Specificity of 1.00%. Specifically, it outperforms the other models in Recall by 0.44-2.14%, in F1-score by 0.52-2.87%, and in Accuracy by 0.52-2.75%, while it achieves equal Precision and Specificity scores with GoogLeNet, ResNet50, and AlexNet. Additionally, GoogLeNet, WideResNet-50-2, AlexNet, and MobileNetV2 achieve Accuracy and F1-scores over 99.00%. AlexNet constitutes the second best performing model obtaining an Accuracy of 99.37% and a F1-score of 99.37%. EfficientNet-B4 achieves the worst evaluation results reaching Accuracy and F1-score up to 97.14% and 97.02% respectively.

TABLE II: Performance comparison among proposed models on the Cresci’17 dataset (using only images based on the type of the tweet). Reported values are mean \pm standard deviation. Results are averaged across five runs. Best results per evaluation metric are in bold.

Architecture	Evaluation metrics				
	Precision	Recall	F1-score	Accuracy	Specificity
digital DNA (type of tweets)					
<i>GoogLeNet (Inception v1)</i>	100.00 ± 0.00	97.63 ± 1.39	98.79 ± 0.72	98.73 ± 0.79	100.00 ± 0.00
<i>ResNet50</i>	100.00 ± 0.00	97.83 ± 1.00	98.90 ± 0.52	98.94 ± 0.47	100.00 ± 0.00
<i>WideResNet-50-2</i>	99.78 ± 0.44	99.35 ± 1.29	99.57 ± 0.87	99.58 ± 0.85	99.79 ± 0.42
<i>AlexNet</i>	99.59 ± 0.81	97.12 ± 1.23	98.34 ± 0.88	98.31 ± 0.91	99.55 ± 0.91
<i>SqueezeNet1_0</i>	99.17 ± 1.02	97.23 ± 1.41	98.17 ± 0.55	98.20 ± 0.54	99.14 ± 1.05
<i>DenseNet-201</i>	99.37 ± 0.51	98.59 ± 0.98	98.98 ± 0.55	98.94 ± 0.58	99.35 ± 0.53
<i>MobileNetV2</i>	99.58 ± 0.52	97.70 ± 1.28	98.63 ± 0.76	98.62 ± 0.72	99.57 ± 0.53
<i>ResNeXt-50 32 \times 4d</i>	99.79 ± 0.43	97.59 ± 1.64	98.67 ± 1.01	98.62 ± 1.04	99.78 ± 0.44
<i>VGG16</i>	99.78 ± 0.44	99.55 ± 0.54	99.67 ± 0.44	99.68 ± 0.42	99.80 ± 0.41
<i>EfficientNet-B4</i>	99.20 ± 1.14	96.31 ± 0.60	97.73 ± 0.47	97.67 ± 0.54	99.08 ± 1.35

b) *Twibot-20 Dataset*: The results of our proposed approaches are reported in Tables IV and V.

TABLE III: Performance comparison among proposed models on the Cresci’17 dataset (using only images based on the content of the tweet). Reported values are mean \pm standard deviation. Results are averaged across five runs. Best results per evaluation metric are in bold.

Architecture	Evaluation metrics				
	Precision	Recall	F1-score	Accuracy	Specificity
digital DNA (content of tweets)					
<i>GoogLeNet (Inception v1)</i>	100.00 ± 0.00	98.09 ± 1.27	99.03 ± 0.65	99.05 ± 0.62	100.00 ± 0.00
<i>ResNet50</i>	100.00 ± 0.00	97.81 ± 1.13	98.89 ± 0.58	98.84 ± 0.62	100.00 ± 0.00
<i>WideResNet-50-2</i>	98.93 ± 0.71	99.34 ± 0.89	99.14 ± 0.69	99.15 ± 0.63	98.95 ± 0.64
<i>AlexNet</i>	100.00 ± 0.00	98.75 ± 1.20	99.37 ± 0.61	99.37 ± 0.62	100.00 ± 0.00
<i>SqueezeNet1_0</i>	99.36 ± 0.53	97.64 ± 1.06	98.49 ± 0.64	98.52 ± 0.62	99.36 ± 0.52
<i>DenseNet-201</i>	99.59 ± 0.49	98.98 ± 1.37	99.28 ± 0.58	99.26 ± 0.54	99.56 ± 0.54
<i>MobileNetV2</i>	99.59 ± 0.50	98.98 ± 1.27	99.28 ± 0.74	99.26 ± 0.79	99.56 ± 0.54
<i>ResNeXt-50 32 \times 4d</i>	99.35 ± 0.53	98.43 ± 1.31	98.88 ± 0.52	98.94 ± 0.47	99.38 ± 0.51
<i>VGG16</i>	100.00 ± 0.00	99.78 ± 0.43	99.89 ± 0.22	99.89 ± 0.21	100.00 ± 0.00
<i>EfficientNet-B4</i>	96.13 ± 2.26	97.98 ± 2.01	97.02 ± 1.38	97.14 ± 1.28	96.36 ± 2.14

As one can easily observe in Table IV, AlexNet constitutes the best performing model in terms of F1-score and Accuracy. Specifically, it outperforms the other models in F1-score by 0.48-1.34% and in Accuracy by 0.06-0.78%. Although other models outperform AlexNet in Precision and Recall, F1-score is a more important metric, since it is a weighted average of precision and recall. DenseNet-201 and MobileNetV2 achieve almost equal Accuracy scores accounting for 66.35% and 66.34% respectively. The second highest F1-score is achieved by SqueezeNet1_0 and is equal to 66.63%. GoogLeNet obtains the worst performance reaching Accuracy and F1-score up to 65.63% and 65.71% respectively.

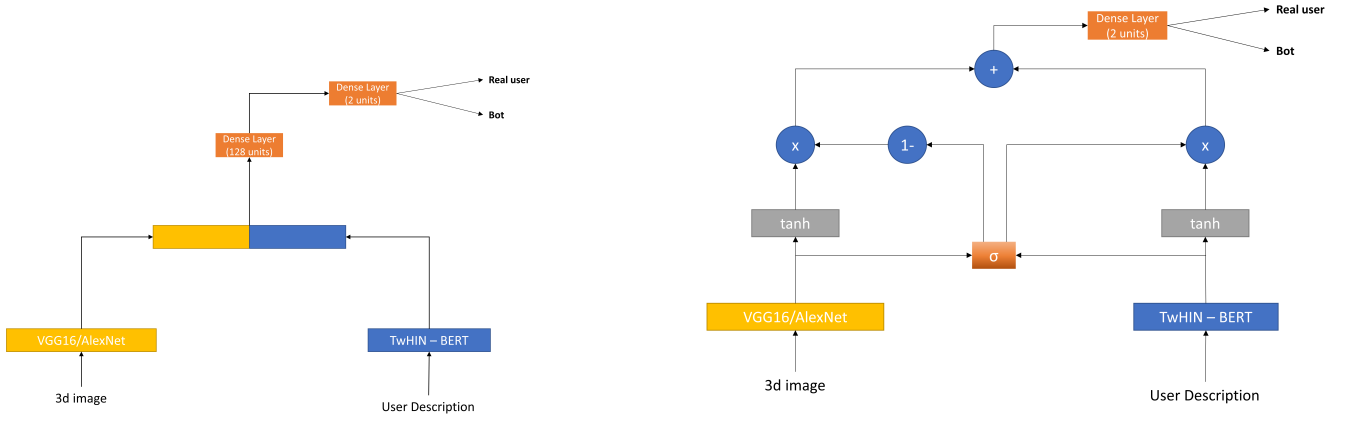
As one can easily observe in Table V, AlexNet is the best performing model outperforming the rest of the models in Precision by 2.11-3.16%, in F1-score by 0.78-3.06%, in Accuracy by 1.60-2.50%, and in Specificity by 1.13-4.58%. GoogLeNet, SqueezeNet1_0, MobileNetV2, VGG16, and EfficientNet-B4 obtain almost equal Accuracy scores ranging from 65.45% to 65.59%, with MobileNetV2 outperforming these models in terms of both Accuracy and F1-score. ResNet50 obtains the worst performance reaching Accuracy and F1-score up to 64.69% and 63.74% respectively.

