

# BotArtist: Generic approach for bot detection in Twitter via semi-automatic machine learning pipeline

Alexander Shevtsov<sup>1,2,3</sup>, Despoina Antonakaki<sup>1,2</sup>, Ioannis Lamprou<sup>1</sup>, Polyvios Pratikakis<sup>3</sup>, Sotiris Ioannidis<sup>1,2</sup>

<sup>1</sup>Technical University of Crete

<sup>2</sup>Foundation for Research and Technology Hellas - FORTH

<sup>3</sup>University of Crete, Computer Science Department

shevtsov@csd.uoc.gr, dantonakaki@tuc.gr, ilamprou1@tuc.gr

## Abstract

Twitter, as one of the most popular social networks, provides a platform for communication and online discourse. Unfortunately, it has also become a target for bots and fake accounts, resulting in the spread of false information and manipulation. This paper introduces a semi-automatic machine learning pipeline (SAMPLP) designed to address the challenges associated with machine learning model development. Through this pipeline, we develop a comprehensive bot detection model named BotArtist, based on user profile features. SAMLP leverages nine distinct publicly available datasets to train the BotArtist model. To assess BotArtist’s performance against current state-of-the-art solutions, we evaluate 35 existing Twitter bot detection methods, each utilizing a diverse range of features. Our comparative evaluation of BotArtist and these existing methods, conducted across nine public datasets under standardized conditions, reveals that the proposed model outperforms existing solutions by almost 10% in terms of F1-score, achieving an average score of 83.19% and 68.5% over specific and general approaches, respectively. As a result of this research, we provide one of the largest labeled Twitter bot datasets. The dataset contains extracted features combined with BotArtist predictions for 10,929,533 Twitter user profiles, collected via Twitter API during the 2022 Russo-Ukrainian War over a 16-month period. This dataset was created based on (Shevtsov et al. 2022b) where the original authors share anonymized tweets discussing the Russo-Ukrainian war, totaling 127,275,386 tweets. The combination of the existing textual dataset and the provided labeled bot and human profiles will enable future development of more advanced bot detection large language models in the post-Twitter API era.

## Introduction

Online social media has become an essential part of everyday life. During the past decade, social networks have transformed the communication routine of our daily lives. Due to their growing popularity, online social media gained millions of daily active users not only consuming the information but also creating a space for content creators. The main reason behind the success of online social networks is the real-time access to unlimited information, where registered users can share their comments and personal opinions on popular topics.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Twitter, one of the most popular social networks, with millions of active users, is used for news dissemination, political discussions, and social interactions. However, the platform has also been plagued by bots and fake accounts that are used to manipulate and spread false information. According to the research community, the usage of manipulation techniques implemented with the use of bot accounts is registered during diverse popular topic discussions. More specifically, studies show that bot accounts are involved in diverse political discussions, for instance, around 2016 in countries like the US, Germany, Sweden, France, Spain, etc. (Golovchenko et al. 2020; Badawy, Ferrara, and Lerman 2018; Howard, Kollanyi, and Woolley 2016; Shevtsov et al. 2022a, 2023; Neudert, Kollanyi, and Howard 2017; Pastor-Galindo et al. 2020; Bradshaw et al. 2017; Fernquist, Kaati, and Schroeder 2018; Castillo et al. 2019; Rossi et al. 2020). Furthermore, propagandizing and false news is identified during the vaccination debate (Broniatowski et al. 2018) and more recent examples of the COVID-19 pandemic (Shahi, Dirkson, and Majchrzak 2021; Ferrara 2020; Yang et al. 2021). This high activity of bot accounts raises concerns in the research community and online social media platforms about the integrity of the shared information.

Currently, the research community offers a variety of ML and NN-based bot detection methods. Each of these provided methods excels at addressing specific bot detection scenarios, yielding optimal solutions for particular use cases. Unfortunately, the performance of these provided methods significantly diminishes in more general bot detection scenarios, encompassing different periods, discussion topics, and languages, among others. Recent studies have shown that existing methods still do not achieve flawless Twitter bot detection, regardless of the number of user characteristics that they incorporate (Feng et al. 2022b). In addition, a large portion of existing ML bot detection methods ignore well-known optimization procedures (including feature selection and dynamic hyper-parameter fine-tuning) providing significant space for improvement.

In our research, we address this challenge by developing a semi-automatic machine learning pipeline (SAMPLP<sup>1</sup>) for constructing a generic bot detection model (BotArtist<sup>2</sup>).

<sup>1</sup>GitHub repository: <https://github.com/alexdrk14/SAMPLP>

<sup>2</sup>GitHub repository: <https://github.com/alexdrk14/BotArtist>

The presented pipeline resolves popular issues during the creation of the classification model, including recursive hyperparameter fine-tuning during the feature selection phase. This approach guarantees the reduction of non-relevant features which sequentially reduces the noise of initial data. Additionally, our pipeline considers the class imbalance during feature selection and model fine-tuning, making the model suitable for real-case applications. Next, we compare this developed model with 35 currently available state-of-the-art approaches using nine publicly available datasets. This research yields two distinct evaluation scenarios: one specific to each dataset, where each model is trained and tested separately on each of them, and a more generic approach, where the models are trained and tested on the combined dataset comprising all nine unique datasets.

As a result of our research, we introduce BotArtist, a finely tuned general bot detection model that achieves the highest F1-score performance on three of the nine datasets, with an average F1-score performance of 83.19. Furthermore, the presented approach achieves greater precision in general evaluation, with an overall improvement in F1 score of more than 9% and an average improvement of 5% compared to the best-performing existing bot detection approaches. Moreover, BotArtist relies on only a limited set of profile features, allowing for the prediction of historical data independently of the Twitter API.

Except from the source code of SAMLP and BotArtist, we make available the model outcome usable for research. For this purpose, we collect a large volume of users correlated to the already existing textual shared dataset and add an additional layer to textual information via user prediction. We then release 10.929.533 user profiles with BotArtist predictions<sup>3</sup>, correlated to 127.275.386 tweets (Shevtsov et al. 2022b).

