Socialbots on Fire: Modeling Adversarial Behaviors of Socialbots via Multi-Agent Hierarchical Reinforcement Learning

Thai Le tql3@psu.edu Penn State University, USA Long Tran-Thanh long.tran-thanh@warwick.ac.uk University of Warwick, UK

Dongwon Lee dongwon@psu.edu Penn State University, USA

ABSTRACT

Socialbots are software-driven user accounts on social platforms, acting autonomously (mimicking human behavior), with the aims to influence the opinions of other users or spread targeted misinformation for particular goals. As socialbots undermine the ecosystem of social platforms, they are often considered harmful. As such, there have been several computational efforts to auto-detect the socialbots. However, to our best knowledge, the adversarial nature of these socialbots has not yet been studied. This begs a question "can adversaries, controlling socialbots, exploit AI techniques to their advantage?" To this question, we successfully demonstrate that indeed it is possible for adversaries to exploit computational learning mechanism such as reinforcement learning (RL) to maximize the influence of socialbots while avoiding being detected. We first formulate the adversarial socialbot learning as a cooperative game between two functional hierarchical RL agents. While one agent curates a sequence of activities that can avoid the detection, the other agent aims to maximize network influence by selectively connecting with right users. Our proposed policy networks train with a vast amount of synthetic graphs and generalize better than baselines on unseen real-life graphs both in terms of maximizing network influence (up to +18%) and sustainable stealthiness (up to +40% undetectability) under a strong bot detector (with 90% detection accuracy). During inference, the complexity of our approach scales linearly, independent of a network's structure and the virality of news. This makes our approach a practical adversarial attack when deployed in a real-life setting.

CCS CONCEPTS

• Computing methodologies \rightarrow Sequential decision making.

KEYWORDS

socialbot, social bot, adversarial, reinforcement learning

ACM Reference Format:

Thai Le, Long Tran-Thanh, and Dongwon Lee. 2022. Socialbots on Fire: Modeling Adversarial Behaviors of Socialbots via Multi-Agent Hierarchical Reinforcement Learning. In *Proceedings of the ACM Web Conference 2022* (*WWW '22*), April 25–29, 2022, Virtual Event, Lyon, France. ACM, New York, NY, USA, 10 pages. https://doi.org/10.1145/3485447.3512215

WWW '22, April 25–29, 2022, Virtual Event, Lyon, France © 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9096-5/22/04...\$15.00 https://doi.org/10.1145/3485447.3512215

1 INTRODUCTION

Socialbots refer to automated user accounts on social platforms that attempt to behave like real human accounts, often controlled by either automatic software, human, or a combination of both-i.e., cyborgs [4]. Different from traditional spambots, which may not have proper profiles or can be easily differentiated from regular accounts, socialbots often mimic the profiles and behaviors of real-life users by using a stolen profile picture or biography, building legitimate followships, replying to others, etc. [4]. Socialbots are often blamed for spreading divisive messages-e.g., hate speech, disinformation, and other low-credibility contents that have been shown to widen political divides and distrust among both online and offline communities [4, 20, 30]. To mitigate such harmful proliferation of socialbots, therefore, there has been extensive research, most of which focus on how to effectively detect them [10, 37, 54]. However, these works usually follow the cat-and-mouse game where they passively wait for socialbot evasion to happen before they can react and develop a suitable detector [8]. Instead of following such a reactive scheme, however, proactively modeling socialbots and their adversarial behaviors on social platforms can better advance the next bot detection research.

In particular, we pose a question "<u>Can socialbots exploit compu-</u> tational learning mechanism such as reinforcement learning to their <u>advantage?</u>" To our best knowledge, *adversarial* nature of socialbots has not yet been fully explored and studied. However, it is plausible that adversaries who own a farm of socialbots operate their socialbots according to certain strategies (or algorithms). Therefore, proactively simulating such a computational learning mechanism and understanding adversarial aspect of socialbots better would greatly benefit future research on socialbot detection.

In general, a socialbot has two main objectives that are adversarial in nature: (i) to facilitate mass propaganda propagation through social networks and (ii) to evade and survive under socialbot detectors. The first goal can be modeled as an NP-Hard influence maximization (IM) problem [25] where the bot needs to build up its network of followers-i.e., seed users, overtime such that any new messages propagated from the bot through these users can effectively spread out and influence many other people. Simultaneously, it also needs to systematically constrain its online behaviors such that it will not easily expose itself to socialbot detectors. Although the IM problem has been widely studied by several works [3, 24, 25, 27], they only focus on maximizing the network influence given a fixed and static budget # of seed nodes (that is relatively small) and they assume that every node is equally acquirable. However, these assumptions are not practical in our context. Not only a socialbot needs to continuously select the next best seed node or follower over a long temporal horizon-i.e., potentially large budget of seed nodes, it also needs to consider that gaining the followship from

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

a very influential actor–e.g., Elon Musk, is practically much more challenging than from a normal user. At the same time, a socialbot that optimizes its network of followers must also refrain from making suspicious behaviors–e.g., constantly following others, that can trigger the attention of bot detectors. Thus, learning how to navigate a socialbot is a very practical yet challenging task with two intertwined goals that cannot be separately optimized. Toward this challenge, in this paper, we formulate the **Adversarial Socialbot Learning (ASL)** problem and design a multi-agent *hierarchical*

reinforcement learning (HRL) framework to tackle it.

3Our main contributions are as follows.

- First, we formulate a novel ASL problem as an optimization problem with constraints.
- Second, we propose a solution to the ASL problem by framing it as a cooperative game of two HRL agents that represent two distinctive functions of a socialbot, namely (i) selecting the next best activity-e.g., tweet, retweet, reply, mention, and (ii) selecting the next best follower. We carefully design the RL agents and exploit unsupervised graph representation learning to minimize the potential computational cost resulted from a long time horizon and a large graph structure.
- Third, we demonstrate that such RL agents can *learn from synthetic graphs yet generalize well on real unseen graphs*. Specifically, our experiments on a real-life dataset show that the learned socialbot outperforms baselines in terms of influence maximization while sustaining its longevity by continuously evading a strong black-box socialbot detector of 90% detection accuracy. During inference, in addition, the complexity of our approach scales linearly and is independent of a network's structure and the virality of news.
- Four, we release an environment under the Open AI's *gym* [1] library. This enables researchers to simulate various adversarial behaviors of socialbots and develop novel bot detectors in a *proactive* manner.