V. OUR PROPOSED MULTIMODAL MODELS FOR DETECTING BOTS IN TWITTER

In this section, we describe our proposed models employing both textual and vision modalities. Our introduced models are illustrated in Fig. 1.

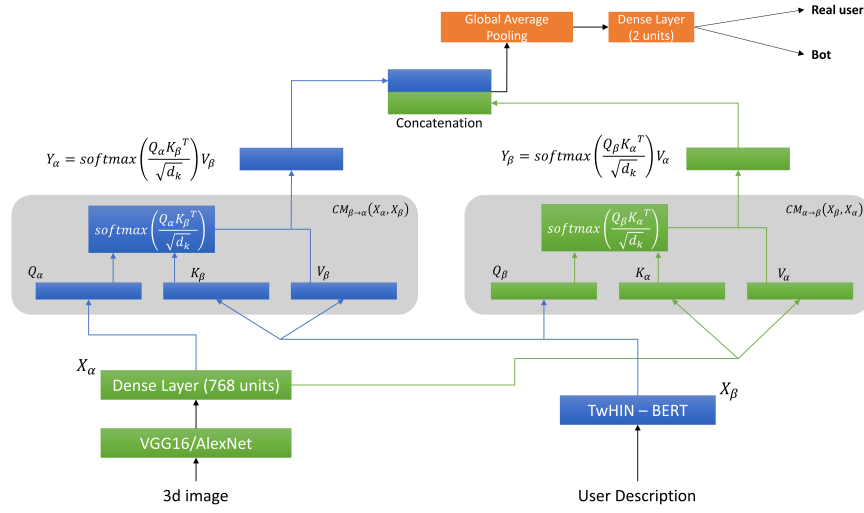
A. Concatenation

In terms of the textual modality, we pass the user description field through a TwHIN-BERT model [16] and extract the [CLS] token. Let $f^t \in \mathbb{R}^{d_t}$ denote the representation vector of the textual modality. d_t denotes the dimensionality and is equal to 768. In terms of the visual modality, we create a 3d image as described in Section IV. Then, we exploit our best performing model, namely VGG16 for Cresci’17 dataset and AlexNet for



(a) Concatenation

(b) Gated Multimodal Unit



(c) Crossmodal Attention

Fig. 1: Our introduced approaches

TABLE IV: Performance comparison among proposed models on the TwiBot-20 dataset (using only images based on the type of the tweet). Reported values are mean \pm standard deviation. Results are averaged across five runs. Best results per evaluation metric are in bold.

Architecture	Evaluation metrics				
	Precision	Recall	F1-score	Accuracy	Specificity
digital DNA (type of tweets)					
<i>GoogLeNet (Inception v1)</i>	68.12 ± 0.61	63.53 ± 2.58	65.71 ± 1.29	65.63 ± 0.67	67.89 ± 1.84
<i>ResNet50</i>	67.90 ± 1.44	65.29 ± 4.30	66.46 ± 1.90	65.90 ± 0.94	66.56 ± 3.98
<i>WideResNet-50-2</i>	67.53 ± 0.73	65.63 ± 2.85	66.52 ± 1.30	65.76 ± 0.63	65.91 ± 2.31
<i>AlexNet</i>	68.33 ± 1.26	65.89 ± 2.05	67.05 ± 0.72	66.41 ± 0.62	66.96 ± 2.79
<i>SqueezeNet1_0</i>	67.52 ± 2.13	66.19 ± 5.54	66.63 ± 1.72	65.74 ± 0.26	65.26 ± 6.27
<i>DenseNet-201</i>	68.86 ± 1.03	64.32 ± 3.10	66.45 ± 1.31	66.35 ± 0.52	68.54 ± 2.83
<i>MobileNetV2</i>	69.04 ± 0.79	64.09 ± 2.26	66.44 ± 0.91	66.34 ± 0.37	68.95 ± 2.20
<i>ResNeXt-50 32 \times 4d</i>	68.24 ± 1.15	64.95 ± 3.13	66.49 ± 1.22	66.08 ± 0.43	67.29 ± 3.14
<i>VGG16</i>	67.85 ± 2.35	65.78 ± 5.69	66.57 ± 1.81	65.86 ± 0.67	65.95 ± 6.68
<i>EfficientNet-B4</i>	67.87 ± 1.44	64.65 ± 4.17	66.11 ± 1.54	65.69 ± 0.36	66.80 ± 4.25

TABLE V: Performance comparison among proposed models on the TwiBot-20 dataset (using only images based on the content of the tweet). Reported values are mean \pm standard deviation. Results are averaged across five runs. Best results per evaluation metric are in bold.

Architecture	Evaluation metrics				
	Precision	Recall	F1-score	Accuracy	Specificity
digital DNA (content of tweets)					
<i>GoogLeNet (Inception v1)</i>	67.53 ± 0.78	64.54 ± 2.12	65.97 ± 0.94	65.47 ± 0.50	66.48 ± 2.07
<i>ResNet50</i>	68.25 ± 0.91	59.89 ± 2.86	63.74 ± 1.39	64.69 ± 0.52	69.88 ± 2.54
<i>WideResNet-50-2</i>	67.57 ± 2.19	62.55 ± 4.80	64.78 ± 1.85	64.83 ± 0.81	67.29 ± 5.59
<i>AlexNet</i>	70.36 ± 1.11	63.64 ± 1.53	66.80 ± 0.37	67.19 ± 0.28	71.01 ± 2.20
<i>SqueezeNet1_0</i>	68.08 ± 1.36	63.38 ± 4.16	65.52 ± 1.66	65.49 ± 0.31	67.77 ± 4.01
<i>DenseNet-201</i>	67.20 ± 1.02	63.60 ± 2.39	65.31 ± 0.96	64.97 ± 0.50	66.43 ± 2.70
<i>MobileNetV2</i>	67.60 ± 0.46	64.47 ± 2.38	66.02 ± 1.14	65.59 ± 0.52	66.80 ± 1.71
<i>ResNeXt-50 32 \times 4d</i>	67.84 ± 2.00	61.65 ± 5.19	64.40 ± 1.77	64.77 ± 0.40	68.14 ± 5.86
<i>VGG16</i>	68.13 ± 1.61	63.19 ± 4.95	65.41 ± 2.06	65.45 ± 0.64	67.89 ± 4.74
<i>EfficientNet-B4</i>	67.56 ± 0.90	64.50 ± 2.09	65.96 ± 0.66	65.47 ± 0.13	66.52 ± 2.46

Twibot-20 dataset, as mentioned in Section IV-A3. Specifically, we remove the last layer of VGG16 (and AlexNet) and replace the next-to-last dense layer consisting of 4096 units with one dense layer consisting of 768 units. Let $f^v \in \mathbb{R}^{d_v}$ represent the visual representation vector. d_v is equal to 768.

After having calculated f^t and f^v , we concatenate these two representation vectors into a single vector as following:

$$z = [f^t, f^v] \quad (1)$$

, where $z \in \mathbb{R}^d$. $d = d_t + d_v$ and is equal to 1536.

After this, we pass z through a dense layer consisting of 128 units with a ReLU activation function. Finally, we use a dense layer consisting of two units, which gives the final prediction.

Our introduced model is illustrated in Fig. 1a.

B. Gated Multimodal Unit

In terms of the textual modality, we pass the user description field through a TwHIN-BERT model and extract the [CLS] token. Let $f^t \in \mathbb{R}^{d_t}$ denote the representation vector of the textual modality. d_t denotes the dimensionality and is equal to 768. In terms of the visual modality, we create a 3d image as described in Section IV. Then, we exploit our best performing model, namely VGG16 for Cresci'17 dataset and AlexNet for Twibot-20 dataset, as mentioned in Section IV-A3. Specifically, we remove the last layer of VGG16 (and AlexNet) and replace the next-to-last dense layer consisting of 4096 units with one dense layer consisting of 768 units. Let $f^v \in \mathbb{R}^{d_v}$ represent the visual representation vector. d_v is equal to 768.