## Related Work

Bot detection on Twitter poses a formidable challenging task, primarily due to the increasing sophistication of these bots. Research efforts in this domain can generally be categorized into three main approaches: feature-based, text-based, and graph-based techniques. Each category offers a distinct angle on bot detection, leveraging various extracted characteristics of users and their activities on the social media platform.

### Feature based

Methods falling into this category apply feature engineering based on information extracted from Twitter user profiles and their activity patterns. Researchers use traditional machine learning or neural network classification algorithms for bot detection. These approaches employ different sets of user characteristics as feature sets, including user metadata (Kudugunta and Ferrara 2018), tweets (Miller et al. 2014), user name (Beskow and Carley 2019), description (Hayawi et al. 2022), temporal patterns (Mazza et al. 2019), and follow relationships (Feng et al. 2021a). Moreover, some of

them aim to enhance the scalability of feature-based methods (Yang et al. 2020; Abreu, Ralha, and Gondim 2020), discover unknown bot accounts using correlation-based techniques (Chavoshi, Hamooni, and Mueen 2016), and improve the trade-off between precision and recall in bot detection (Morstatter et al. 2016). However, the creators of bot accounts are increasingly aware of the features utilized by the research community, leading to novel bot implementations designed to evade detection based on known features (Cresci 2020). Consequently, existing feature-based methods face several challenges in the accurate detection of these new bot accounts (Feng et al. 2021a).

### Text based

Text-based methods rely on natural language processing (NLP) techniques for Twitter bot detection, extracting characteristics from posted content (tweets) and user profile descriptions. Approaches in this category include sequence fingerprint (Cresci et al. 2016), word embeddings (Wei and Nguyen 2019), recurrent neural networks (RNN) (Kudugunta and Ferrara 2018), attention mechanisms (Feng et al. 2021a), transformers (Guo et al. 2021a; Liu et al. 2019; Raffel et al. 2020) and the adoption of pre-trained language models for encoding tweets (Đukić, Keča, and Stipić 2020). Researchers have also combined tweet representations with user profile characteristics (Cai, Li, and Zengi 2017; Efthimion, Payne, and Proferes 2018; Kantepe and Ganiz 2017; Miller et al. 2014; Varol et al. 2017; Kouvela, Dimitriadis, and Vakali 2020; Lee, Eoff, and Caverlee 2011; Echeverria et al. 2018), employed unsupervised machine learning models (Feng et al. 2021a), and addressed multilingual context issues (Knauth 2019). Despite extensive existing research studies and impressive performance in text-based approaches, new bot accounts can still evade detection by sharing stolen content from genuine users (Cresci 2020). Additionally, recent work has demonstrated that relying solely on textual information is insufficient for robust and accurate bot detection (Feng et al. 2021c).

### Graph based

The graph-based category of bot detection methods combines geometric deep neural networks with graph analytics. Current implementations leverage techniques such as node centrality (Dehghan et al. 2023), node representation learning (Pham et al. 2022), graph neural networks (GNNs) (Ali Alhosseini et al. 2019; Moghaddam and Abbaspour 2022), and heterogeneous GNNs (Feng et al. 2021c) to conduct graph-based bot detection. In recent research, the authors have borrowed ideas from other categories to merge different approaches, such as combining multiple methods (Guo et al. 2021b; Yang, Ferrara, and Menczer 2022; Knauth 2019; Beskow and Carley 2020; Rodríguez-Ruiz et al. 2020; Magelinski, Beskow, and Carley 2020; Yang, Harkreader, and Gu 2013; Kipf and Welling 2016,?; Veličković et al. 2017; Lv et al. 2021). Furthermore, novel GNN architectures have been proposed to exploit heterogeneity in the Twitter network (Feng et al. 2022a). Generally, these approaches hold significant promise for Twitter bot detection.

<sup>3</sup>Zenodo repository: <https://zenodo.org/records/11203900>

Dataset	C-15	G-17	C-17	M-18	C-S-18	C-R-19	B-F-19	Twibot-20	Twibot-22
# Total User	5,301	2,484	14,368	50,538	13,276	693	518	229,580	1,000,000
# Human	1,950	1,394	3,474	8,092	6,174	340	380	5,237	860,057
# Bot	3,351	1,090	10,894	42,446	7,102	353	138	6,589	139,943
# Total Tweet	2,827,757	0	6,637,615	0	0	0	0	33,488,192	88,217,457
# Human Tweet	2,631,730	0	2,839,361	0	0	0	0	927,292	81,250,102
# Bot Tweet	196,027	0	3,798,254	0	0	0	0	1,072,496	6,967,355
# Graph Edges	7,086,134	0	6,637,615	0	0	0	0	33,488,192	170,185,937

Table 1: Description of selected datasets and information contained in those datasets.

Unfortunately, existing Twitter bot detection methods have their limitations. Simplified feature-based approaches struggle to generalize and provide robust results, while more complex methods based on text or graph characteristics require intensive computational resources and large datasets for model development. Furthermore, as shown in (Shevtsov et al. 2024) the complex embedding-based model tends to perfectly capture the training data patterns, but lacks generalizability and requires frequent retraining over multiple datasets to keep high performance. These limitations, coupled with the recent announcement of the Twitter API monetization (Twitter 2023), make many existing methods expensive to maintain or nearly impossible to operate daily. In this study, we address these challenges and present robust and accurate Twitter bot detection methods, based on a lightweight set of features. Our methods require only a single user profile object (Twitter API v1.1 or v2) to provide user prediction without additional Twitter API requests, thereby significantly reducing operational costs.

## Datasets

For this research paper, we collect nine well-known publicly available datasets (Feng et al. 2021b). All selected datasets already contain ground truth labels, primarily obtained through manual analysis or crowd-sourcing. For simplicity, we label the selected datasets as follows: C-15 (Cresci et al. 2015), G-17 (Gilani et al. 2017), C-17 (Cresci et al. 2017b,a), M-18 (Yang et al. 2020), C-S-18 (Cresci et al. 2018, 2019), C-R-19 (Mazza et al. 2019), B-F-19 (Yang et al. 2019), TwiBot-20 (Feng et al. 2021b), and TwiBot-22 (Feng et al. 2022b). In Table 1, we present the information provided in each dataset, along with the volume of normal and bot accounts.