2 RELATED WORK

2.1 Socialbots Detection

The majority of previous computational works on socialbots within the last decade [2, 10, 37, 39, 42, 52, 54] primarily focus on developing computer models to effectively detect bots on social networks [4, 8]. These models are usually trained on a ground truth dataset using supervised learning algorithms-e.g., Random Forest, Decision Tree, SVM, to classify an individual social media account into a binary label-i.e., bot or legitimate [4]. Moreover, these learning algorithms usually depend on either a set of statistical engineered predictive features such as the number of followers, tweeting frequency, etc. [5, 42, 54], or a deep learning network where the features are automatically learned from unstructured data such as an account's description text. Even though there are many possible features that can be used to detect socialbots, statistical features that can be directly extracted from user metadata provided by official APIs-e.g., Twitter API, are more practical due to their favorable computational speed in practice [54]. In fact, many of the features that are utilized by the popular socialbot detection API botometer fall into this category. Moreover, we later also show that using simple statistical features derived from user metadata can help

train a socialbot detector with around 90% prediction accuracy on a hold-out test set (Sec. 3.1). Regardless of how a socialbot detector extracts its predictive features, they are mainly designed following a *reactive* schema where they learn how to detect socialbots after they appear (thus a training dataset can be collected).

2.2 Adversarial Socialbot Learning

While previous works help us to understand better the detection aspect of socialbots, the learning aspect of them has not been widely studied [8]. Distinguished from learning how to detect socialbots using a stationary snapshot of their features, ASL computationally models the adversarial learning behaviors of socialbots over time. To the best of our knowledge, relevant works on this task are limited to [7]. This work adopts an evolution optimization algorithm to find different adversarial permutations from a *fixed* socialbot' encoded activity sequence-e.g., "tweet→tweet→retweet→reply,...", and examine if such permutations can help improve the detection accuracy of a bot detector. However, such permutations, even though adversarial in nature, are just static snapshots of a socialbot and do not tell a whole story on how the bot evolves. In other words, we are still lacking a general computation framework that models the temporal dynamics of socialbots and their adversarial behaviors. Therefore, this paper aims to formally formulate their behaviors as a Markov Decision Process (MDP) [21] and designs an RL framework to train socialbots that can optimize their adversarial goals on real-life networks.

We investigate two adversarial objectives of a socialbot: influencing people while evading socialbot detection. While the first one can be modeled as an IM task on graph networks, traditional IM algorithms–e.g.,[3, 25, 27], assume that the number of seed nodes is relatively small and all nodes are equally acquirable, all of which are not applicable in the socialbot context as previously described. There have been also a few works–e.g., [33, 49], that utilizes RL to IM task. Yet their scope is still limited to a single constraint on the budget number of seeds. Influence maximization under a temporal constraint–i.e., not to be detected lead to early termination in this case, is a non-trivial problem.

3 PROBLEM FORMULATION

3.1 Social Network Environment

Network Representation and Influence Diffusion Model A social network includes users, their interactions and how they influence each other. We model this network as a directed graph G=(V, E). An edge between two users $u, v \in V$, denoted as $(u, v) \in E$, means u can have influence on v. (u, v) also illustrates a piece of news can spread from u to v-i.e, v follows u (thus u influences v).

As there is no influence model that can perfectly reflect realworld behaviors, to model the influence flow through *G*, we adopt *Independence Cascade Model (ICM)* [16, 17], which is the most commonly used in the context of a social network [22, 26, 34]. ICM was originally proposed to model the "word-of-mouth" behaviors, which resemble the information sharing phenomena online well. In ICM, a node is either active or inactive. Once a node *u* is activated, it has a single opportunity to activate or influence its inactive neighbors $\mathcal{N}(u)$ with an *uniform activation probability p*. At first, every node is inactive except a set of seed nodes *S*. After that, as the

Table 1: Predictive features of the social bot detector \mathcal{F} .

Feature	Description				
#tweets	# of tweets posted by the user				
#replies	# of replies posted by the user				
#retweets	# of retweets posted by the user				
#avg.tweets	average # tweets posted per timestep				
#avg.replies	average # replies posted per timestep				
#avg.retweets	average # retweets posted per timestep				
#retweet.ratio	#retweets/#tweets				
#replies.ratio	#replies/#tweets				
#retweet.replies.ratio	#retweets/#replies				
#mentions.ratio	# unique mentions posted per tweet				

environment rolls out throughout a sequence of discrete timesteps, the influence will propagate from S through the network by activating different nodes in G following E and p. The process ends when there is no additional activated nodes being activated [24, 32]. Hence p is also the virality of news–i.e., how fast a piece of news can travel through G. We then use G=(V, E, p) to denote the social network G.

Let denote by $\sigma(S, G)$ the *spread function* that measures how many nodes in *G* a piece of information–e.g., fake news, can spread from *S* via the ICM model. Given a fixed network structure (*V*, *E*) and the news virality *p*, different *S* will result in different values of $\sigma(S, G)$. Hence, selecting a good *S* is decisive in optimizing the spread of influence on *G*. However, choosing *S* to maximize $\sigma(S, G)$ has already been proven to be an NP-Hard problem [25].

Socialbots. A socialbot is then a vertex in *G* that attempts to mimic human behaviors for various aims–e.g., spreading propaganda or low-credible contents through *G*, [4, 44, 46]. It carries out a sequence of activities \mathcal{A} to *simultaneously* achieve two main objectives:

- **Obj. 1:** Optimizing its influence over G by selectively collecting good seed nodes–i.e., followers, $S \in V$, over time
- **Obj. 2:** Evading bots detectors-i.e., not to be detected and removed

These two goals are often in tension in that improving Obj 1 typically hurts Obj 2 and vice versa. That is while having a good network of followers S enables a socialbot to spread disinformation to a large number of users at any time, having a high undetectability helps it to sustain this advantage over a long period. As socialbots are usually deployed in groups, and later coming socialbots can also easily inherit a previously established network of followers S of a current one. If a bot is detected and removed from G, not only it can lose its followers S and expose itself to be used to develop stronger detectors, it can also risk revealing the identity of other bots–e.g., by way of guilt-by-association [50]. This makes the sustainability achieved through *Obj 2* distinguishably important from previous literature–e.g., [24, 25, 51], where the optimization of S plays a more central role.