Next, we use a gated multimodal unit introduced in [56] for controlling the information flow of the textual and visual modalities. The equations governing the gated multimodal unit are presented below:

$$h^t = \tanh(W^t f^t + b^t) \quad (2)$$

$$h^v = \tanh(W^v f^v + b^v) \quad (3)$$

$$z = \sigma(W^z [f^t; f^v] + b^z) \quad (4)$$

$$h = z * h^t + (1 - z) * h^v \quad (5)$$

$$\Theta = \{W^t, W^v, W^z\} \quad (6)$$

, where h is a weighted combination of the textual and visual information h^t and h^v respectively. Θ denotes the parameters to be learned, while $[\cdot; \cdot]$ indicates the concatenation operation.

Finally, we pass h through a dense layer consisting of two units for getting the final prediction.

Our introduced model is illustrated in Fig. 1b.

C. Crossmodal Attention

In terms of the textual modality, we pass the user description field through a TwHIN-BERT model. Let $f^t \in \mathbb{R}^{N \times d_t}$ denote the textual representation. N indicates the sequence length, while d_t denotes the dimensionality and is equal to 768. In terms of the visual modality, we create a 3d image as described in Section IV. Next:

- In terms of the Cresci'17 dataset, we exploit our best performing model, namely VGG16, as mentioned in Section IV-A3. Specifically, we get the output of the last CNN layer (after max pooling) as the output of the

pretrained VGG16 model. Let $f^v \in \mathbb{R}^{T \times d_v}$ represent the visual representation vector. T and d_v are equal to 64 and 512 respectively. Next, we pass f^v through a dense layer consisting of 768 units.

- In terms of the Twibot-20 dataset, we exploit our best performing model, namely AlexNet, as mentioned in Section IV-A3. Specifically, we get the output of the last CNN layer (after max pooling) as the output of the pretrained AlexNet model. Let $f^v \in \mathbb{R}^{T \times d_v}$ represent the visual representation vector. T and d_v are equal to 49 and 256 respectively. Next, we pass f^v through a dense layer consisting of 768 units.

Motivated by [57]–[59], we exploit two crossmodal attention layers, i.e., one from textual f^t to visual features f^v and one from visual to textual features.

Specifically, we calculate the scaled dot attention [60] as follows:

$$z = \text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \quad (7)$$

, where the textual modality corresponds to the query (Q), while the visual modality corresponds to the key (K) and value (V).

Similarly, we design another scaled dot attention layer, where the visual modality corresponds to the query (Q), while the textual modality corresponds to the key (K) and value (V). Let y be the output of the scaled dot attention layer.

Next, we concatenate z and y and pass the resulting matrix into a global average pooling layer.

Finally, we exploit a dense layer consisting of two units for getting the final output.

Our introduced model is illustrated in Fig. 1c.

VI. EXPERIMENTS

A. Baselines

- Comparison with state-of-the-art approaches
 - Cresci'17 Dataset
 - * DeepSBD [27]: This method extracts profile, temporal, activity, and content information. Finally, a deep neural network is trained consisting of LSTMs, CNNs, and an attention layer.
 - * DNNBD [10]: This method uses only profile information along with SMOTE. We report the results reported in [27].
 - * DBDM [28]: This method models social behavior and content information. A deep neural network consisting of CNNs and BiLSTMs is trained. We report the results reported in [27].
 - * DeBD [61]: This method passes the tweet content and the relationship between them into a CNN. Secondly, it uses LSTM to extract the potential temporal features of the tweet metadata. Finally, the authors concatenate the temporal features with the joint content features for detecting social bots. We report the results reported in [27].
 - * MulBot-Glob_Hier [37]: This method extracts multidimensional temporal features from user's

timeline, employs dimensionality reduction algorithms, extracts global features, and performs the Agglomerative Hierarchical clustering algorithm.

- * DNA - sequence (supervised) [14]: This method extracts the digital DNA sequence per user by using either the type or content of tweets. Then, the authors leverage the longest common substring (LCS) curves for detecting the social spambots.
 - * Ahmed_DBSCAN [37], [62]: This method adopts the approach proposed by [62] and uses DBSCAN to increase the performance. We report the results reported in [37].
 - * Ahmed and Abulaish [62]: This approach exploits the Euclidean distance between feature vectors to create a similarity graph of the accounts. Next, graph clustering and community detection algorithms are used for identifying groups of similar accounts in the graph. We report the results reported in [14].
 - * Botometer [18]: Botometer is a publicly available bot detection tool via the website, Python, or REST APIs. This approach is trained with more than 1,000 features using a Random Forest classifier.
- TwiBot-20 dataset
- * Kudugunta and Ferrara [10]: This method utilizes user metadata features and tweet content features for identifying bots.
 - * Wei and Nguyen [31]: This method uses GloVe embeddings and trains a deep neural network consisting of three layers of BiLSTMs.
 - * Miller et al. [38]: Bot detection task is considered as an anomaly detection problem. Specifically, this method extracts 107 features and modifies two stream clustering algorithms, namely StreamKM++ and DenStream.
 - * Cresci et al. [13]: This method extracts the digital DNA sequence per user by using either the type or content of tweets. Then, the authors leverage the longest common substring (LCS) curves for detecting bots in groups.
 - * Botometer [18].
 - * Alhosseini et al. [36]: This method extracts a set of features per user, including age, favourites_count, statuses_count, friends_count, followers_count, and account length name. Then, a Graph Convolutional Neural Network is trained.

- TwHIN-BERT using only the user description field
- Our best performing model described in Section IV

B. Experimental Setup

All the details are provided in Section IV-A1. We use the TwHIN-BERT-base version from the Transformers library in Python [63].

C. Evaluation Metrics

We use the evaluation metrics, which are described in Section IV-A2.

VII. RESULTS

The results of our proposed approaches are reported in Tables VI and VII. Specifically, Table VI reports the results on the Cresci’17 dataset, while Table VII reports the results on the TwiBot-20 dataset.

A. Cresci’17 Dataset

The results of our proposed approaches described in Section V are reported in Table VI.

TABLE VI: Performance comparison among proposed models and state-of-the-art approaches on the Cresci’17 dataset. Reported values are mean \pm standard deviation. Results are averaged across five runs. Best results per evaluation metric are in bold.

Architecture	Evaluation metrics				
	Precision	Recall	F1-score	Accuracy	Specificity
Comparison with state-of-the-art approaches					
DeepSBD [27]	100.00	-	99.81	99.83	-
DNNBD [10]	77.66	-	75.63	78.20	-
DBDM [28]	100.00	-	98.82	99.32	-
DeBD [61]	97.73	-	97.59	97.74	-
MulBot-Glob_Hier [37]	99.50	99.50	99.00	99.30	-
DNA - sequence (supervised) [14]	98.20	97.70	97.70	97.70	98.10
Ahmed_DBSCAN [37], [62]	93.00	93.00	93.00	92.80	-
Ahmed and Abulaish [62]	94.50	94.40	94.40	94.30	94.50
Botometer [18]	-	-	97.31	95.97	-
Unimodal approaches (only user description)					
TwHIN-BERT	99.59 ± 0.50	99.18 ± 0.75	99.38 ± 0.39	99.37 ± 0.40	99.56 ± 0.54
Unimodal approaches (only images)					
VGG16 (type of tweets)	99.78 ± 0.44	99.55 ± 0.54	99.67 ± 0.44	99.68 ± 0.42	99.80 ± 0.41
VGG16 (content of tweets)	100.00 ± 0.00	99.78 ± 0.43	99.89 ± 0.22	99.89 ± 0.21	100.00 ± 0.00
Proposed Transformer-based models (images based on the type of tweets)					
TwHIN-BERT + VGG16 (Concatenation)	99.78 ± 0.44	99.34 ± 0.88	99.56 ± 0.54	99.58 ± 0.52	99.80 ± 0.41
TwHIN-BERT + VGG16 (GMU)	99.79 ± 0.42	99.79 ± 0.42	99.79 ± 0.42	99.79 ± 0.42	99.78 ± 0.43
TwHIN-BERT + VGG16 (Cross-Modal Attention)	100.00 ± 0.00	99.79 ± 0.41	99.90 ± 0.21	99.89 ± 0.21	100.00 ± 0.00
Proposed Transformer-based models (images based on the content of tweets)					
TwHIN-BERT + VGG16 (Concatenation)	99.77 ± 0.46	99.77 ± 0.46	99.77 ± 0.46	99.79 ± 0.42	99.80 ± 0.39
TwHIN-BERT + VGG16 (GMU)	100.00 ± 0.00	99.58 ± 0.52	99.79 ± 0.26	99.79 ± 0.26	100.00 ± 0.00
TwHIN-BERT + VGG16 (Cross-Modal Attention)	100.00 ± 0.00	99.96 ± 0.08	99.98 ± 0.04	99.98 ± 0.04	100.00 ± 0.00