Furthermore, in collaboration with (Shevtsov et al. 2022b) collect 10,929,533 Twitter profiles correlated with 127,275,386 publicly available tweets related to the public discussion topic of the 2022 Russo-Ukrainian War. Our collection of user profiles is based on the monitoring of selected topics starting from February 23, 2022, till June 23, 2023. The shared datasets contain a set of extracted features in an anonymized form of a CSV file and contain only preprocessed numerical features to protect user information. The provided dataset also provides anonymized user IDs which are identically correlated with publicly available user tweet dataset (Shevtsov et al. 2022b).

## Methodology

The proposed approach is based on the development of a semi-automatic machine learning pipeline (SAMLP). We choose this implementation due to its simplicity for further usage and its ability to avoid many trivial mistakes during the ML model development such as data processing, feature selection, hyper-parameter fine-tuning, and model evaluation. The steps of the developed pipeline are presented in Figure 1, where the initial step involves data splitting.

During data splitting, it is crucial to maintain a class balance between the training/validation and testing data portions. For this purpose, we apply a stratified data split with a 70:30 ratio for training, validation, and testing portions respectively. Such an approach allows us to preserve the natural class distribution difference. The testing data portion is kept hidden and is only utilized during the final testing of the model to prevent information leakage between train and test portions

## Feature Selection

After the data splitting, our pipeline applies a feature selection procedure over the train/validation data portion during the K-Fold cross validation. Although K-Fold cross-validation may lead to slightly over-optimistic performance estimations, in our case, it is not relevant since we use this procedure for  $\alpha$  hyper-parameter fine-tuning of the Lasso model, and we are interested in identifying the best  $\alpha$  parameter.

Additionally, since the Lasso regression model does not perform well over imbalanced datasets, we also take into account the class imbalance of the dataset. In case of data imbalance, we utilize under-sampling of the majority class, making the data perfectly balanced and improving the prediction performance of the Lasso model. Due to the high-class imbalance and the under-sampling of the majority class, there is a very high probability of losing some important samples, which may lead to inaccurate feature selection. To address this, we utilize 10 consecutive repetitions to capture as many samples of the majority class as possible.

During these repetition rounds, we store the best  $\alpha$  parameters based on the square mean error metric from the entire stratified k-Fold cross-validation results, where  $K = 5$ . At the end of the procedure, we compute the most frequent best  $\alpha$  parameter and train the Lasso model with the entire train/validation dataset, keeping features with non-zero coefficients. Additionally, we check if the selected  $\alpha$  param-

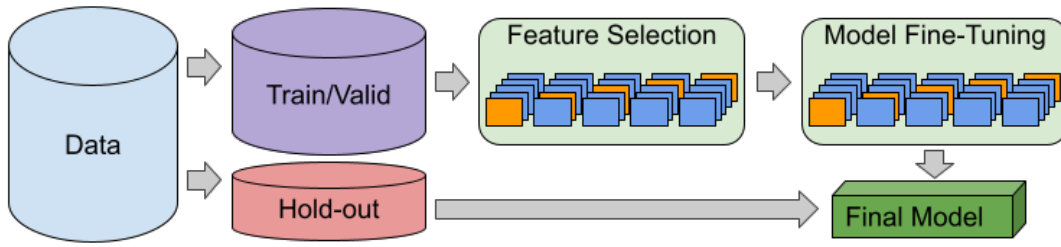


Figure 1: Machine learning pipeline used for *BotArtist* model creation including feature selection and model fine-tuning; each step is executed during separate K-Fold cross-validation.

ter falls within the defined limits of the searching area. If it falls on the minimum or maximum values, we create a new search area based on the original site to ensure the selected  $\alpha$  is genuinely the best.

This approach allows us to utilize feature selection without additional knowledge of the original data and modifications, as the procedure automatically finds the best  $\alpha$  values, manages class balance and imbalance via under-sampling and repeated executions, and processes both perfectly balanced and highly imbalanced datasets without information loss.

### Model Fine-Tuning

Having reduced the problem’s dimensionality through feature selection, we now focus on identifying the model and its configuration (hyper-parameter set) that can provide accurate predictions over unseen data samples. We select three well-known machine learning classification models: SVM (Boser, Guyon, and Vapnik 1992), RandomForest (Breiman 2001), and state-of-the-art XGBoost (Extreme Gradient Boosting) (Chen et al. 2015).

For each model, we define large hyper-parameter ranges to cover a variety of configurations. To reduce complexity and execution time, we develop a sampling method that randomly selects  $C$  unique configurations per model (in our case,  $C = 50$ ). This method allows us to estimate different possible configurations without sacrificing execution time. Additionally, since the dataset may be affected by class imbalance, we compute class-specific weights to construct selected classifiers with consideration of class imbalance.

To evaluate each selected configuration, we use a stratified K-Fold cross-validation approach, allowing us to utilize the entire train/validation data portion and evaluate each configuration over multiple validation data portions. During the evaluation, we compute the average configuration performance measured in the F1-score, which can effectively measure model performance over both binary and multi-class classification tasks, considering class imbalance.

Afterward, we select the best classifier and configuration based on the average validation (F1-score). This model configuration becomes the final model, which will be evaluated over the hold-out (testing) data portion. In the final stage of the model development, we provide model explainability using the SHAP game-theoretical approach (Shapley 1953). This allows us to describe the reasons behind the model’s

decisions and highlight the distinctions between classes.

Additionally, in the case of binary classification, we need to adjust the decision threshold, as binary classification models are not perfectly aligned with the 50% decision threshold. To do this, we utilize the precision vs. recall curve to identify the decision threshold that maximizes the model’s performance. The developed model is then trained over the entire dataset, including the train/validation and testing sets, and can be utilized for real-world applications with the identified decision threshold.

Described experiments were conducted in the single machine with AMD Ryzen 9 CPU (16 cores/32 threads) with a total of 64 GB DDR4 system memory. For the graphical unit based computations the single unit of Nvidia RTX 3080 where used with total memory of 12 GB.