Relationship between \mathcal{A} and \mathcal{S} . \mathcal{A} denotes the activity sequencei.e., the DNA of the bot [6]. \mathcal{A} includes four possible types of actions to be made at every timestep *t*, namely *tweet*, *retweet*, *reply or mention*, and only the *last three* of which can directly interact with others to expand \mathcal{S} . Despite these actions are in the *Twitter* context, other platforms also provide similar functions–e.g., tweet->post, retweet->share, reply->comment, mention->tag on Facebook. In practice, *not* every node requires an equal effort to convert to a follower. For example, a bot needs to accumulate its reputability over time and interact more frequently to have an influencer–e.g., Elon Musk, rather than a normal user to become its follower. Since a real model underlining such observation is unknown, we model it using a simple heuristic:

$$g_Q(u,t) = \max(1, Qf(u,t)) \quad \text{where}$$

$$f(u,t) \sim \text{Bernoulli}(1 - \frac{1 + |S_t|}{1 + |\mathcal{N}(u)|}), \quad (1)$$

where $g_Q(u, t)$ with hyper-parameter $Q \ge 1$, is the **number of times** the socialbot is *required* by the environment to continuously interact with an influencer u-i.e., high $\mathcal{N}(u)$, for it to become a follower at t. Intuitively, a bot with a good reputation overtime-i.e., a high number of followers at the timestep t-i.e., $|S_t|$, can influence others to follow itself more effortlessly than a newly created bot. Overall, \mathcal{A} encodes when and what type of interaction-i.e., *retweet, reply* or *mention*, to use to acquire a new follower $s \in S$, sthen decides the *frequency* of such interaction in \mathcal{A} . Thus, \mathcal{A} and S is temporally co-dependent.

Socialbot Detection Model. Bot detectors are responsible for detecting and removing socialbots from G. Let $\mathcal{F}(\mathcal{A}_t) \in \{0, 1\}$ denote a model that predicts whether or not an account is a socialbot based on its activity sequence up to the timestep $t(\mathcal{A}_t)$. This sequence of ordered activity is then usually represented as an unordered list of statistical features such as number of replies, tweets per day, by socialbot detectors [10, 37, 54]. In this paper, \mathcal{F} extracts and adopts several features (Table 1) from previous works for detection. Most of the features are utilized by the popular bot detection API Botometer [9]. We train \mathcal{F} using the Random Forest [47] algorithm with supervised learning on a *publicly available dataset* [36, 53]¹ of nearly 15K Twitter accounts, half of which is labelled as socialbots. This dataset is not exposed to the socialbots. Here we also assume that $\mathcal{F}(\cdot)$ is a *black-box model*-i.e., we do not have access to its parameters. $\mathcal F$ achieves nearly 90% in F1 score on an unseen test set following the standard 5-fold cross validation (train and test with 80%/20% data). Since \mathcal{A} and \mathcal{S} are co-dependent, we can easily see that S also has effects on the detectability of a social bot. Note that to focus on the study of the adversarial aspect of socialbots, we had to resort to a certain combination of account features and the socialbot detection model. 90% in F1 score is also in line with SOTA detectors on a similar set of features [37].

3.2 The ASL Problem and Objective Function

From the above analysis, this paper proposes to study the Adversarial Socialbot Learning (ASL) problem to achieve both *Obj 1* and *Obj 2*. In other words, we aim to solve the following problem.

PROBLEM: Adversarial Socialbot Learning (ASL) aims to develop an automatic socialbot that can exercise adversarial behaviors against a black-box bot detector \mathcal{F} while at the same time maximizing its influence on *G* through a set of selective followers \mathcal{S} .

¹https://botometer.osome.iu.edu/bot-repository/



Figure 1: An example of ACORN HRL framework. As the environment rolls out, AGENTI (π_1) decides which type of activity (T, R, A or M) to perform. Whenever an interactive action (R, A, M) is selected, AGENTII (π_2) then selects a new follower. Since the selected user u_4 at t_4 is an influencer, π_2 needs perform not once but Q=3 times of action "A" to acquire u_4 (blue arrow). Whenever $|\mathcal{A}|$ reaches an interval of K=7, the bot detector $\mathcal{F}(\mathcal{A}_t)$ is triggered (red arrow).

Specifically, we formulate this task as an optimization problem with the objective function as follows.

OBJECTIVE FUNCTION: Given a black-box bot detection model \mathcal{F} and a social network environment what is characterized by G=(V, E, p), K, Q, we want to optimize the objective function:

$$\max_{\mathcal{S}_t, \mathcal{A}_t} \quad \boldsymbol{R}^* = \boldsymbol{\sigma}(\mathcal{S}_{T^*}, G)(1 + T^*) \quad \text{subject to} \qquad (2a)$$

$$T^* = \min_{T^*} \left[\mathcal{F}(\mathcal{A}_{T^*}) = 1 \land \mathcal{F}(\mathcal{A}_t) = 0 \right]$$
(2b)

$$\forall 1 < t < T^*, |\mathcal{A}_t| \mod K = |\mathcal{A}_{T^*}| \mod K = 0$$
 (2c)

$$g_Q(u,t) = max(1,Qf(u,t)) \quad \forall \ 1 < t < T^*$$
(2d)

Socialbot detector \mathcal{F} can run prediction on the socialbot every time it performs a new activity. However, \mathcal{A}_u and |V| can potentially be very large. Thus, we assume that \mathcal{F} only runs detection every time K new activities is added to \mathcal{A} (Eqn. 2b). This makes T^* the earliest interval timestep at which a socialbot is detected and removed by \mathcal{F} (Eqn. 2b,c). Since \mathbb{R}^* is *monotonically increasing* on both $V \ge \sigma(\mathcal{S}_b, G) \ge 0$ and $T^* \ge 1$, to maximize \mathbb{R}^* , a socialbot cannot focus *only* either on *Obj 1* or *Obj 2*. In other words, Eqn. (2d) encourages the socialbot to simultaneously optimize both objectives.

4 THE PROPOSED METHOD: ACORN

4.1 Markov Decision Process Formulation

The ASL problem can be formulated as an MDP process which consists of a state set *S*, an action set *A*, a transition function \mathcal{P} , a reward function *R*, a discount factor $\gamma \in [0, 1]$ and the horizon *T*. Since the space requirement for *A* can be very large–i.e., 4|V| for 4 possible activities and |V| possible seed nodes, especially on a large network, this can make the task much more challenging to optimize due to potential sparse reward problem. To overcome this, we transformed this into a HRL framework of two functional agents, AGENTI and AGENTII, with a *global reward* (Figure 1). We call this ACORN (Adversarial so<u>CialbOts leaRniNg</u>) framework. While AGENTI is responsible for deciding which *type* of activity among {tweet, retweet, reply, mention} to perform at each timestep

t, AGENTII is mainly responsible for S-i.e., to select which follower to accumulate, *only when* AGENTI chooses to do so-i.e., retweet, reply, mention. This reduces the overall space of A to only |V|+4. Since \mathcal{A} and \mathcal{S} are co-dependent (Sec. 3.1), the two agents need to continuously cooperate to optimize both influence maximization and undetectability. It is noted that the Markov assumption behind this MDP is not violated because both influence function $\sigma(\cdot)$ and detection probability \mathcal{F} at time t only depends on statistical snapshot of the two agents at t-1. This HRL task is then described in detail as follows.