Regarding our proposed approaches (transformer-based models based on the type of tweets), one can observe that TwHIN-BERT + VGG16 (Cross-Modal Attention) constitutes the best performing model achieving an Accuracy of 99.89%, a Precision of 100.00%, a Recall of 99.79%, a F1-score of 99.90%, and a Specificity of 100.00%. It outperforms the other introduced models in Accuracy by 0.10-0.31%, in Precision by 0.21-0.22%, in F1-score by 0.11-0.34%, and in Specificity by 0.20-0.22%. Also, it outperforms both TwHIN-BERT and VGG16 (type of tweet). Specifically, it surpasses TwHIN-BERT in Precision by 0.41%, in Recall by 0.61%, in F1-score by 0.52%, in Accuracy by 0.52%, and in Specificity by 0.44%. It outperforms also VGG16 (type of tweet) in Precision by 0.22%, in Recall by 0.24%, in F1-score by 0.23%, in Accuracy by 0.21%, and in Specificity by 0.20%. Additionally, TwHIN-BERT + VGG16 (GMU) outperforms TwHIN-BERT + VGG16 (Concatenation) in Precision by 0.01%, in Recall by 0.45%, in F1-score by 0.23%, and in Accuracy by 0.22%. Also, it outperforms both TwHIN-BERT and VGG16 (type of tweets) in Accuracy by 0.42% and 0.11% respectively. In terms of TwHIN-BERT + VGG16 (Concatenation), one can observe

that this model obtains worse performance than VGG16 (type of tweets). Specifically, VGG16 (type of tweets) outperforms TwHIN-BERT (Concatenation) in Recall by 0.21%, in F1-score by 0.11%, and in Accuracy by 0.10%. We speculate that this decrease in performance is attributable to the concatenation operation, which assigns equal importance to each modality ignoring the inherent correlations between the two modalities. Overall, we believe that TwHIN-BERT + VGG16 (Cross-Modal Attention) obtains better performance than the other two multimodal models, since it captures the crossmodal interactions. On the other hand, the Gated Multimodal Unit controls the information flow from each modality, while the concatenation operation neglects the inherent correlations.

In terms of our proposed transformer-based models (images based on the content of tweets), one can observe that TwHIN-BERT + VGG16 (Cross-Modal Attention) constitutes the best performing model outperforming both TwHIN-BERT + VGG16 (GMU) and TwHIN-BERT + VGG16 (Concatenation). Specifically, TwHIN-BERT + VGG16 (Cross-Modal Attention) outperforms TwHIN-BERT + VGG16 (GMU) in Recall by 0.42%, in F1-score by 0.19%, and in Accuracy by 0.19%. Although equal Precision scores are achieved, TwHIN-BERT + VGG16 (Cross-Modal Attention) yields a better F1-score, which constitutes the weighted average of Precision and Recall. Additionally, TwHIN-BERT + VGG16 (Cross-Modal Attention) outperforms TwHIN-BERT + VGG16 (Concatenation) in Precision by 0.23%, in Recall by 0.19%, in F1-score by 0.21%, in Accuracy by 0.19%, and in Specificity by 0.20%. In comparison with unimodal approaches employing either text or images, we observe that TwHIN-BERT + VGG16 (Crossmodal Attention) outperforms both TwHIN-BERT and VGG16 (content of tweets). Finally, TwHIN-BERT + VGG16 (Cross-Modal Attention) with images based on the content of tweets outperforms all the introduced models exploiting images based on the type of tweets. Therefore, TwHIN-BERT + VGG16 (Cross-Modal Attention) with images based on the content of tweets constitutes our best performing model.

In comparison with the state-of-the-art approaches, one can observe that our best performing model, namely TwHIN-BERT + VGG16 (Cross-Modal Attention) with images based on the content of the tweets, outperforms these approaches in Precision by 0.50-22.34% (except for DeepSBD and DBDM), in Recall by 0.46-6.96%, in F1-score by 0.17-24.35%, in Accuracy by 0.15-21.78%, and in Specificity by 1.90-5.50%. Similarly, TwHIN-BERT + VGG16 (Cross-Modal Attention) with images based on the type of tweets, surpasses the existing research initiatives in Precision by 0.50-22.34% (except for DeepSBD and DBDM), in Recall by 0.29-6.79%, in F1-score by 0.09-24.27%, in Accuracy by 0.06-21.69%, and in Specificity by 1.90-5.50%. Although our best performing model outperforms DeepSBD by a small margin of 0.17% and 0.15% in F1-score and Accuracy respectively, our best performing model has multiple advantages over DeepSBD. Firstly, DeepSBD processes the tweets of each user, extracts GloVe embeddings, and creates a 3d matrix which is given as input to a 6-layer CNN increasing the computational demands. Additionally, this method extracts a set of features per user. On the contrary, our method seems to be simpler and more effective, since it does

not rely on a feature extraction strategy; it relies only on the user description and the sequence of actions performed by the user. Additionally, the authors use GloVe embeddings instead of using a language model based on transformers. Thus, this approach inherits the limitations of the GloVe embeddings. In terms of the training time per epoch, 120 to 180 seconds are required for training by our model, while the training time of DeepSBD ranges from 248 to 634 seconds.

B. TwiBot-20 Dataset

The results of our proposed approaches described in Section V are reported in Table VII.

In terms of the proposed transformer-based models using images based on the type of tweets posted by the users, one can observe that TwHIN - BERT + AlexNet (Cross-Modal Attention) constitutes the best performing model in terms of Recall, F1-score, and Accuracy. Specifically, it outperforms the other introduced multimodal models in Recall by 2.55-2.85%, in F1-score by 0.83-1.26%, and in Accuracy by 0.28-0.76%. Additionally, TwHIN-BERT + AlexNet (Cross-Modal Attention) with images based on the type of tweets, surpasses the performance achieved by unimodal models, i.e., TwHIN-BERT (using user description) and AlexNet (using images). Similar differences in models' performance are observed, when the images based on the content of tweets are used as input to proposed approaches. To be more precise, the usage of the crossmodal attention as a fusion method yields the best results surpassing the performance achieved by concatenation and gated multimodal unit as fusion approaches. TwHIN-BERT + AlexNet (Cross-modal Attention) outperforms the other multimodal models in Recall, F1-score, and Accuracy by 3.16-4.43%, 1.26-2.03%, and 0.70-1.36% respectively. TwHIN-BERT + AlexNet (Cross-Modal Attention) outperforms also the unimodal models employing either text or images. Specifically, it surpasses TwHIN-BERT in F1-score and Accuracy by 3.18% and 3.25% respectively, while it also outperforms AlexNet in F1-score and Accuracy by 9.50% and 7.47% respectively. Overall, we observe that the usage of multiple modalities improves bots' detection performance. Cross-modal Attention boosts performance, since it models the cross-modal interactions. On the other hand, the gated multimodal unit refers to a weighting strategy, which controls the contribution of each modality to the bots' prediction task. The concatenation operation is not capable of capturing the interactions of the two modalities or controlling the contribution of each modality, since equal importance to text and image information is assigned.

In comparison with state-of-the-art approaches, our best performing model outperforms these approaches in Accuracy and F1-score by 3.40-26.73% and 0.97-65.58% respectively. We observe that although Botometer was capable of identifying the social spambots in Cresci'17 dataset, the performance drops significantly in TwiBot-20 dataset, which verifies the fact that bots in Twitter have changed their behavior and have managed to evade previous detection methods. Similarly, we observe that the other approaches employing feature extraction techniques fail to detect bots. On the contrary, our approach relying only

on the user description and the sequence of actions performed by the user based on the content of tweets seems to be an effective method.