### Feature Extraction

As mentioned earlier, in this study we pursue the creation of a simple model that in comparison with other bot detection systems- does not require a large volume of information or features for accurate prediction. In our case, the number of Twitter profiles triggers our attention. The majority of social media bot accounts are created for content promotion and increasing the attention of other registered users. To maximize this goal, bot accounts aim for the expansion of social audience. Social expansion is very difficult to achieve without high activity, such as tweets, re-tweets, comments, etc. In comparison with profile-based features, the extraction and utilization of textual or graph features can result in instability, causing a significant performance drop across various periods or discussion topic datasets (Shevtsov et al. 2024). Consequently, we need to extract as much information as possible from the user profiles to identify extremely active accounts. The difference between very popular social media accounts (also known as celebrities), and bot accounts is that they try to post and share as much content as possible during the day. Due to the differences between Twitter API v1.1 and v2, we have devised a method for constructing features in a way that allows the extraction of selected information from both v1.1 and v2 without any additional manipulation. The profile features we have extracted are categorized into three distinct groups: count, boolean, and real-valued.

The initial count feature categories encompass raw Twitter API object values, including followers, followings, status, and number of subscribed Twitter lists. In addition to

Feature	Type	Feature	Type	Calculation
name_len	count	screen_name_sim	real-valued	Jaccard of (screen_name, user_name)
screen_name_len	count	foll_friends	real-valued	follower / friends
description_len	count	age	real-valued	created_at / collection date
listed	count	listed_by_age	real-valued	listed / account age
statuses	count	statuses_by_age	real-valued	statuses / account age
followers	count	followers_by_age	real-valued	followers / account age
following	count	following_by_age	real-valued	friends / account age
name_upper_len	count	name_upper_pcmt	real-valued	percentage of upper case
name_lower_len	count	name_lower_pcmt	real-valued	percentage of lower case
name_digits_len	count	name_digits_pcmt	real-valued	percentage of digits
name_special_len	count	name_special_pcmt	real-valued	percentage of other characters
screen_name_upper_len	count	screen_name_upper_pcmt	real-valued	percentage of upper case
screen_name_lower_len	count	screen_name_lower_pcmt	real-valued	percentage of lower case
screen_name_digits_len	count	screen_name_digits_pcmt	real-valued	percentage of digits
screen_name_special_len	count	screen_name_special_pcmt	real-valued	percentage of other characters
description_upper_len	count	description_upper_pcmt	real-valued	percentage of upper case
description_lower_len	count	description_lower_pcmt	real-valued	percentage of lower case
description_digits_len	count	description_digits_pcmt	real-valued	percentage of digits
description_special_len	count	description_special_pcmt	real-valued	percentage of other characters
description_urls	count	name_entropy	real-valued	entropy of user name
description_mentions	count	screen_name_entropy	real-valued	entropy of the screen name
description_hashtags	count	has_location	boolean	
total_urls	count	has_profile_image	boolean	
protected	boolean	has_profile_url	boolean	
verified	boolean			

Table 2: List of the profile features extracted from the Twitter API user objects.

these raw values, we calculate statistics for textual fields such as user name, screen\_name, and description. For these features, we quantify the number of distinct characters, including uppercase letters, lowercase letters, digits, and special characters that do not fall into any of the previous categories. Furthermore, we leverage the profile description field, where users can provide additional information such as mentions of other users/organizations, hashtags, and additional URLs. In this context, we have counted the volume of the mentioned entities in the description field as raw values.

In contrast to the count categories, where values are represented as integers, the real-valued categories contain floating-point values. We conduct various comparisons and measurements across the profile fields in this category. Initially, we compute the age of the accounts, which allows us to distinguish highly active users with millions of posts spanning several years from users with similar activity but with an active period of only a few days. Account age is measured in days, calculated as the duration between the account creation date, provided by the Twitter API, and the collection date. Based on the account’s age, we can compute metrics similar to those in the count values category but adjusted for the account’s age perspective.

Given that many bot accounts are mass-created by automated scripts, they often have very similar or even trivial user names and screen names. Therefore, we apply the Jaccard similarity measure between user names and screen

names to calculate the intersection of characters used in both fields to the differences between them.

Additionally, we compute the entropy of the user name and screen\_name, enabling us to assess the randomness of these selected fields. In addition to the count values derived from the profile text fields (user name, screen\_name, and description), we compute the percentage of each specific character category, including lowercase, uppercase, digits, and special characters, relative to the length of the text field. These measurements offer a more precise means of identifying textual preferences more flexible to varying text lengths since they are computed as percentages of the field length.

The final category of extracted features consists of Boolean values, used to identify information that may or may not be provided by the account creator. For instance, we determine if the user account is protected, and verified, as well as extract information about the provided user location, profile URL, and whether the profile uses the default profile image.

By employing this comprehensive feature engineering approach, we maximize the feature extraction while working with the limited profile information available through Twitter API v2. In total, we extract 49 unique profile features Table 2, representing the full set of features used before the feature selection process in the BotArtist model.