State. Following [12, 29, 35], we assume that the state space *S* can be factorized into bot-specific S_{DNA} and network-specific S_{ENV} , and $s^I \in S_{\text{DNA}}$, $s^{II} \in S_{\text{ENV}}$, where s^I , s^{II} is the state space of AGENTI and AGENTII, respectively. Specifically, s^I_t encodes (i) the number of followers |S| of the bot and (ii) a snapshot of \mathcal{A}_t at timestep *t*. While s^I_t can directly store the actual \mathcal{A}_t sequence, this potentially induces a computational and space overhead especially when *t* becomes very large. Instead, we compact \mathcal{A}_t into a fixed vector summarizing the frequency of each *tweet*, *retweet*, *reply*, and *mention* action up to *t*. This effectively limits the space complexity of $s^I \in \mathbb{R}^5$ to O(1). Similarly, $s^{II}_t \in \mathbb{R}^{4+|V|(k+1)}$ comprises of (i) node2vec(*G*) [19] which encodes the structure of *G* to |V| vectors of size *k*, (ii) a statistical snapshot of \mathcal{A}_t and (iii) information regarding \mathcal{S}_t , encoded as:

$$(\mathbb{1}(u \notin S_t) \frac{1 + |S_t|}{1 + |\mathcal{N}(u)|})_{u=0}^{|V|} \in \mathbb{R}^{|V|}$$
(3)

Previous works have often encoded the network structures ([15, 55]) via a parameterized Graph Neural Network (GCN) [28] as part of the policy network. As this approach requires frequent parameter updates during training, instead, we adopt *node2vec*(*G*) as an alternative unsupervised method which requires the calculation *only once*. While S_t can be encoded as a one-hot vector $(\mathbbm{1}(u \notin S_t))_{u=0}^{|V|}$, we enrich it by multiplying it with the binary f(u, t) condition $\frac{1+|S_t|}{1+|N(u)|}$ (Sec. 3.1), which then results in Eq. (3). This enables AGEN-TII to select nodes accordingly with the current reputation of the bot $|S_t|$.

Socialbots on Fire: Modeling Adversarial Behaviors of Socialbots via Multi-Agent Hierarchical Reinforcement Learning WWW '22, April 25-29, 2022, Virtual Event, Lyon, France



Figure 2: Examples of a real (Left) and synthetic (Right) news propagation networks on Twitter with a similar star-like shape structure.

Action and Policy. Similarly, we factor A into two different action spaces a^{I}, a^{II} for AGENTI and AGENTII, respectively. $a^{I} \in \mathbb{R}^{4}$, $a^{II} \in \mathbb{R}^{|V|}$ are both encoded as one-hot vectors, representing one of four available activities and one of potential followers, respectively. We then have two policies $\pi_1 = (a^{I}|s^{I}), \pi_2 = (a^{II}|s^{II}, a^{I})$ that control AGENTI and AGENTII, respectively.

Reward. Even though we can directly reward the RL agents with $\sigma(S_t, G) \ge 1.0$ at every timestep $t \le T^*$, this calculation will incur large computational cost, especially when T^* becomes large. Instead, therefore, we design an accumulative reward function R that consists of a *step reward* and a *delayed reward* to incentivize RL agents to maximize \mathbf{R}^* (Eqn. 2) as follows.

$$\boldsymbol{R}_{\text{step}}(t) = \boldsymbol{\sigma}(\mathcal{S}_t \setminus \mathcal{S}_{t-1}, G)$$

$$\boldsymbol{R}_{\text{delaved}}(T^*) = \boldsymbol{\sigma}(\mathcal{S}_{T^*}, G)$$
(4)

where $T^* \leq T$ is the interval timestep at which the bot is detected and the episode is terminated. The step reward \mathbf{R}_{step} , which can be efficiently computed, is the *marginal* gain on the network influence given a new follower selected at *t*. Using the step reward with a discount factor, $\gamma_{step} < 1.0$, helps avoid the sparse reward problem and encourages good follower selection *early* during an episode. Since $\mathbf{R}_{step} \geq 1.0$, it also encourages the bot to survive against bot detection longer–i.e., to maximize T^* . In other words, as long as the socialbot survives–i.e., T^* increases, in other to make new friendship, it will be able to influence more people. However, since $\sigma(\cdot)$ is *subadditive*–i.e., $\sigma(\{u\}, G)+\sigma(\{v\}, G)\geq\sigma(\{u, v\}, G) \forall u, v \in V$, we then introduce the delayed reward $\mathbf{R}_{delayed}$ at the end of each episode with a discounted factor $\gamma_{delayed} < 1.0$ as a reward adjustments for each node selection step.

4.2 Parameterization

A policy network π_1 is a Multi-Layer Perceptron (MLP) followed by a softmax function that projects s^I to a probability distribution of 4 possible activities. We can then sample a^I from such a distribution. A policy network π_2 utilizes Convolutional Neural Network [23] (CNN) to efficiently extract useful *spatial* features from the stack of representation vectors of all vertex $u \in V$ calculated by *node2vec*(*G*) (Sec. 4.1), and MLP to extract features from the rest of the components of s^{II} . The resulted vectors are then concatenated as the final feature vector. Instead of directly projecting this feature on the



Figure 3: We generate synthetic networks that ensemble real networks' structures on the fly to train ACORN and test it with real networks.

original action space of a^{II} using an MLP, we adopt the parametricaction technique [14, 38] with invalid actions at each timestep *t*-i.e., already *chosen* node, being masked.

4.3 Learning Paradigm

Learning algorithm. We train π_1 , π_2 using the *actor-critic* Proximal Policy Optimization (PPO) algorithm [43]. It has a theoretical guarantee and is known to be versatile in various scenarios [15, 43, 45, 55]. The actor refers to π_1 and π_2 , as described above. Their critics share the same network structure but output a single scalar value as the estimated accumulated reward at *t*.