TABLE VII: Performance comparison among proposed models and state-of-the-art approaches on the TwiBot-20 dataset. Reported values are mean \pm standard deviation. Results are averaged across five runs. Best results per evaluation metric are in bold.

Architecture	Evaluation metrics				
	Precision	Recall	F1-score	Accuracy	
Comparison with state-of-the-art approaches					
Kudugunta and Ferrara [10]	-	-	47.26	59.59	-
Wei and Nguyen [31]	-	-	75.33	71.26	-
Miller et al. [38]	-	-	62.66	48.01	-
Cresci et al. [13]	-	-	10.72	47.93	-
Botometer [18]	-	-	48.92	55.84	-
Alhosseini et al. [36]	-	-	73.18	68.13	-
Unimodal approaches (only user description)					
<i>TwHIN-BERT</i>	71.41	75.12	73.12	71.41	67.41
	± 1.64	± 4.15	± 1.61	± 1.07	± 3.76
Unimodal approaches (only images)					
<i>AlexNet (type of tweets)</i>	68.33	65.89	67.05	66.41	66.96
	± 1.26	± 2.05	± 0.72	± 0.62	± 2.79
<i>AlexNet (content of tweets)</i>	70.36	63.64	66.80	67.19	71.01
	± 1.11	± 1.53	± 0.37	± 0.28	± 2.20
Proposed Transformer-based models (images based on the type of tweets)					
<i>TwHIN-BERT + AlexNet</i>	74.22	75.01	74.52	73.42	71.70
(Concatenation)	± 1.93	± 3.35	± 0.88	± 0.53	± 3.98
<i>TwHIN-BERT + AlexNet</i>	74.70	75.31	74.95	73.90	72.39
(GMU)	± 1.56	± 3.02	± 1.26	± 1.00	± 2.91
<i>TwHIN-BERT + AlexNet</i>	73.98	77.86	75.78	74.18	70.20
(Cross-Modal Attention)	± 2.40	± 3.02	± 0.48	± 0.78	± 4.59
Proposed Transformer-based models (images based on the content of the tweets)					
<i>TwHIN-BERT + AlexNet</i>	74.31	74.26	74.27	73.30	72.27
(Concatenation)	± 0.90	± 1.39	± 0.35	± 0.29	± 1.78
<i>TwHIN-BERT + AlexNet</i>	74.74	75.53	75.04	73.96	72.27
(GMU)	± 1.88	± 3.58	± 0.93	± 0.31	± 3.93
<i>TwHIN-BERT + AlexNet</i>	74.16	78.69	76.30	74.66	71.61
(Cross-Modal Attention)	± 1.14	± 2.97	± 0.83	± 0.31	± 3.16

VIII. QUALITATIVE AND ERROR ANALYSIS

In this section, we perform a qualitative and error analysis of the best performing model on the TwiBot-20 dataset, namely TwHIN-BERT + AlexNet (Cross-Modal Attention) using images based on the content of tweets. Table VIII reports the predictions on some sample instances. Specifically, rows 1 and 2 refer to correct classifications as real users. Similarly, rows 3-7 refer to correct predictions as bots. Rows 8-10 refer to incorrect predictions made by our best performing model. Specifically, our model predicts these instances as belonging to a bot, while the real label corresponds to a genuine account. Similarly, rows 11 and 12 indicate incorrect classifications, where the actual label is a bot, while the predicted label is a real user.

Regarding correct classifications corresponding to genuine accounts, we observe that the Digital DNA of a real user includes a great number of consecutive mentions. Thus, our model learns that consecutive mentions indicate a real user, since mentions indicate often interactions with other users. Additionally, we observe that the user description field refers to a simple description of the biography of the user.

In terms of the correct classifications referring to bots, we observe in row 3 that the digital DNA includes a lot of consecutive hashtags. Although Digital DNAs of rows 4-7 include a lot of mentions, we observe that the user description fields have common characteristics. Specifically, we observe that the users' descriptions of rows 4-6 include many words separated by commas (,). The user description of row 7 includes

many hashtags (8 in number) and emojis. Therefore, both user description and the activity of users are important towards the final prediction.

Next, we examine reasons of misclassifications, i.e., rows 8-10. Specifically, we hypothesize that our model misclassifies row 8 as bot, since the Digital DNA consists of a great number of consecutive URLs. Also, we believe that row 9 is predicted as bot, since the user description includes many mentions. Finally, we observe that row 10 is similar to rows 4-6, and thus this row is predicted falsely as bot. Specifically, we observe that the user description of row 10 includes many words separated by commas.

Finally, we investigate the reasons of misclassifications in rows 11 and 12. We observe that Digital DNAs contain a lot of consecutive mentions, which refer to a genuine account (see rows 1 and 2). This verifies the fact that bots have found ways to evade detection approaches and mimic human behavior.

IX. DISCUSSION

A. Limitations

This study has some limitations. Firstly, we did not apply explainability approaches to explain the predictions made by our proposed approaches. Additionally, this study requires labelled data. On the contrary, self-supervised learning approaches address the issue of labels' scarcity.

B. Scalability and Real-time detection

An often neglected concern among researchers introducing methods for identifying bots in Twitter is the reliance on the platform APIs [64]. Specifically, data must be retrieved from social media platforms, so as to ensure that the researchers perform the necessary experiments.

However, in 2023, Twitter decided to end free access to their APIs. As a result of this modification, several research projects³ have been cancelled or suspended for a period of time. After that, legal obligations in the bloc's Digital Services Act (DSA)⁴ require larger platforms, including Meta, Twitter, and Google, to provide data access to external researchers doing public interest research on systemic risks. However, the access to Twitter APIs still remains uncertain.

Modifications also to the API affect the bot detection algorithms. For example, the field `geo_enabled` was removed in 2019 due to privacy reasons. Therefore, algorithms trained with this features, should be modified.

Data access is inextricably linked with scalability. Scalability refers to the analysis of streaming data with limited computing resources [65]. Specifically, the speed at which a method processes a group of accounts depends highly on the possibility for access to the Twitter API. For instance, although the construction of large graphs leads to detection of bots with increasing evaluation performance, fetching such information is not feasible. The study in [66] claims that the following

³<https://www.reuters.com/technology/elon-musks-x-restructuring-curtails-disinformation-research-spurs-legal-fears-2023-11-06/>

⁴<https://techcrunch.com/2023/11/17/change-in-xs-terms-indicate-eu-researchers-will-get-api-access/>