Method	Type	C-15	G-17	C-17	M-18	C-S-18	C-R-19	B-F-19	TB-20	TB-22	Average
SGBot	F	77.9	72.1	94.6	99.5	82.3	82.7	49.6	84.9	36.6	75.57
Kudugunta et al.	F	75.3	49.8	91.7	94.5	50.9	49.2	49.6	47.3	51.7	62.22
Hayawi et al.	F	85.6	34.7	93.8	91.5	60.8	60.9	20.5	77.1	24.7	61.06
BotHunter	F	97.2	69.2	91.6	<u>99.6</u>	82.2	82.9	49.6	79.1	23.5	74.98
NameBot	F	83.4	44.8	85.7	<u>91.6</u>	61.1	67.5	38.5	65.1	0.5	59.80
Abreu et al.	F	76.4	66.7	95.0	97.9	76.9	<u>83.5</u>	<u>53.8</u>	77.1	53.4	75.63
BotArtist	F	98.3	<u>76.1</u>	97.0	<b>99.7</b>	80.6	<b>88.3</b>	<b>68.4</b>	82.2	<u>58.2</u>	<b>83.19</b>
Cresci	T	1.17	-	22.8	-	-	-	-	13.7	-	-
Wei	T	82.7	-	78.4	-	-	-	-	57.3	53.6	-
BGSRD	T	90.8	35.7	86.3	90.5	58.2	41.1	13.0	70.0	21.1	56.30
RoBERTa	T	95.8	-	94.3	-	-	-	-	73.1	20.5	-
T5	T	89.3	-	92.3	-	-	-	-	70.5	20.2	-
Efthimion	FT	94.1	5.2	91.8	95.9	68.2	71.7	0.0	67.2	27.5	57.95
Kantepe	FT	78.2	-	79.4	-	-	-	-	62.2	<b>58.7</b>	-
Miller	FT	83.8	59.9	86.8	91.1	56.8	43.6	0.0	74.8	45.3	60.23
Varol	FT	94.7	-	-	-	-	-	-	81.1	27.5	-
Kouvela	FT	98.2	66.6	<u>99.1</u>	98.2	80.4	81.1	28.1	86.5	30.0	74.24
Santos	FT	78.8	14.5	83.0	92.4	65.2	75.7	21.0	60.3	-	-
Lee	FT	<u>98.6</u>	67.8	<b>99.3</b>	97.9	<u>82.5</u>	82.7	50.3	80.0	30.4	<u>76.61</u>
LOBO	FT	<b>98.8</b>	-	97.7	-	-	-	-	80.8	38.6	-
Moghaddam	FG	73.9	-	-	-	-	-	-	79.9	32.1	-
Alhosseini	FG	92.2	-	-	-	-	-	-	72.0	38.1	-
Knauth	FTG	91.2	39.1	93.4	91.3	<b>94.0</b>	54.2	41.3	85.2	37.1	69.64
FriendBot	FTG	97.6	-	87.4	-	-	-	-	80.0	-	-
SATAR	FTG	95.0	-	-	-	-	-	-	86.1	-	-
Botometer	FTG	66.9	<b>77.4</b>	96.1	46.0	79.6	79.0	30.8	53.1	42.8	63.5
Rodríguez-Ruiz	FTG	87.7	-	85.7	-	-	-	-	63.1	56.6	-
GraphHist	FTG	84.5	-	-	-	-	-	-	67.6	-	-
EvolveBot	FTG	90.1	-	-	-	-	-	-	69.7	14.1	-
Dehghan	FTG	88.3	-	-	-	-	-	-	76.2	-	-
GCN	FTG	97.2	-	-	-	-	-	-	80.8	54.9	-
GAT	FTG	97.6	-	-	-	-	-	-	85.2	55.8	-
HGT	FTG	96.9	-	-	-	-	-	-	<b>88.2</b>	39.6	-
SimpleHGN	FTG	97.5	-	-	-	-	-	-	<b>88.2</b>	45.4	-
BotRGCN	FTG	97.3	-	-	-	-	-	-	87.3	57.5	-
RGT	FTG	97.8	-	-	-	-	-	-	<u>88.0</u>	42.9	-

Table 3: The performance of each selected bot detection model, as reported in the (Feng et al. 2021b) paper, is compared with that of BotArtist. Performance is measured using the F1-score. In this benchmark, each model is trained and tested on each dataset separately.

## Experimental Results

According to the methodology of SAMLP described earlier, we successfully develop a semi-automated machine learning pipeline capable of constructing classification models (both binary and multiclass) for both balanced and imbalanced datasets. Using this pipeline, we create and test our custom implementation called BotArtist for Twitter bot detection, utilizing nine well-known datasets. To evaluate BotArtist’s performance, in comparison to 35 other existing research approaches, we employ the most comprehensive Twitter bot detection benchmark available (Feng et al. 2022b), where various methods are compared using identical training/validation and testing data portions. This benchmarking allows us to make a fair comparison between our approach and established bot detection methods.

In our comparison, we not only evaluate BotArtist against bot detection approaches based on similar categories of features but also compare it to existing approaches developed based on different sets of characteristics. The selected methods are categorized into five different groups: feature-based (F), textual (T), graph (G), and any combination of them.

### Models comparison

Our goal with BotArtist is to create a model with the highest possible level of generalizability. To achieve this, we design two separate comparison scenarios. Initially, our focus is on assessing the fine-tuned BotArtist’s ability to capture general patterns within each dataset in our collection. As a result, we conduct training and testing for each of the selected models, including BotArtist, on each of the nine chosen datasets

Method	C-15	G-17	C-17	M-18	C-S-18	C-R-19	B-F-19	TB-20	TB-22	Total	Average
BotArtist	82.7	<u>39.9</u>	87.3	<u>99.0</u>	<b>80.6</b>	<u>73.8</u>	16.6	<b>80.3</b>	<b>56.9</b>	<b>63.7</b>	<b>68.5</b>
Lee	82.3	0.0	83.6	97.7	78.2	67.7	20.0	8.5	42.4	52.9	53.3
Abreu	<u>84.4</u>	0.3	80.1	88.4	67.1	40.8	11.7	15.6	29.0	40.4	46.3
SGBot	75.0	3.6	79.8	<b>99.2</b>	76.7	68.9	0.0	15.2	<u>43.3</u>	<u>53.8</u>	51.5
BotHunter	73.4	7.1	76.0	<b>99.2</b>	76.0	44.8	11.1	14.7	28.0	43.1	47.8
Kouvela	<b>95.5</b>	20.4	<u>94.7</u>	98.1	78.4	71.4	<u>21.0</u>	28.5	36.0	52.0	60.5
Botometer	66.9	<b>77.4</b>	<b>96.1</b>	46.0	<u>79.6</u>	<b>79.0</b>	<b>30.8</b>	<u>53.1</u>	42.8	45.3	<u>63.5</u>

Table 4: The measurement of performance in the case of general bot detection approaches, involves training and testing models on all datasets. Performance is assessed using the F1-score.

individually. This comparative approach allows us to evaluate how well these selected models perform on each specific dataset, serving as a more rigorous bot detection implementation.

