# AI-Powered Social Bots

Terrence Adams
*tea331@g.harvard.edu*

*Abstract*—This paper gives an overview of impersonation bots that generate output in one or possibly, multiple modalities. We also discuss rapidly advancing areas of machine learning and artificial intelligence that could lead to frighteningly powerful new multi-modal social bots. Our main conclusion is that most commonly known bots are one dimensional (i.e., chatterbot), and far from deceiving serious interrogators. However, using recent advances in machine learning, it is possible to unleash incredibly powerful, human-like armies of social bots, in potentially well coordinated campaigns of deception and influence.

*Keywords*-bot, social bot, twitter bot, chatterbot, botnet, machine learning, artificial intelligence, Turing test

## I. INTRODUCTION

The term bot is a shortening of robot and refers to a software agent, or software robot device. It is gaining more attention very recently due to the prominence and influence of social media networks, as well as the ability to influence large groups through coordinated information operations. The potential of deploying large influence campaigns has already been realized, most notably during the 2016 U.S. election or the preceding Bexit vote. However, the influence of political (twitter) bots has been observed for years in several elections globally, including Turkey, Mexico and India. During the 2014 elections in India, Narendra Modi reached nearly 4 million twitter followers. However, when *Time* started monitoring Twitter for its Person of the Year award, the local media soon spotted a pattern. Thousands of Modi's followers were tweeting the same phrase: "I think Narendra Modi should be #TIMEPOY" at regular intervals, 24 hours a day - while a rival army of bots tweeted the opposite [4].

While social bots and specifically Twitter bots have demonstrated influence in political elections, and at the same time, in other areas (i.e., marketing), these bots are not known to show the sophistication that is quickly becoming a possibility. As social media has grown, pushing out the impact of traditional print journalism, there has been tremendous growth and effort on new technological advances related to machine learning. This change has been fueled by the availability of commodity-based distributed computing, a trend toward open source software development, and collection/curation of large-scale datasets. This has enabled major advances in machine learning where just about any competition in human language technologies (speech, language, visual) is dominated by "black box" deep neural network architectures.

*Why now?*

Deep neural networks are not new to machine learning. Frank Rosenblatt's original publication on the perceptron appeared in 1958 [5]. The first work detailing a learning algorithm for supervised deep feedforward multilayer perceptrons was developed by Ivakhnenko and Lapa in 1965 [6]. Other deep networks date back to 1971 [8], 1979 [9] and 1989 [11]. An effective backpropagation apparently first appeared in 1970 [7]. Recurrent neural networks have also been studied such as the Boltzmann machine, and in 1997, the Long Short-Term Memory (LSTM) network of Hochreiter and Schmidhuber [36]. LSTM is the basis for many current state of the art machine learning algorithms. What changed the equation is not that deep neural networks didn't exist or weren't studied. Instead, computing (e.g., GPGPU) and the creation of large-scale labeled datasets could be used to effectively train deep neural networks and learn a rich non-linear mapping from raw data (e.g., 224x224 images) to a low-dimensional representation and output layer (e.g., cat/not cat).

## II. A BRIEF HISTORY OF BOTS

According to *Botnets: The Killer Web Applications* [25], "bots were originally developed as a virtual individual that could sit on an IRC channel and do things for its owner while the owner was busy elsewhere". The first IRC bot called GM was developed in 1989, as a benevolent bot that would play a game of Hunt the Wumpus with IRC users.

It was a decade later that the first botnets became known. A botnet is a potentially large collection of bots that communicate with each other or communicate with a single botmaster. Pretty Park was first seen in March 1999, and introduced many of the features that were used for more than a decade after being introduced. It had the ability to report computer specifications, search email addresses, retrieve passwords, update its functionality, transfer files, redirect traffic, perform DoS attacks and communicate with an IRC server. SubSeven (May 1999) was the first remote controlled malware. The SubSeven trojan created a backdoor on the victim machine (zombie) by running the SubSeven server. IRC remote control started in a later version, when the SubSeven server was able to receive commands via IRC. Many more botnets have been deployed over the

years, including GTBot, SDBot, SpyBot, AgoBot, Rbot and Polybot.

Botnets have continued to evolve to avoid detection, now using alternate communication channels, http protocol (web pages or email sites that serve up commands), or peer-to-peer (P2P) protocols. They can modify DNS records and incorporate increasingly sophisticated techniques for command and control (C&C). Botnets exhibit new forms of malicious activity, including ransomware, instant message spam, blog spam and distribution of disinformation and fake news.