TABLE VIII: Sample Instances for Analysis

ID	User Description	Digital DNA	Real	Pred.
1	Sports Broadcaster & Event Host @skysports	MXMMMNUUNUMXUMUMMMUUMXXXXUMMMMMMMMMXXUXMMMMXU NUXXMMUMUMXUUXXXMMUUXUXXUMUXUMMMUXMXXUXUMM- MUMX UUXUMMUXXXMMXMUUUMMXUUMMMMMXXUUMMXUMM- MUMMUXXMMUMUXXUUMMUMU UMMUXXXUUNUUMXMM- MXMXXUMXMUUUXMMMMMUUMXMUMX	Real User	Real User
2	Former NASA and Apple engineer. Current YouTuber and friend of science. <URL>	XMUMMMXMXMMXXMXXMXXUMMXMUMXUUMMMMMXMMMMMM MMMMXMMMMMUUUMMXMMMXMMXXUNXMMUMMMMMMM- MMUMUXMMNXUUMXUMXMUMM MMMMMXMMUUMMUMM- MMMMMUUMMMUMMMMMMMMNMMNUUMMXMMNMMXMM MMXMMXUMMMMMMUUNXUMMMMMUMMMUUXMXX	Real User	Real User
3	Fun, outgoing guy that loves to meet new people. Casual gamer with a competitive nature.	HXXXXXXXXMXXMXXMXXHXXHHHHMXXXXXXXXXXXXXXXXXX MHHXXUXMUXXXUXXXXUHMHMXXMMMXHHHHHHHHHHXHHH- HXHHHHHHHH HXHHHHHHMXXMMMXHXXHHHHHHHHXHHH- MXXXXHMMHHHHHHHHH HMXHMHNNHHHHHHHHHHHHHH- HUHHHXHHXXHM	Bot	Bot
4	progressive, vegan, animal rights, social justice, anti racism, resistance, anti trump, lock him up, the SDNY will get the whole crime family <raised fist emoji> #FBR	MXMMMMXMMUMMMMMXUMMMMMMMUMUMMMUMXMXM XMMMM- MXMMMMMXMXXMMMMMMMMMMMMMMMMMMMMMUMMMMM MUXMXXU- MUMMMMMXMMMMMMMXMXXMMMMXUMMMM MMXMXUMUMM- MMXMMMMMMMMMXMXXUMMMMMUMXMMXMMMXMMXMXM- MUMXMMUMMMMXMMUMMXUXM XMMMXMXXMM	Bot	Bot
5	wife, mom, gma, Dem, liberal, truth, justice, equality, free press, free speech, LGBTQ, unions, proChoice, singlePayer, science#Resistance, #VoteBluetoSaveAmerica	MXMMMMXMMUMMMMMXUMMMMMMMUMUMMMUMXMXM XMMMM- MMMMMXMXXMMMMMXMMMXMXXMXXMXXMMMMMMMXMM- MMXMMXXMXXMMUUMMMU XMXMXXMXXMMMMXMMMMMM- MUMXMMXMMUMUMMMXMMUMXMMUM MMMXMMMMUMXMM- MXMMMMMMMXMXXMMMMMMMMMXMXXMXX	Bot	Bot
6	Christian, Conservative, Army,Life member NRA,Believe in God my Savior and Lord.Believe in America and believe in our great State of TEXAS !!!Senator Cruz III%	XMMMMXMXMXXMMMMXMMNMMXMMXMMXMMXMXM XMMMM- MMMMMMMMMMMMMMXMMMMMMMMMMXMMNMMUMM XMMMM- MMXMMMMXMMMMUMXMMUMXMMUMMMM XMXMMMM- MXMXXMMMMMMMMMMMMMXMMMMMMMM MXMXXMXXMM- MXMMMMMMMMXMM	Bot	Bot
7	I AM HERE 4 POTUS #GGOAT #MomentsMatter <glowing star emoji> #BELIEVING=PEACE <folded hands emoji> <paw prints emoji> #BeachBabeBum <desert island emoji> <eyes emoji> <fox emoji> #TrumpTrain <locomotive emoji> #MAGA #KAG #RedWomenRealWomen ArkTrumper <anchor emoji> SADLY_Disabled	MMXXMMMXMXXMMMMMMXMXMXXMMMMMMMMMMMMMXM MMMNUXMMXMMXMMMMMMMMMXMXXMXXMXXMXXMXXM- MXXXXXMM XMMMMXMMXXMMMMMXMXXMMXXMMMMMXM- MXMMMXMXXMXXMXXMXXM MMMMXMMMMMMMMMM- MXXMMXMMMMXMMXMXMXXMXXMXXMXXMXXMXXMXXM	Bot	Bot
8	Pasta lover. I don't tweet much. My new Netflix series Master of None is now streaming on Netflix. I wrote a book called Modern Romance.	XXXXXXXXXXUUUUUUUUUUUUUUUUUUUUUUUUUUUUUUUUUU UUUUUUUUUUUUUUUNUXUUUUUUUUUUUUUUUUUUUUUUUU UUUMNXUUU MUUUUUUXUNXXUUUUUUUUUUUUUUUUUUUU UU UUU	Real User	Bot
9	@mu_foundation @veteransgarage & @SSAFA ambassador. @helpforheroes @rmchcharity and @manchesterpride patron.	XMUXXMXUMXUUUMMMNMUMUMMUU MUUMMMNXUMMM- NUM UMUXXUUUMMUUXXNUMMUMXUMXMMXXUNUMMUUU- UMUUUUUM MXNUMXUUUMXUUUXXMXXMXXMUMMMXUM- MMXNXMMMM MXUNMNNMXXUXUNXUMXXUUMMHUUUN- NXXMMMMMXUMXUUNMM	Real User	Bot
10	news/media, music/guitar, ECB, EU, Frankfurt, Greece, politics. Opinions my own, not necessarily my employer's (=ECB). RT's often mean I agree, but not always.	MXMXXMXUMMMUMUMXMXMXXUXXMNUMUXXMUMXMMUMXXX MMXUMXMMNMUMXMMXMMMMMMMMMMXMMUMMMXUNNN- MMXXUMMM XMMUMXMMXMMXMMXMMXMMXMMXMMXMM- MXXMMXMMMMXUMXMX XXMMXXMXXMXXMXXMXXMXXM- MMXNXMXXMXXMMXMMMMXU XXXUU	Real User	Bot
11	Political Booking Producer at @nbc-news @todayshow @nbcnightlynews	MMXNXNXMMMXMMXMMMMXMMMMXMMMMMMMMXMXMXXM XMMXMMXMMMMMMMMXMMXMMXMMMMMMXMMXMM- MMXXMMXUMXMMMMMMMMXMM MMXMXHMMMMMMMM- MMXMMMMXMMXXNXMMXXMMXMMXMMMMMMMMXMM- MXNXMXXMXXMM MXMMXMMXMMXMMMMMMXMMX	Bot	Real User
12	Lets Hope and pray for The Best	MMMMMXMMMMMMMMMMMMMMMMMMMMMMMMMMMMXMMXMMUMM MM	Bot	Real User

methods are required for improving the scalability of bot detection algorithms:

- Model compression and distillation
- Incremental Learning and online algorithms
- Parallel and distributed processing
- Stream-based processing and data reduction

C. Generalization of our approach on other social media

The introduced approaches in this paper can be generalized on other social media, including Facebook and Reddit. Similar to Twitter, users on the Facebook platform make posts, reposts, and reply to other users. At the same time, their posts may include mentions to other users, hashtags, and URLs. Additionally, each user may include a description field in his/her profile.

However, platforms' policy about data access is critical. For instance, investigating bot detection approaches on the Facebook platform is a very difficult task due to the unwillingness of Facebook to share individual account data [64].

X. CONCLUSION AND FUTURE WORK

In this paper, we present the first study introducing multimodal and cross-modal models for detecting bots in Twitter by exploiting only the user description and 3d images, which represent the actions of each user. Firstly, we create two digital DNA sequences based on both the type and content of the tweets each user posts. Next, we apply a DNA-to-image conversion algorithm and create two 3d images per user based on the two digital DNA sequences. Finally, we fine-tune several pretrained models of the vision domain and show that VGG16 achieves the highest evaluation results. Next, we introduce multimodal and crossmodal models. First, we pass each user description field into a TwHIN-BERT and obtain a textual representation. We pass each 3d image through a VGG16 model and obtain a visual representation. Finally, we compare three methods for fusing the textual and visual representations, including the concatenation operation, the gated multimodal unit, and the cross-modal attention. Findings show that the crossmodal attention outperforms the other introduced approaches, be it either unimodal or multimodal approaches, obtaining also comparable performance with the state-of-the-art approaches.

In the future, we aim to apply our models in a federated learning framework. Applying also continual learning approaches is one of our future plans. Finally, we plan to exploit explainability techniques for rendering the introduced approaches explainable.