Additionally, the designed experiments provide us with a broader understanding of the overall capabilities of the selected methods when applied in more general scenarios. While implementing these experiments, we take into account that the performance would likely be higher than that of a generic approach. Therefore, our primary focus is on discerning the differences between the models to select the most accurate bot detection methods for the broader, more generalized comparison. Based on the designed experiment we identified that (see Table 3) the BotArtist model emerged as the only bot detection method that achieves superior performance across three different datasets (M-18, C-R-19, B-F-19). Furthermore, the BotArtist model demonstrates the highest average F1-score, reaching 83.19. This represents a substantial improvement of 6.5% over the most accurate of the existing methods. As shown in Table 3, some datasets posed challenges for more complex methods that rely on text and graph features. These methods struggle to utilize datasets without text or graph information due to their limitations in accessing essential information required for the model.

Previously, we had evaluated and compared models on each isolated dataset and calculated their average performance. While this approach provided us with an initial insight into model performance across diverse datasets, it lacked an assessment of model generalizability. To address this, we design additional comparisons where we assess the model’s ability to capture general patterns across different datasets. In these cases, we combine the training data from the nine datasets into a single dataset for model training. Combining data from multiple datasets, each with varying periods, discussion topics, and communities, demands a high level of model generalization to avoid overfitting.

To measure the performance of the trained methods, we evaluate them by testing on portions of each dataset and also on the union of all the dataset testing portions. These approaches allow us to not only identify the general performance of each bot detection method but also gain insights into performance variations across each specific dataset. To achieve this, we selectively choose the most accurate models from Table 3, in addition to the Botometer model (Yang, Ferrara, and Menczer 2022). The Botometer model is included

because it has already been trained on a diverse range of datasets and is expected to perform well in general-case scenarios.

The results, as presented in Table 4, indicate that the most generalizable models are BotArtist, Botometer, and SGBot. BotArtist demonstrates superior performance over existing methods in terms of both total and average scores, surpassing almost 10% the existing methods in total. These results confirm that fine-tuned models which are based on the limited set of features are capable of capturing the general differences between normal and bot accounts.

## Conclusions and Future Work

In this research paper, we introduce a semi-automatic machine learning pipeline (SAMPLP) that addresses multiple challenges in creating machine learning models, including feature selection, hyperparameter fine-tuning, model evaluation, decision threshold optimization for binary classification, and provides model explainability through the SHAP game-theoretical approach. We apply this developed pipeline to create BotArtist, a versatile bot detection model based on user profile features. Our approach is trained and evaluated alongside current state-of-the-art solutions using nine different datasets. As demonstrated in our experiments, BotArtist surpasses existing, more complex methods in terms of generalization, achieving almost a 10% increase in the total F1 score in comparisons of multiple data sets and a 6.5% increase in comparisons of separated data sets. Additionally, we offer insights into the final model’s decision-making process and the patterns it captures, based on the SHAP game-theoretical approach (available at the repository).

One of the limitations of the presented methodology is the relatively limited feature set, which might make it susceptible to evasion by future bot accounts that specifically aim to avoid detection based on these features. Further research conducted over a longer period is required to assess the model’s effectiveness in detecting bot accounts in the evolving landscape of social networks.

For further research, we provide one of the largest labeled datasets containing user profiles correlated with another publicly available dataset of the 127M user posts. Such large labeled data will allow the development of state-of-the-art LLM models to detect bot-generated content in the post-Twitter API era.