## III. TYPES OF BOTS

The website botnerds.com [28] lists several types of bots, broken down into good and bad bots:

- Chatbots
- Crawlers
- Transactional bots
- Informational bots
- Etertainment bots
- Hackers
- Spammers
- Scrapers
- Impersonators

Chatbots or chatterbots are bots that can communicate, generally through text messages, with humans. It is the type of bot typically associated with two of the best known Turing tests (Loebner Prize, University of Reading competition). Possibly, the earliest chatbot was ELIZA which was created by Joseph Weizenbaum of the MIT Artificial Intelligence Lab from 1964 - 1966. Other well known chatbots include Cleverbot and Tay.

Intelligent personal assistants such as Siri and Alexa act in a manner similar to a chatbot, but are typically much more sophisticated, communicate through audio/speech, and are presented as a service to users. These digital personal assistants may be considered the closest broadly available artificial intelligence (AI) technology for providing human-like or human-friendly information delivery. There continues to be a strong effort to teach these assistants personal skills (e.g., whispers, pauses, emotion [13]). IBM Watson was originally developed as a distributed question-and-answering system, but more recently is being tailored to act as a chatterbot for children's toys. Other chatterbots such as Eugene Goostman saw "15 minutes of fame" for winning the 2014 University of Reading competition. However, this bot, which acted the part of a 13-year old Ukranian boy, was quickly dismissed when exposed to the public through an Amazon Web Service. Figure 1 shows conversations conducted by Goostman. Note that Figure 1 was created by G. Fariello for his 2017 course on comparisons between natural intelligence, artificial intelligence, and the potential wide-ranging implications of rapidly advancing AI [1].

User friendly applications and APIs are being developed for integrated delivery of chatbot services [29], [30]. Also, on the bright side, a non-profit research company, OpenAI,



Figure 1. [1] Eugene Goostman bot posing as 13-year old Ukranian boy

was established for "discovering and enacting" a path to safe artificial intelligence [32]. Other self-proclaimed developers of friendly bots come from the Botwiki community [33], [34]. Also, Microsoft offers an open source bot builder SDK [31].

In this paper, we are most concerned with social bots that interact with social media platforms, and have the ability to impersonate users, while deceiving real users. As social bots adopt greater levels of AI, their behaviors become increasingly difficult to separate from actions of real human users.

## IV. INFORMATION OPERATIONS

There is a growing concern over the ability to unleash strategically mapped out and well coordinated influence campaigns using social bots. In particular, Facebook appears to be responding to this threat (in panic mode). On 27 April 2017, Facebook released version 1.0 of *Information Operations and Facebook* [14]. This guide starts by quoting founder Mark Zuckerberg (Feb 2017):

> It is our responsibility to amplify the good effects and mitigate the bad – to continue increasing diversity while strengthening our common understanding so our community can create the greatest positive impact on the world.

The document continues to define several terms bandied in the media such as: false news, false amplifiers, disinformation, misinformation, intent, medium, amplification. The article lists three major features of online information operations that Facebook assesses have been attempted on Facebook.

- Targeted Data Collection
- Content Creation
- False Amplification

The article primarily focuses on the first and third items. It is the second item which has the potential to advance rapidly based on new machine learning breakthroughs, and we focus on this topic in later sections.

While information operations are gaining an increased concern globally, it is believed by many that this is a global threat that is taken far too lightly. Also, on 27 April 2017, the US Senate Committee on Armed Services, subcommittee on cybersecurity, held a hearing on cyber-enabled information operations. Presenting at this meeting were John C. Inglis (Former Deputy Director, NSA), Honorable Michael D. Lumpkin (Principal at Neptune Computer Inc and former Acting Under Secretary of Defense for Policy), Dr. Rand Waltzman (Senior Information Scientist, RAND Corp) and Mr. Clint Watts (Robert A. Fox Fellow, Foreign Policy Research Institute). The testimony is too long to present here, but contains alarming facts demonstrating the extent of information operations directed against the US and other nations. It calls for major changes to the US effort defending against information operations, and the influence campaigns conducted against our citizens.

### Twitter Bot Challenge

In February/March 2015, DARPA held a 4-week competition to identify a set of previously identified "influence bots" serving as ground truth on a specific topic within Twitter [2], [19]. The topic was pro-vaccination/anti-vaccination. The company Sentimetrix scored at the top when considering the metric,

$$FinalScore = Hits - 0.25 * Misses + Speed.$$

There were a total of 39 ground truth influence bots, and Sentimetrix had only 1 false hit. Participant USC correctly guessed all 39 bots with no false hits. There were some clear indicators of bots (e.g., the variation of sentiment from the bots tends to be lower than the variation from humans). Humans run more hot-and-cold. In a paper summarizing the results of the challenge, the following types of features were shown to be of interest [19]:

- Tweet Syntax
- Tweet Semantics
- Temporal Behavior
- User Profile
- Network

The Twitter Bot Challenge laid ground-work for future evaluations of bot detection in social media. We think it's important to consider much more sophisticated bots that are active across multiple media platforms.