REFERENCES

- [1] L. Ilias, S. Mouzakitis, and D. Askounis, "Calibration of transformer-based models for identifying stress and depression in social media," *IEEE Transactions on Computational Social Systems*, vol. 11, no. 2, pp. 1979–1990, 2024.
- [2] S. Cresci, "Detecting malicious social bots: story of a never-ending clash," in *Disinformation in Open Online Media: First Multidisciplinary International Symposium, MISDOOM 2019, Hamburg, Germany, February 27–March 1, 2019, Revised Selected Papers 1*. Springer, 2020, pp. 77–88.
- [3] A. Bessi and E. Ferrara, "Social bots distort the 2016 us presidential election online discussion," *First monday*, vol. 21, no. 11-7, 2016.
- [4] F. Brachten, S. Stieglitz, L. Hofeditz, K. Kloppenborg, and A. Reimann, "Strategies and influence of social bots in a 2017 german state election-a case study on twitter," *arXiv preprint arXiv:1710.07562*, 2017.
- [5] M. Shafahi, L. Kempers, and H. Afsarmanesh, "Phishing through social bots on twitter," in *2016 IEEE International Conference on Big Data (Big Data)*, 2016, pp. 3703–3712.
- [6] C. Shao, G. L. Ciampaglia, O. Varol, K.-C. Yang, A. Flammini, and F. Menczer, "The spread of low-credibility content by social bots," *Nature communications*, vol. 9, no. 1, pp. 1–9, 2018.
- [7] J. Uyheng, D. Bellutta, and K. M. Carley, "Bots amplify and redirect hate speech in online discourse about racism during the covid-19 pandemic," *Social Media + Society*, vol. 8, no. 3, p. 20563051221104749, 2022. [Online]. Available: <https://doi.org/10.1177/20563051221104749>
- [8] L. Ilias and I. Roussaki, "Detecting malicious activity in twitter using deep learning techniques," *Applied Soft Computing*, vol. 107, p. 107360, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1568494621002830>
- [9] S. Ouni, F. Fkih, and M. N. Omri, "Bots and gender detection on twitter using stylistic features," in *International Conference on Computational Collective Intelligence*. Springer, 2022, pp. 650–660.
- [10] S. Kudugunta and E. Ferrara, "Deep neural networks for bot detection," *Information Sciences*, vol. 467, pp. 312–322, 2018. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0020025518306248>
- [11] K. Hayawi, S. Mathew, N. Venugopal, M. M. Masud, and P.-H. Ho, "Deepprot: a hybrid deep neural network model for social bot detection based on user profile data," *Social Network Analysis and Mining*, vol. 12, no. 1, p. 43, 2022.
- [12] S. Feng, H. Wan, N. Wang, and M. Luo, "Botrgcn: Twitter bot detection with relational graph convolutional networks," in *Proceedings of the 2021 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, ser. ASONAM '21. New York, NY, USA: Association for Computing Machinery, 2022, p. 236–239. [Online]. Available: <https://doi.org/10.1145/3487351.3488336>
- [13] S. Cresci, R. Di Pietro, M. Petrocchi, A. Spognardi, and M. Tesconi, "Dna-inspired online behavioral modeling and its application to spambot detection," *IEEE Intelligent Systems*, vol. 31, no. 5, pp. 58–64, 2016.
- [14] S. Cresci, R. D. Pietro, M. Petrocchi, A. Spognardi, and M. Tesconi, "Social fingerprinting: Detection of spambot groups through dna-inspired behavioral modeling," *IEEE Transactions on Dependable and Secure Computing*, vol. 15, no. 4, pp. 561–576, 2018.
- [15] E. Di Paolo, M. Petrocchi, and A. Spognardi, "From online behaviours to images: A novel approach to social bot detection," in *Computational Science – ICCS 2023*, J. Mikyška, C. de Mulatier, M. Paszynski, V. V. Krzhizhanovskaya, J. J. Dongarra, and P. M. Sloot, Eds. Cham: Springer Nature Switzerland, 2023, pp. 593–607.
- [16] X. Zhang, Y. Malkov, O. Florez, S. Park, B. McWilliams, J. Han, and A. El-Kishky, "Twhin-bert: A socially-enriched pre-trained language model for multilingual tweet representations," *arXiv preprint arXiv:2209.07562*, 2022.
- [17] S. B. Jr, G. F. C. Campos, G. M. Tavares, R. A. Igawa, M. L. P. Jr, and R. C. Guido, "Detection of human, legitimate bot, and malicious bot in online social networks based on wavelets," *ACM Trans. Multimedia Comput. Commun. Appl.*, vol. 14, no. 1s, mar 2018. [Online]. Available: <https://doi.org/10.1145/3183506>
- [18] C. A. Davis, O. Varol, E. Ferrara, A. Flammini, and F. Menczer, "Botnot: A system to evaluate social bots," in *Proceedings of the 25th International Conference Companion on World Wide Web*, ser. WWW '16 Companion. Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee, 2016, p. 273–274. [Online]. Available: <https://doi.org/10.1145/2872518.2889302>
- [19] E. Arin and M. Kutlu, "Deep learning based social bot detection on twitter," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 1763–1772, 2023.
- [20] S. Najari, M. Salehi, and R. Farahbakhsh, "Ganbot: A gan-based framework for social bot detection," *Social Network Analysis and Mining*, vol. 12, pp. 1–11, 2022.
- [21] L. Yu, W. Zhang, J. Wang, and Y. Yu, "Seqgan: Sequence generative adversarial nets with policy gradient," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, ser. AAAI'17. AAAI Press, 2017, p. 2852–2858.
- [22] B. Wu, L. Liu, Y. Yang, K. Zheng, and X. Wang, "Using improved conditional generative adversarial networks to detect social bots on twitter," *IEEE Access*, vol. 8, pp. 36 664–36 680, 2020.
- [23] D. Martín-Gutiérrez, G. Hernández-Peñaloza, A. B. Hernández, A. Lozano-Diez, and F. Álvarez, "A deep learning approach for robust detection of bots in twitter using transformers," *IEEE Access*, vol. 9, pp. 54 591–54 601, 2021.