Learning on synthetic and evaluating on real networks. We evaluate our method on real world data. To make our RL model generalize well on unseen *real* networks (Figure 2, Left) with different possible configurations of G=(V, E), it is important to train our model on a sufficient number of diverse scenarios–i.e., training graphs. However, collecting such a train dataset often requires much time and efforts. Hence, we propose to train our model on synthetic graphs, which can be efficiently generated on the fly during the training [24]. To avoid distribution shifts between train and test graphs, we first collect a *seed dataset* of several news propagation networks and use their statistical properties (p_{intra}, p_{inter}) to *spontaneously* generate a *synthetic* graph (Figure 2, Right) for each training iteration. We describe this in detail in Section 5.

5 EXPERIMENT

5.1 Set-Up

Datasets. We collected a total of top-100 trending articles on Twitter from January 2021 to April 2021 and their corresponding propagation networks with a maximum of 1.5K nodes using the public Hoaxy API². All the downloaded data is free from user-identifiable information. The majority of these articles are relevant to the events surrounding the 2020 U.S. presidential election and the COVID-19 pandemic. We also share the same observation with previous literature [24, 41] such that retweet networks tend to have star-like shapes. These networks have a high p_{intra} and a low p_{inter} value, suggesting multiple separate star-shape communities with few connections among them. Therefore, viral news usually originates from a few very influential actors in social networks and quickly propagates to their followers.

Training and Testing Set. Figure 3 illustrates how to utilize synthetic data during training. Since we observe that our framework

²https://rapidapi.com/truthy/api/hoaxy



Figure 4: Performance comparison of a single socialbot under bot detection constraint.

generalizes better when trained with more complex graphs-i.e., more edges with high intra-community (p_{intra}) and inter-community (pinter) edge probabilities, We first selected 10% of the collected real networks with the highest *pintra* and *pinter* as initial seed graphse.g., Figure 2, Left, to generate the training set and use the rest as the test set. Then, during training, we used the average statistics (pintra, pinter, # of communities and their sizes) of the seed graphs to generate a stochastic, synthetic graph for each training episode of a maximum T timesteps-e.g., Figure 2, Right. These two statistics are selected because they well capture the star-like shapes of a typical retweet network. Since the real activation probabilities *p* of the collected networks are unknown, we found that using a fixed high p value during training achieves the best results. We then reported the averaged results across 5 different random seeds on the remaining 90 real test networks with varied p values and on a much longer horizon than T. Note that this number of testing networks is much larger and more extensive than those of previous studies [24, 25, 51].

Baselines. Since there are no previous works that address the ASL problem, we combined different approximation and heuristic approaches for the IM task with the socialbot detector evasion feature that is provided by *learned* AGENTI as baselines:

- AGENTI+C. This baseline extends the *Cost Effective Lazy Forward (CELF)* [31] and exploits the submodularity of the spread function $\sigma(\cdot)$ to become the first substantial improvement over the traditional GREEDY method [25] in terms of computational complexity. IM is a standard baseline in influence maximization literature.
- AGENTI+H. Since *G* consists of several star-like communities, we also used a *heuristic approach* DEGREE [3, 25] that always selects the node with the largest out-degree that is available–i.e., user with the largest # of followers.
- AGENTI*+C and AGENTI*+H train the first-level agent *independently* from the second-level agent and combined it with CELF or the heuristic approach DEGREE, respectively. These are introduced to examine the dependency between the trained AGENTI and AGENTII

Since the GREEDY approach does not scale well with a large number of seeds, however, we excluded it from our experiments.

Models and Configurations. We used a fixed hyper-parameter setting. During training, we set $K \leftarrow 20$, $Q \leftarrow 3$, $T \leftarrow 60$, $p \leftarrow 0.8$, and γ_{step} , $\gamma_{\text{delayed}} \leftarrow 0.99$. We refer the readers to the appendix for detailed configurations for RL agents. We ran all experiments on the machines with Ubuntu OS (v18.04), 20-Core Intel(R) Xeon(R) Silver 4114 CPU @ 2.20GHz, 93GB of RAM and a Titan Xp GPU 16GB. All implementations are written in Python (v3.8) with Pytorch (v1.5.1).

5.2 Main Results

Network Influence Ratio. Figure 4 shows the network influence ratio–i.e., network influence over total number of users, *under a bot detection environment* given different number of budget seeds |S| and p values:

$$\boldsymbol{\sigma}(\mathcal{S}, G) / |V| \le 1.0 \tag{5}$$

A high network influence ratio requires both (i) efficient followship selection and (ii) efficient detection evasion strategy. Overall, ACORN outperforms all baselines with different news virality (*p* values). However, ACORN underperforms when |S| is low–e.g., |S|=50 in Figure 4. This is because AGENTII learns not to connect with the most influential nodes early in the process. This can help prevent disrupting the sequence \mathcal{A} and lead to early detection, especially when it gets closer to the next prediction interval of \mathcal{F} .

The larger the *p* value, the further–i.e., more hoops, a news can propagate through *G*. Hence, as *p* increases–i.e., the more viral a piece of news, utilizing the network structure to make new connections is crucial and more effective than simply selecting the most influential users. This is reflected in the inferior performance of AGENTI+H when compared with AGENTI+C, ACORN in Figure 4, *p*=0.75. This means that ACORN is able to utilize the network structured capture by *node2vec* and postpone short-term incentives– i.e., makes friends with influential users, for the sake of long-term rewards. Overall, ACORN also behaves more predictably than baselines in terms of the influence ratio's deviation across several runs.

Survival Timesteps. We then evaluated if a trained socialbot can survive even after collecting all followers. Table 2 shows that while we train a socialbot with a finite horizon T=60, it can live on the network for a much longer period during testing. However, other baselines were detected very early. Since only three out of four activities–i.e., tweet, retweet, reply, and mention, allow to collect

Socialbots on Fire: Modeling Adversarial Behaviors of Socialbots via Multi-Agent Hierarchical Reinforcement Learning WWW '22, April 25-29, 2022, Virtual Event, Lyon, France



Figure 5: Empirical comparison of running time between CELF and AGENTII (ACORN).

Distributions of Degrees (# of Followers) of Selected Seed Nodes



Figure 6: Insights on the learned policies.