## References

- Abreu, J. V. F.; Ralha, C. G.; and Gondim, J. J. C. 2020. Twitter bot detection with reduced feature set. In *2020 IEEE International Conference on Intelligence and Security Informatics (ISI)*, 1–6. IEEE.
- Ali Alhosseini, S.; Bin Tareaf, R.; Najafi, P.; and Meinel, C. 2019. Detect me if you can: Spam bot detection using inductive representation learning. In *Companion proceedings of the 2019 world wide web conference*, 148–153.
- Badawy, A.; Ferrara, E.; and Lerman, K. 2018. Analyzing the digital traces of political manipulation: The 2016 Russian interference Twitter campaign. In *2018 IEEE/ACM international conference on advances in social networks analysis and mining (ASONAM)*, 258–265. IEEE.
- Beskow, D. M.; and Carley, K. M. 2019. Its all in a name: detecting and labeling bots by their name. *Computational and mathematical organization theory*, 25: 24–35.
- Beskow, D. M.; and Carley, K. M. 2020. You are known by your friends: Leveraging network metrics for bot detection in twitter. *Open Source Intelligence and Cyber Crime: Social Media Analytics*, 53–88.
- Boser, B. E.; Guyon, I. M.; and Vapnik, V. N. 1992. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, 144–152.
- Bradshaw, S.; Kollanyi, B.; Desigaud, C.; and Bolsover, G. 2017. Junk news and bots during the French presidential election: What are French voters sharing over Twitter?
- Breiman, L. 2001. Random forests. *Machine learning*, 45: 5–32.
- Broniatowski, D. A.; Jamison, A. M.; Qi, S.; AlKulaib, L.; Chen, T.; Benton, A.; Quinn, S. C.; and Dredze, M. 2018. Weaponized health communication: Twitter bots and Russian trolls amplify the vaccine debate. *American journal of public health*, 108(10): 1378–1384.
- Cai, C.; Li, L.; and Zengi, D. 2017. Behavior enhanced deep bot detection in social media. In *2017 IEEE International Conference on Intelligence and Security Informatics (ISI)*, 128–130. IEEE.
- Castillo, S.; Allende-Cid, H.; Palma, W.; Alfaro, R.; Ramos, H. S.; Gonzalez, C.; Elortegui, C.; and Santander, P. 2019. Detection of Bots and Cyborgs in Twitter: A study on the Chilean Presidential Election in 2017. In *Social Computing and Social Media. Design, Human Behavior and Analytics: 11th International Conference, SCSM 2019, Held as Part of the 21st HCI International Conference, HCII 2019, Orlando, FL, USA, July 26-31, 2019, Proceedings, Part I 21*, 311–323. Springer.
- Chavoshi, N.; Hamooni, H.; and Mueen, A. 2016. Debot: Twitter bot detection via warped correlation. In *Icdm*, volume 18, 28–65.
- Chen, T.; He, T.; Benesty, M.; Khotilovich, V.; Tang, Y.; Cho, H.; Chen, K.; Mitchell, R.; Cano, I.; Zhou, T.; et al. 2015. Xgboost: extreme gradient boosting. *R package version 0.4-2*, 1(4): 1–4.
- Cresci, S. 2020. A decade of social bot detection. *Communications of the ACM*, 63(10): 72–83.
- Cresci, S.; Di Pietro, R.; Petrocchi, M.; Spognardi, A.; and Tesconi, M. 2015. Fame for sale: Efficient detection of fake Twitter followers. *Decision Support Systems*, 80: 56–71.
- Cresci, S.; Di Pietro, R.; Petrocchi, M.; Spognardi, A.; and Tesconi, M. 2016. DNA-inspired online behavioral modeling and its application to spambot detection. *IEEE Intelligent Systems*, 31(5): 58–64.
- Cresci, S.; Di Pietro, R.; Petrocchi, M.; Spognardi, A.; and Tesconi, M. 2017a. 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*, 963–972.
- Cresci, S.; Di Pietro, R.; Petrocchi, M.; Spognardi, A.; and Tesconi, M. 2017b. Social fingerprinting: detection of spambot groups through DNA-inspired behavioral modeling. *IEEE Transactions on Dependable and Secure Computing*, 15(4): 561–576.
- Cresci, S.; Lillo, F.; Regoli, D.; Tardelli, S.; and Tesconi, M. 2018. FAKE: Evidence of spam and bot activity in stock microblogs on Twitter. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 12.
- Cresci, S.; Lillo, F.; Regoli, D.; Tardelli, S.; and Tesconi, M. 2019. Cashtag piggybacking: Uncovering spam and bot activity in stock microblogs on Twitter. *ACM Transactions on the Web (TWEB)*, 13(2): 1–27.
- Dehghan, A.; Siuta, K.; Skorupka, A.; Dubey, A.; Betlen, A.; Miller, D.; Xu, W.; Kamiński, B.; and Prafat, P. 2023. Detecting bots in social-networks using node and structural embeddings. *Journal of Big Data*, 10(1): 1–37.
- Dukić, D.; Keča, D.; and Stipić, D. 2020. Are you human? Detecting bots on Twitter Using BERT. In *2020 IEEE 7th International Conference on Data Science and Advanced Analytics (DSAA)*, 631–636. IEEE.
- Echeverria, J.; De Cristofaro, E.; Kourtellis, N.; Leontiadis, I.; Stringhini, G.; and Zhou, S. 2018. LOBO: Evaluation of generalization deficiencies in Twitter bot classifiers. In *Proceedings of the 34th annual computer security applications conference*, 137–146.
- Efthimion, P. G.; Payne, S.; and Proferes, N. 2018. Supervised machine learning bot detection techniques to identify social twitter bots. *SMU Data Science Review*, 1(2): 5.
- Feng, S.; Tan, Z.; Li, R.; and Luo, M. 2022a. Heterogeneity-aware twitter bot detection with relational graph transformers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 3977–3985.
- Feng, S.; Tan, Z.; Wan, H.; Wang, N.; Chen, Z.; Zhang, B.; Zheng, Q.; Zhang, W.; Lei, Z.; Yang, S.; et al. 2022b. TwiBot-22: Towards graph-based Twitter bot detection. *Advances in Neural Information Processing Systems*, 35: 35254–35269.
- Feng, S.; Wan, H.; Wang, N.; Li, J.; and Luo, M. 2021a. Satar: A self-supervised approach to twitter account representation learning and its application in bot detection. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 3808–3817.



- Feng, S.; Wan, H.; Wang, N.; Li, J.; and Luo, M. 2021b. Twibot-20: A comprehensive twitter bot detection benchmark. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 4485–4494.
- Feng, S.; Wan, H.; Wang, N.; and Luo, M. 2021c. 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*, 236–239.
- Fernquist, J.; Kaati, L.; and Schroeder, R. 2018. Political bots and the Swedish general election. In *2018 IEEE International Conference on Intelligence and Security Informatics (ISI)*, 124–129. IEEE.
- Ferrara, E. 2020. What types of COVID-19 conspiracies are populated by Twitter bots? *arXiv preprint arXiv:2004.09531*.
- Gilani, Z.; Farahbakhsh, R.; Tyson, G.; Wang, L.; and Crowcroft, J. 2017. Of bots and humans (on twitter). In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*, 349–354.
- Golovchenko, Y.; Buntain, C.; Eady, G.; Brown, M. A.; and Tucker, J. A. 2020. Cross-platform state propaganda: Russian trolls on Twitter and YouTube during the 2016 US presidential election. *The International Journal of Press/Politics*, 25(3): 357–389.
- Guo, Q.; Xie, H.; Li, Y.; Ma, W.; and Zhang, C. 2021a. Social bots detection via fusing bert and graph convolutional networks. *Symmetry*, 14(1): 30.
- Guo, W.; Su, R.; Tan, R.; Guo, H.; Zhang, Y.; Liu, Z.; Tang, R.; and He, X. 2021b. Dual graph enhanced embedding neural network for CTR prediction. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, 496–504.
- Hayawi, K.; Mathew, S.; Venugopal, N.; Masud, M. M.; and Ho, P.-H. 2022. DeeProBot: a hybrid deep neural network model for social bot detection based on user profile data. *Social Network Analysis and Mining*, 12(1): 43.
- Howard, P. N.; Kollanyi, B.; and Woolley, S. 2016. Bots and Automation over Twitter during the US Election. *Computational propaganda project: Working paper series*, 21(8).
- Kantepe, M.; and Ganiz, M. C. 2017. Preprocessing framework for Twitter bot detection. In *2017 International conference on computer science and engineering (ubmk)*, 630–634. IEEE.
- Kipf, T. N.; and Welling, M. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.
- Knauth, J. 2019. Language-agnostic twitter-bot detection. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)*, 550–558.
- Kouvela, M.; Dimitriadis, I.; and Vakali, A. 2020. Bot-Detective: An explainable Twitter bot detection service with crowdsourcing functionalities. In *Proceedings of the 12th International Conference on Management of Digital EcoSystems*, 55–63.
- Kudugunta, S.; and Ferrara, E. 2018. Deep neural networks for bot detection. *Information Sciences*, 467: 312–322.
- Lee, K.; Eoff, B.; and Caverlee, J. 2011. Seven months with the devils: A long-term study of content polluters on twitter. In *Proceedings of the international AAAI conference on web and social media*, volume 5, 185–192.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Lv, Q.; Ding, M.; Liu, Q.; Chen, Y.; Feng, W.; He, S.; Zhou, C.; Jiang, J.; Dong, Y.; and Tang, J. 2021. Are we really making much progress? revisiting, benchmarking and refining heterogeneous graph neural networks. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, 1150–1160.
- Magelinski, T.; Beskow, D.; and Carley, K. M. 2020. Graph-hist: Graph classification from latent feature histograms with application to bot detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 5134–5141.
- Mazza, M.; Cresci, S.; Avvenuti, M.; Quattrociocchi, W.; and Tesconi, M. 2019. Rtbust: Exploiting temporal patterns for botnet detection on twitter. In *Proceedings of the 10th ACM conference on web science*, 183–192.
- Miller, Z.; Dickinson, B.; Deitrick, W.; Hu, W.; and Wang, A. H. 2014. Twitter spammer detection using data stream clustering. *Information Sciences*, 260: 64–73.
- Moghaddam, S. H.; and Abbaspour, M. 2022. Friendship preference: Scalable and robust category of features for social bot detection. *IEEE Transactions on Dependable and Secure Computing*, 20(2): 1516–1528.
- Morstatter, F.; Wu, L.; Nazer, T. H.; Carley, K. M.; and Liu, H. 2016. A new approach to bot detection: striking the balance between precision and recall. In *2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 533–540. IEEE.
- Neudert, L.; Kollanyi, B.; and Howard, P. N. 2017. Junk news and bots during the German parliamentary election: What are German voters sharing over Twitter?
- Pastor-Galindo, J.; Zago, M.; Nespoli, P.; Bernal, S. L.; Celdrán, A. H.; Pérez, M. G.; Ruipérez-Valiente, J. A.; Pérez, G. M.; and Mármol, F. G. 2020. Spotting political social bots in Twitter: A use case of the 2019 Spanish general election. *IEEE Transactions on Network and Service Management*, 17(4): 2156–2170.
- Pham, P.; Nguyen, L. T.; Vo, B.; and Yun, U. 2022. Bot2Vec: A general approach of intra-community oriented representation learning for bot detection in different types of social networks. *Information Systems*, 103: 101771.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1): 5485–5551.