- [24] V. Chawla and Y. Kapoor, “A hybrid framework for bot detection on twitter: Fusing digital dna with bert,” *Multimedia Tools and Applications*, pp. 1–24, 2023.
- [25] A. S. Alhassun and M. A. Rassam, “A combined text-based and metadata-based deep-learning framework for the detection of spam accounts on the social media platform twitter,” *Processes*, vol. 10, no. 3, 2022. [Online]. Available: <https://www.mdpi.com/2227-9717/10/3/439>
- [26] Y. Wu, Y. Fang, S. Shang, J. Jin, L. Wei, and H. Wang, “A novel framework for detecting social bots with deep neural networks and active learning,” *Knowledge-Based Systems*, vol. 211, p. 106525, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950705120306547>
- [27] M. Fazil, A. K. Sah, and M. Abulaish, “Deepssbd: A deep neural network model with attention mechanism for socialbot detection,” *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 4211–4223, 2021.
- [28] C. Cai, L. Li, and D. Zeng, “Detecting social bots by jointly modeling deep behavior and content information,” in *Proceedings of the 2017 ACM Conference on Information and Knowledge Management*, ser. CIKM ’17. New York, NY, USA: Association for Computing Machinery, 2017, p. 1995–1998. [Online]. Available: <https://doi.org/10.1145/3132847.3133050>
- [29] F. Wei and U. T. Nguyen, “Twitter bot detection using neural networks and linguistic embeddings,” *IEEE Open Journal of the Computer Society*, pp. 1–12, 2023.
- [30] F. Wei and T. Nguyen, “A lightweight deep neural model for sms spam detection,” in *2020 International Symposium on Networks, Computers and Communications (ISNCC)*, 2020, pp. 1–6.
- [31] F. Wei and U. T. Nguyen, “Twitter bot detection using bidirectional long short-term memory neural networks and word embeddings,” in *2019 First IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications (TPS-ISA)*, 2019, pp. 101–109.
- [32] S. Feng, Z. Tan, R. Li, and M. Luo, “Heterogeneity-aware twitter bot detection with relational graph transformers,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 4, pp. 3977–3985, Jun. 2022. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/20314>
- [33] Z. Lei, H. Wan, W. Zhang, S. Feng, Z. Chen, J. Li, Q. Zheng, and M. Luo, “Bic: Twitter bot detection with text-graph interaction and semantic consistency,” 2023.
- [34] E. Alothali, M. Salih, K. Hayawi, and H. Alashwal, “Bot-mgat: A transfer learning model based on a multi-view graph attention network to detect social bots,” *Applied Sciences*, vol. 12, no. 16, 2022. [Online]. Available: <https://www.mdpi.com/2076-3417/12/16/8117>
- [35] H. Peng, Y. Zhang, H. Sun, X. Bai, Y. Li, and S. Wang, “Domain-aware federated social bot detection with multi-relational graph neural networks,” in *2022 International Joint Conference on Neural Networks (IJCNN)*, 2022, pp. 1–8.
- [36] S. Ali Alhosseini, R. Bin Tareaf, P. Najafi, and C. Meinel, “Detect me if you can: Spam bot detection using inductive representation learning,” in *Companion Proceedings of The 2019 World Wide Web Conference*, ser. WWW ’19. New York, NY, USA: Association for Computing Machinery, 2019, p. 148–153. [Online]. Available: <https://doi.org/10.1145/3308560.3316504>
- [37] L. Mannocci, S. Cresci, A. Monreale, A. Vakali, and M. Tesconi, “Mulbot: Unsupervised bot detection based on multivariate time series,” in *2022 IEEE International Conference on Big Data (Big Data)*. Los Alamitos, CA, USA: IEEE Computer Society, dec 2022, pp. 1485–1494. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/BigData55660.2022.10020363>
- [38] Z. Miller, B. Dickinson, W. Deitrick, W. Hu, and A. H. Wang, “Twitter spammer detection using data stream clustering,” *Information Sciences*, vol. 260, pp. 64–73, 2014. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0020025513008037>
- [39] Y. Yang, R. Yang, Y. Li, K. Cui, Z. Yang, Y. Wang, J. Xu, and H. Xie, “Rosgas: Adaptive social bot detection with reinforced self-supervised gnn architecture search,” *ACM Trans. Web*, vol. 17, no. 3, may 2023. [Online]. Available: <https://doi.org/10.1145/3572403>
- [40] G. Lingam, R. R. Rout, and D. V. Somayajulu, “Adaptive deep q-learning model for detecting social bots and influential users in online social networks,” *Applied Intelligence*, vol. 49, pp. 3947–3964, 2019.
- [41] R. R. Rout, G. Lingam, and D. V. L. N. Somayajulu, “Detection of malicious social bots using learning automata with url features in twitter network,” *IEEE Transactions on Computational Social Systems*, vol. 7, no. 4, pp. 1004–1018, 2020.
- [42] S. Cresci, R. Di Pietro, M. Petrocchi, A. Spognardi, and M. Tesconi, “The paradigm-shift of social spambots: Evidence, theories, and tools for the arms race,” in *Proceedings of the 26th International Conference on World Wide Web Companion*, ser. WWW ’17 Companion. Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee, 2017, p. 963–972. [Online]. Available: <https://doi.org/10.1145/3041021.3055135>
- [43] M. Avvenuti, S. Bellomo, S. Cresci, M. N. La Polla, and M. Tesconi, “Hybrid crowdsensing: A novel paradigm to combine the strengths of opportunistic and participatory crowdsensing,” in *Proceedings of the 26th International Conference on World Wide Web Companion*, ser. WWW ’17 Companion. Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee, 2017, p. 1413–1421. [Online]. Available: <https://doi.org/10.1145/3041021.3051155>
- [44] S. Feng, H. Wan, N. Wang, J. Li, and M. Luo, “Twibot-20: A comprehensive twitter bot detection benchmark,” in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, ser. CIKM ’21. New York, NY, USA: Association for Computing Machinery, 2021, p. 4485–4494. [Online]. Available: <https://doi.org/10.1145/3459637.3482019>
- [45] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1–9.
- [46] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [47] S. Zagoruyko and N. Komodakis, “Wide residual networks,” *arXiv preprint arXiv:1605.07146*, 2016.
- [48] A. Krizhevsky, “One weird trick for parallelizing convolutional neural networks,” *arXiv preprint arXiv:1404.5997*, 2014.
- [49] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, “Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size,” *arXiv preprint arXiv:1602.07360*, 2016.
- [50] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2261–2269.
- [51] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.
- [52] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5987–5995.
- [53] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [54] M. Tan and Q. Le, “EfficientNet: Rethinking model scaling for convolutional neural networks,” in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 09–15 Jun 2019, pp. 6105–6114. [Online]. Available: <https://proceedings.mlr.press/v97/tan19a.html>
- [55] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [56] J. Arevalo, T. Solorio, M. Montes-y Gomez, and F. A. González, “Gated multimodal networks,” *Neural Computing and Applications*, pp. 1–20, 2020.
- [57] Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov, “Multimodal transformer for unaligned multimodal language sequences,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 6558–6569. [Online]. Available: <https://aclanthology.org/P19-1656>
- [58] L. Ilias, D. Askounis, and J. Psarras, “Detecting dementia from speech and transcripts using transformers,” *Computer Speech & Language*, vol. 79, p. 101485, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0885230823000049>
- [59] D. Sánchez Villegas and N. Aletras, “Point-of-interest type prediction using text and images,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Online and

Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 7785–7797. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.614>

- [60] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS’17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 6000–6010.
- [61] H. Ping and S. Qin, “A social bots detection model based on deep learning algorithm,” in *2018 IEEE 18th International Conference on Communication Technology (ICCT)*, 2018, pp. 1435–1439.
- [62] F. Ahmed and M. Abulaish, “A generic statistical approach for spam detection in online social networks,” *Computer Communications*, vol. 36, no. 10, pp. 1120–1129, 2013. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0140366413001047>
- [63] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush, “Transformers: State-of-the-art natural language processing,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics, Oct. 2020, pp. 38–45. [Online]. Available: <https://www.aclweb.org/anthology/2020.emnlp-demos.6>
- [64] K.-C. Yang, O. Varol, A. C. Nwala, M. Sayyadiharikandeh, E. Ferrara, A. Flammini, and F. Menczer, “Social bots: Detection and challenges. in: Yasserli, t. (ed.),” in *Handbook of Computational Social Science*. Edward Elgar Publishing Ltd, 2024.
- [65] K.-C. Yang, O. Varol, P.-M. Hui, and F. Menczer, “Scalable and generalizable social bot detection through data selection,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 01, pp. 1096–1103, Apr. 2020. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/5460>
- [66] E. Ferrara, “Social bot detection in the age of chatgpt: Challenges and opportunities,” *First Monday*, vol. 28, no. 6, Jun. 2023. [Online]. Available: <https://firstmonday.org/ojs/index.php/fm/article/view/13185>



Ioannis Michail Kazelidis received the integrated master’s degree from the School of Electrical and Computer Engineering (SECE), National Technical University of Athens (NTUA), Athens, Greece, in October 2023. His research interests include deep learning and social media analysis.



Dimitris Askounis was the Scientific Director of over 50 European research projects in the above areas (FP7, Horizon2020, and so on). For a number of years, he was an Advisor to the Minister of Justice and the Special Secretary for Digital Convergence for the introduction of information and communication technologies in public administration. Since June 2019, he has been the President of the Information Society SA, Kallithea, Greece. He is currently a Professor at the School of Electrical and Computer Engineering, National Technical University of Athens (NTUA), Athens, Greece, and the Deputy Director of the Decision Support Systems Laboratory. He has over 25 years of experience in decision support systems, intelligent information systems and manufacturing, e-business, e-government, open and linked data, big data analytics, Artificial Intelligence (AI) algorithms, and the application of modern Information Technology (IT) techniques in the management of companies and organizations.

XI. BIOGRAPHY SECTION



Loukas Ilias received the integrated master’s degree from the School of Electrical and Computer Engineering (SECE), National Technical University of Athens (NTUA), Athens, Greece, in June 2020, where he is currently pursuing the Ph.D. degree with the Decision Support Systems (DSS) Laboratory, SECE. He has completed a Research Internship with University College London (UCL), London, U.K.

He is a Researcher with the DSS Laboratory, NTUA, where he is involved in EU-funded research projects. He has published in numerous journals, including *IEEE Journal of Biomedical and Health Informatics*, *IEEE Transactions on Computational Social Systems*, *Knowledge-Based Systems (Elsevier)*, *Expert Systems With Applications (Elsevier)*, *Applied Soft Computing (Elsevier)*, *Online Social Networks and Media (Elsevier)*, *Computer Speech and Language (Elsevier)*, *IEEE Access*, *Frontiers in Aging Neuroscience*, and *Frontiers in Big Data*. His research has also been accepted for presentation at international conferences, including the *IEEE-EMBS International Conference on Biomedical and Health Informatics (BHI’22)* and the *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP) 2023*. His research interests include speech processing, natural language processing, social media analysis, and the detection of complex brain disorders.