Table 2: Total survival timesteps v.s. network influence ratio after reaching |S| = |V|

	<i>p</i> = 0.25		p = 0.50		<i>p</i> = 0.75	
	%↑	Steps↑	%↑	Steps↑	%↑	Steps↑
AgentI+H	0.63 ± 0.43	$1.2\mathrm{K}\pm1\mathrm{K}$	0.68 ± 0.36	$1.2\text{K} \pm 1\text{K}$	0.73 ± 0.31	$1.2K \pm 1K$
AgentI+C	0.73 ± 0.41	1.5K±968	0.71 ± 0.36	1.3K±1K	0.77 ± 0.30	1.3K±1.1K
ACORN	$\textbf{0.99} \pm \textbf{0.10}$	$2.1K\pm254$	$\textbf{0.99} \pm \textbf{0.10}$	$\textbf{2.0K} \pm \textbf{276}$	$\textbf{0.99} \pm \textbf{0.10}$	$2.0K\pm305$

new followers, it is natural that socialbots need to survive much longer than |V| steps–e.g., around 2.0K in Table 2, to accumulate all followers. This corresponds to 98%, 64%, and 56% of socialbots surviving-i.e., not detected, after reaching |S|=V for ACORN, AGENTI+C and AGENTI+H, respectively. Our trained socialbot can also sustain much longer if we keep it going during testing, even with different detection intervals K>20. This implies that AGENTI can generalize its adversarial activities against ${\mathcal F}$ toward unseen real-life scenarios.

Dependency between RL Agents. The above results also demonstrate the effects of co-training AGENTI and AGENTII. First, the heuristic and CELF method when paired with the learned AGENTI (blue & green lines, Figure 4) performs much better than when paired with an independently trained (without AGENTII) AGENTI (yellow & black lines, Figure 4). This shows that AGENTI, when trained with AGENTII, becomes more versatile and can help a socialbot survive a much longer period of time, especially even when the socialbot only uses a heuristic node selection. However, AGENTI performs the best when paired with AGENTII. This shows that two RL agents successfully learn to collaborate, not only to evade the

socialbot detection but also to effectively maximize its network influence. This further reflects the co-dependency between the roles of \mathcal{A} and \mathcal{S} as analyzed in Sec. 3.1.

Computational Analysis. We compared the computational complexity of AGENTII specifically with the CELF algorithm during inference. Even though CELF significantly improves from the traditional GREEDY [25] IM algorithm with the computational complexity of O(|S||V|m) [48] (assuming each call of σ takes O(m) and only one round of Monte Carlo simulation is needed), its computation greatly depends on $\boldsymbol{\sigma}(\cdot)$, the size of the graph and becomes only computationally practical when |S| is small. This is also similar to other traditional IM algorithms such as CELF++ [18], TIM [48], and ASIM [13]. To illustrate, CELF takes much more time to compute as |S| increases especially with large *p*-i.e., more nodes need to be reached when computing $\sigma(\cdot)$ (Figure 5). However, with the O(1) complexity of the forward pass through π_2 , AGENTII is able to scale linearly O(|S|) regardless of the network structure and the virality of the news during inference. Even though our framework requires to calculate the graph representation using node2vec, it is specifically designed to be scalable to be able to process large graphs [40] and we only need to run it once.

Insights on the Learned Policies. We summarized the node selection strategies of all methods in Figure 6. We observed that both heuristic and CELF selects very influential nodes with many followers (high out-degrees) very early. Alternatively, AGENTII acquires an array of normal users (low out-degrees) before connecting with influential ones. This results in early detection and removal of the

WWW '22, April 25-29, 2022, Virtual Event, Lyon, France



Figure 7: Performance of multiple socialbots under bot detection constraint on a large network.

baselines and sustainable survival of our approach. This shows that AGENTII can learn to cope with the relationship constraint (Eqn. (1)) between \mathcal{A} and \mathcal{S} imposed by the environment. Moreover, the degrees of selected users by ACORN has a *right* long-tail distribution, which means that ACORN overall still tries to maximize its network influence early in the process.

5.3 Multiple Socialbots Results

We have evaluated our approach on different real-life news propagation graphs. These networks can be considered as sub-graphs of a much larger social network. In practice, different sub-graphs can represent different communities of special interests–e.g., politics, COVID-19 news, or different characteristics–e.g., political orientation. Since socialbots usually target to influence a specific group of users–e.g., anti-vaxxer, it is practical to deploy several bots working in tandem on different sub-graphs. To evaluate this scenario, we aggregated all 90 test sub-graphs into a large network of 135K nodes and used each learned socialbot for each sub-graph. Figure 7 shows that ACORN still outperforms other baselines especially later in the time horizon. Moreover, ACORN can efficiently scale to a real-life setting thanks to its linear running time and highly parallel architecture.

6 DISCUSSION AND LIMITATION

Our contribution goes beyond our demonstration such that one can train adversarial socialbots to effectively navigate real-life networks using an HRL framework. We will also publish a multi-agent RL environment for the ASL task under the gym library [1]. This environment will facilitate researchers to test different RL agents, examine and evaluate assumptions regarding the behaviors of socialbots, bot detection models, and the underlying influence diffusion models on synthetic and real-life news propagation networks. It remains a possibility that our proposed framework could be deliberately exploited to train and deploy socialbots to spread low-credibility content on social networks without being detected. To reduce any potential misuse of our work, we have also refrained from evaluating our framework with an actual socialbot detector API such as Botometer³. However, ultimately, such misuse can occur (as much as the misuse of the latest AI techniques such as GAN or GPT is unavoidable). Yet, we firmly believe that the benefits of our

framework in demonstrating the possibility of adversarial nature of socialbots, and enabling researchers to understand and develop better socialbot detection models far outweigh the possibility of misuse for developing "smarter" socialbots. In fact, by learning and simulating various adversarial behaviors of socialbots, we can now analyze the weakness of the current detectors. Moreover, we can also incorporate these adversarial behaviors to advance the development of novel bot detection models in a *proactive* manner [4]. Time-wise, this gives us a great advantage over the traditional *reactive* flow of developing socialbot detectors where researchers and network administrators are always one step behind the malicious bots developers [4].

One limitation of our current approach is that we only considered statistical features of a bot detector that are relevant to four activities–i.e., tweet, retweet, reply, and mention (Table 1). While these features help achieve 90% of detection accuracy in F1 score on a real-life dataset, we hope to lay the foundation for further works to consider more complex network and content-based features [11, 36, 37, 54].