- Rodríguez-Ruiz, J.; Mata-Sánchez, J. I.; Monroy, R.; Loyola-Gonzalez, O.; and López-Cuevas, A. 2020. A one-class classification approach for bot detection on Twitter. *Computers & Security*, 91: 101715.
- Rossi, S.; Rossi, M.; Upreti, B. R.; and Liu, Y. 2020. Detecting political bots on Twitter during the 2019 Finnish parliamentary election. In *Annual Hawaii International Conference on System Sciences*, 2430–2439. Hawaii International Conference on System Sciences.
- Shahi, G. K.; Dirkson, A.; and Majchrzak, T. A. 2021. An exploratory study of COVID-19 misinformation on Twitter. *Online social networks and media*, 22: 100104.
- Shapley, L. S. 1953. A value for n-person games. *Contrib. Theory Games*, 2: 307–317.
- Shevtsov, A.; Antonakaki, D.; Lamprou, I.; Kontogiorgakis, I.; Pratikakis, P.; and Ioannidis, S. 2024. Russo-Ukrainian War: Prediction and explanation of Twitter suspension. In *Proceedings of the 2023 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, 348–355.
- Shevtsov, A.; Oikonomidou, M.; Antonakaki, D.; Pratikakis, P.; and Ioannidis, S. 2023. What Tweets and YouTube comments have in common? Sentiment and graph analysis on data related to US elections 2020. *Plos one*, 18(1): e0270542.
- Shevtsov, A.; Tzagkarakis, C.; Antonakaki, D.; and Ioannidis, S. 2022a. Identification of Twitter Bots Based on an Explainable Machine Learning Framework: The US 2020 Elections Case Study. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 16, 956–967.
- Shevtsov, A.; Tzagkarakis, C.; Antonakaki, D.; Pratikakis, P.; and Ioannidis, S. 2022b. Twitter dataset on the Russo-Ukrainian war. *arXiv preprint arXiv:2204.08530*.
- Twitter. 2023. Twitter API price list. <https://developer.twitter.com/en/products/twitter-api>. Accessed: 2023-09-09.
- Varol, O.; Ferrara, E.; Davis, C.; Menczer, F.; and Flammini, A. 2017. Online human-bot interactions: Detection, estimation, and characterization. In *Proceedings of the international AAAI conference on web and social media*, volume 11, 280–289.
- Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; and Bengio, Y. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*.
- Weï, F.; and Nguyen, U. T. 2019. 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)*, 101–109. IEEE.
- Yang, C.; Harkreader, R.; and Gu, G. 2013. Empirical evaluation and new design for fighting evolving twitter spammers. *IEEE Transactions on Information Forensics and Security*, 8(8): 1280–1293.
- Yang, K.-C.; Ferrara, E.; and Menczer, F. 2022. Botometer 101: Social bot practicum for computational social scientists. *Journal of Computational Social Science*, 5(2): 1511–1528.
- Yang, K.-C.; Pierri, F.; Hui, P.-M.; Axelrod, D.; Torres-Lugo, C.; Bryden, J.; and Menczer, F. 2021. The COVID-19 infodemic: twitter versus facebook. *Big Data & Society*, 8(1): 20539517211013861.
- Yang, K.-C.; Varol, O.; Davis, C. A.; Ferrara, E.; Flammini, A.; and Menczer, F. 2019. Arming the public with artificial intelligence to counter social bots. *Human Behavior and Emerging Technologies*, 1(1): 48–61.
- Yang, K.-C.; Varol, O.; Hui, P.-M.; and Menczer, F. 2020. Scalable and generalizable social bot detection through data selection. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 1096–1103.