7 CONCLUSION AND FUTURE WORK

This paper proposes a novel adversarial socialbot learning (ASL) problem where a socialbot needs to simultaneously maximize its influence on social networks and minimize the detectability of a strong black-box bot detector. We carefully designed and formulated this task as a cooperative game between two functional hierarchical reinforcement learning agents with a global reward. We demonstrated that the learned socialbots can sustain their presence on unseen real-life networks over a long period while outperforming other baselines in terms of network influence. During inference, the complexity of our approach also scales linearly with the number of followers and is independent of a network's structures and the virality of the news. Our research is also the first step towards developing more complex adversarial socialbot learning settings where multiple socialbots can work together to obtain a common goal [4]. By simulating the learning of these socialbots under various realistic assumptions, we also hope to analyze their adversarial behaviors to develop effective detection models against more advanced socialbots in the future.⁴

³https://botometer.osome.iu.edu/

 $^{^4\}mathrm{The}$ work was in part supported by NSF awards #1820609, #1940076, and #1909702

Socialbots on Fire: Modeling Adversarial Behaviors of Socialbots via Multi-Agent Hierarchical Reinforcement Learning WWW '22, April 25-29, 2022, Virtual Event, Lyon, France

REFERENCES

- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. OpenAI Gym. arXiv:arXiv:1606.01540
- [2] Chiyu Cai, Linjing Li, and Daniel Zengi. 2017. Behavior enhanced deep bot detection in social media. In 2017 IEEE International Conference on Intelligence and Security Informatics (ISI). IEEE, 128–130.
- [3] Wei Chen, Yajun Wang, and Siyu Yang. 2009. Efficient influence maximization in social networks. In Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining. 199–208.
- [4] Stefano Cresci. 2020. A decade of social bot detection. Commun. ACM 63, 10 (2020), 72–83.
- [5] Stefano Cresci, Roberto Di Pietro, Marinella Petrocchi, Angelo Spognardi, and Maurizio Tesconi. 2017. 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.
- [6] Stefano Cresci, Roberto Di Pietro, Marinella Petrocchi, Angelo Spognardi, and Maurizio Tesconi. 2017. Social fingerprinting: detection of spambot groups through DNA-inspired behavioral modeling. *IEEE Transactions on Dependable* and Secure Computing 15, 4 (2017), 561–576.
- [7] Stefano Cresci, Marinella Petrocchi, Angelo Spognardi, and Stefano Tognazzi. 2019. Better safe than sorry: An adversarial approach to improve social bot detection. In Proceedings of the 10th ACM Conference on Web Science. 47–56.
- [8] Stefano Cresci, Marinella Petrocchi, Angelo Spognardi, and Stefano Tognazzi. 2021. The coming age of adversarial social bot detection. *First Monday* (2021).
- [9] Clayton Allen Davis, Onur Varol, Emilio Ferrara, Alessandro Flammini, and Filippo Menczer. 2016. Botornot: A system to evaluate social bots. In Proceedings of the 25th international conference companion on world wide web. 273–274.
- [10] Guozhu Dong and Huan Liu. 2018. Feature engineering for machine learning and data analytics. CRC Press.
- [11] Phillip George Efthimion, Scott Payne, and Nicholas Proferes. 2018. Supervised machine learning bot detection techniques to identify social twitter bots. SMU Data Science Review 1, 2 (2018), 5.
- [12] Carlos Florensa, Yan Duan, and Pieter Abbeel. 2017. Stochastic Neural Networks for Hierarchical Reinforcement Learning. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017.
- [13] Sainyam Galhotra, Akhil Arora, Srinivas Virinchi, and Shourya Roy. 2015. Asim: A scalable algorithm for influence maximization under the independent cascade model. In *Proceedings of the 24th International Conference on World Wide Web*. 35–36.
- [14] Jason Gauci, Edoardo Conti, Yitao Liang, Kittipat Virochsiri, Yuchen He, Zachary Kaden, Vivek Narayanan, Xiaohui Ye, Zhengxing Chen, and Scott Fujimoto. 2018. Horizon: Facebook's open source applied reinforcement learning platform. arXiv preprint arXiv:1811.00260 (2018).
- [15] Tarun Gogineni, Ziping Xu, Exequiel Punzalan, Runxuan Jiang, Joshua Kammeraad, Ambuj Tewari, and Paul Zimmerman. 2020. TorsionNet: A Reinforcement Learning Approach to Sequential Conformer Search. In Advances in Neural Information Processing Systems, Vol. 33.
- [16] Jacob Goldenberg, Barak Libai, and Eitan Muller. 2001. Talk of the network: A complex systems look at the underlying process of word-of-mouth. *Marketing letters* 12, 3 (2001), 211–223.
- [17] Jacob Goldenberg, Barak Libai, and Eitan Muller. 2001. Using complex systems analysis to advance marketing theory development: Modeling heterogeneity effects on new product growth through stochastic cellular automata. Academy of Marketing Science Review 9, 3 (2001), 1–18.
- [18] Amit Goyal, Wei Lu, and Laks VS Lakshmanan. 2011. Celf++ optimizing the greedy algorithm for influence maximization in social networks. In Proceedings of the 20th international conference companion on World wide web. 47–48.
- [19] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining. 855–864.
- [20] Matthew Hindman and Vlad Barash. 2018. Disinformation, and influence campaigns on twitter. *Knight Foundation: George Washington University* (2018).
- [21] Ronald A Howard. 1960. Dynamic programming and markov processes. (1960).
- [22] Siwar Jendoubi, Arnaud Martin, Ludovic Liétard, Hend Ben Hadji, and Boutheina Ben Yaghlane. 2017. Two evidential data based models for influence maximization in twitter. *Knowledge-Based Systems* 121 (2017), 58–70.
- [23] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A Convolutional Neural Network for Modelling Sentences. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Baltimore, Maryland, 655–665. https: //doi.org/10.3115/v1/P14-1062
- [24] Harshavardhan Kamarthi, Priyesh Vijayan, Bryan Wilder, Balaraman Ravindran, and Milind Tambe. 2020. Influence Maximization in Unknown Social Networks: Learning Policies for Effective Graph Sampling. In Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (Auckland, New Zealand) (AAMAS '20). International Foundation for Autonomous Agents

and Multiagent Systems, Richland, SC, 575-583.

- [25] David Kempe, Jon Kleinberg, and Éva Tardos. 2003. Maximizing the spread of influence through a social network. In Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining. 137–146.
- [26] Masahiro Kimura and Kazumi Saito. 2006. Tractable models for information diffusion in social networks. In European conference on principles of data mining and knowledge discovery. Springer, 259–271.
- [27] Hautahi Kingi, Li-An Daniel Wang, Tom Shafer, Minh Huynh, Mike Trinh, Aaron Heuser, George Rochester, and Antonio Paredes. 2020. A numerical evaluation of the accuracy of influence maximization algorithms. *Social Network Analysis* and Mining 10, 1 (2020), 1–10.
- [28] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In International Conference on Learning Representations (ICLR).
- [29] George Dimitri Konidaris and Andrew G Barto. 2007. Building Portable Options: Skill Transfer in Reinforcement Learning.. In IJCAI, Vol. 7. 895–900.
- [30] Thai Le, Suhang Wang, and Dongwon Lee. 2020. MALCOM: Generating Malicious Comments to Attack Neural Fake News Detection Models. In 2020 IEEE International Conference on Data Mining (ICDM). 282–291. https://doi.org/10. 1109/ICDM50108.2020.00037
- [31] Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne Van-Briesen, and Natalie Glance. 2007. Cost-effective outbreak detection in networks. In Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining. 420–429.
- [32] Dexun Li, Meghna Lowalekar, and Pradeep Varakantham. 2021. CLAIM: Curriculum Learning Policy for Influence Maximization in Unknown Social Networks. arXiv preprint arXiv:2107.03603 (2021).
- [33] Hui Li, Mengting Xu, Sourav S Bhowmick, Changsheng Sun, Zhongyuan Jiang, and Jiangtao Cui. 2019. Disco: Influence maximization meets network embedding and deep learning. arXiv preprint arXiv:1906.07378 (2019).
- [34] Mei Li, Xiang Wang, Kai Gao, and Shanshan Zhang. 2017. A survey on information diffusion in online social networks: Models and methods. *Information* 8, 4 (2017), 118.
- [35] Siyuan Li, Rui Wang, Minxue Tang, and Chongjie Zhang. [n. d.]. Hierarchical Reinforcement Learning with Advantage-Based Auxiliary Rewards. In Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada.
- [36] Michele Mazza, Stefano Cresci, Marco Avvenuti, Walter Quattrociocchi, and Maurizio Tesconi. 2019. Rtbust: Exploiting temporal patterns for botnet detection on twitter. In Proceedings of the 10th ACM Conference on Web Science. 183–192.
- [37] Guanyi Mou and Kyumin Lee. 2020. Malicious Bot Detection in Online Social Networks: Arming Handcrafted Features with Deep Learning. In Social Informatics. 220–236.
- [38] OpenAI, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. 2019. Dota 2 with Large Scale Deep Reinforcement Learning. (2019). arXiv:1912.06680 https://arxiv.org/abs/1912.06680
- [39] Jorge Rodríguez-Ruiz, Javier Israel Mata-Sánchez, Raul Monroy, Octavio Loyola-Gonzalez, and Armando López-Cuevas. 2020. A one-class classification approach for bot detection on twitter. *Computers & Security* 91 (2020), 101715.
- [40] Ryan A Rossi, Rong Zhou, and Nesreen K Ahmed. 2018. Deep inductive graph representation learning. *IEEE Transactions on Knowledge and Data Engineering* 32, 3 (2018), 438–452.
- [41] Eldar Sadikov, Montserrat Medina, Jure Leskovec, and Hector Garcia-Molina. 2011. Correcting for missing data in information cascades. In Proceedings of the fourth ACM international conference on Web search and data mining. 55–64.
- [42] Mohsen Sayyadiharikandeh, Onur Varol, Kai-Cheng Yang, Alessandro Flammini, and Filippo Menczer. 2020. Detection of novel social bots by ensembles of specialized classifiers. In Proceedings of the 29th ACM International Conference on Information & Knowledge Management. 2725–2732.
- [43] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017).
- [44] Chengcheng Shao, Giovanni Luca Ciampaglia, Onur Varol, Kai-Cheng Yang, Alessandro Flammini, and Filippo Menczer. 2018. The spread of low-credibility content by social bots. *Nature communications* 9, 1 (2018), 1–9.
- [45] Yunlong Song, Mats Steinweg, Elia Kaufmann, and Davide Scaramuzza. 2021. Autonomous Drone Racing with Deep Reinforcement Learning. arXiv preprint arXiv:2103.08624 (2021).
- [46] Venkatramanan S Subrahmanian, Amos Azaria, Skylar Durst, Vadim Kagan, Aram Galstyan, Kristina Lerman, Linhong Zhu, Emilio Ferrara, Alessandro Flammini, and Filippo Menczer. 2016. The DARPA Twitter bot challenge. *Computer* 49, 6 (2016), 38–46.

WWW '22, April 25-29, 2022, Virtual Event, Lyon, France

- [47] Vladimir Svetnik, Andy Liaw, Christopher Tong, J Christopher Culberson, Robert P Sheridan, and Bradley P Feuston. 2003. Random forest: a classification and regression tool for compound classification and QSAR modeling. *Journal of chemical information and computer sciences* 43, 6 (2003), 1947–1958.
- [48] Youze Tang, Xiaokui Xiao, and Yanchen Shi. 2014. Influence maximization: Nearoptimal time complexity meets practical efficiency. In Proceedings of the 2014 ACM SIGMOD international conference on Management of data. 75–86.
- [49] Shan Tian, Songsong Mo, Liwei Wang, and Zhiyong Peng. 2020. Deep reinforcement learning-based approach to tackle topic-aware influence maximization. *Data Science and Engineering* 5, 1 (2020), 1–11.
- [50] Binghui Wang, Neil Zhenqiang Gong, and Hao Fu. 2017. GANG: Detecting fraudulent users in online social networks via guilt-by-association on directed graphs. In 2017 IEEE International Conference on Data Mining (ICDM). IEEE, 465–474.
- [51] Zheng Wen, Branislav Kveton, Michal Valko, and Sharan Vaswani. 2017. Online Influence Maximization under Independent Cascade Model with Semi-Bandit

Feedback. In Advances in Neural Information Processing Systems, Vol. 30.

- [52] Yuhao Wu, Yuzhou Fang, Shuaikang Shang, Jing Jin, Lai Wei, and Haizhou Wang. 2021. A novel framework for detecting social bots with deep neural networks and active learning. *Knowledge-Based Systems* 211 (2021), 106525.
- [53] Kai-Cheng Yang, Onur Varol, Clayton A Davis, Emilio Ferrara, Alessandro Flammini, and Filippo Menczer. 2019. Arming the public with artificial intelligence to counter social bots. *Human Behavior and Emerging Technologies* 1, 1 (2019), 48–61.
- [54] Kai-Cheng Yang, Onur Varol, Pik-Mai Hui, and Filippo Menczer. 2020. Scalable and generalizable social bot detection through data selection. In *Proceedings of* the AAAI Conference on Artificial Intelligence, Vol. 34. 1096–1103.
- [55] Cong Zhang, Wen Song, Zhiguang Cao, Jie Zhang, Puay Siew Tan, and Xu Chi. 2020. Learning to Dispatch for Job Shop Scheduling via Deep Reinforcement Learning. In Advances in Neural Information Processing Systems, Vol. 33. 1621– 1632.