## V. Deep Learning

The scary part of this story is that there is no evidence that today's influence campaigns are tapping into the rapidly advancing machine learning or AI technologies. The amount of research into generating signals that mimic human behavior, as well as, complex deep neural netowrk architectures that digest, understand and continually update based on human interaction, is staggering.

It was only in 2012 that Krizhevsky, Sutskever and Hinton published *ImageNet Classification with Deep Convolutional Neural Networks* [12] that presented a ground-breaking leap in image recognition using a neural network with 60 million parameters and 500,000 neurons. Since then, performance has continued to improve steadily, using increasingly complex networks, and increasingly larger datasets. As an example, MSRA won the ImageNet challenge in 2015 using a network with 154 layers. Google admitted to utilizing a database of 500 million facial images to train their face recognizer for Megaface. More recently, Google has announced a significant leap in machine translation using stacked recurrent neural networks (LSTMs) with both bi-drectional and uni-directional layers, as well as residual connections between layers [17].

### A. Reinforcement Learning

The deep learning improvements highlighted above are supervised techniques depending on huge amounts of labeled data (text, image, speech or video). However, there have been significant advances in unsupervised or semi-supervised machine learning. Techniques that utilize reinforcement learning accept goal oriented responses, and learn to optimize their decisions in complex environments. Probably the most recent success reported along these lines, is the ability of a Google DeepMind AlphaGo system to beat a 9-dan professional without handicaps. Initially, the system was developed using supervised learning from human play, but then it was optimized using reinforcement learning, by being set to play against other instances of itself. Duplicating this success does not appear very easy. Facebook's Darkforest system has not defeated a professional human Go player. Also, Dwango and the University of Tokyo developed DeepZenGo which has not yet demonstrated dominance over top human players.

### B. Generative Adversarial Networks

A new deep learning architecture known as generative adversarial networks (GANs) [20] has the potential to improve recognition in multiple domains, and at the same time deliver systems which can generate non-human output that resembles human-like characteristics. GANs are constructed from two multilayer perceptrons G and D. The network $G = G(z; \theta_g)$ will take input parameters and map to the data space. The network $D = D(x; \theta_d)$ takes parameters and input from the data space and outputs a single scalar. $D$ is trained to maximize the probability of assigning the correct label to both training examples and samples from $G$. Simultaneously, $G$ is trained to minimize $\log(1 - D(G(z)))$, or alternatively, maximize $\log D(G(z))$.

While the discriminating network $D$ is learned, the network $G$ is improved to produce data that resembles the initial labeled data fed into $D$, and its probability distribution.

3

Figure 2.   Synthetically generated bedrooms taken from [21]

See [21] and [22] for interesting output from some trained generative networks.

## VI. Multimodal Personas & Scary Bots

Deep neural networks are being applied in just about every human language technology to achieve state of the art for recognition. With companion generative networks, it will (or is already) possible to generate output that mimics human behavior, as well as human environments [21]. For instance, a Montreal based company Lyrebird [15] recently demonstrated the ability to take any user's voice and create a speechbot that talks in a similar manner to the given voice. Today, it is still possible for most human ears to detect anomalies with the bot voice that separate it from the real human voice, but it's uncannily close. With more collected data, and a larger neural network, it will only get better. Also, Google DeepMind produced a similar capability called WaveNet [16].

It is becoming commonplace for Hollywood movies to modify facial gestures of actors after shooting has completed, and without the participation of the actual actors. It can be done synthetically, and movie goers have no idea any changes were made. Also, there are efforts to synthetically produce digital actors/actresses which are amazingly close to real-life performers [35].

With the ability to generate text (i.e., chatterbot), speech or imagery, it is a matter of time before someone puts these capabilities together to build full walking, talking, texting social bots. While the replication of physical robots will be limited initially by the ability to manufacture material resources (e.g., metal), the proliferation of various social bots will only be limited by storage and compute power.

## References

[1] Fariello, G., Building the Brain: A Survey of Artificial Intelligence, *Harvard University Extension School*, **CSCI E-86** (2017).

[2] Waltzman, R., Personal communication, *14 April 2017*.

[3] Turing, A.M., Computing machinery and intelligence, *Mind*, **59** 433 - 460 (1950).

[4] Woollacott, E., Why fake Twitter accounts are a political problem, *New Statesman*, 28 May 2014.

[5] Rosenblatt, F., The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain, *Psychological Review*, **65:6**, 386-408 (1958).

[6] Ivakhnenko, A.G. and Lapa, V.G., Cybernetic Predicting Devices, *CCM Information Corporation* (1965).

[7] Linnainmaa, S., The representation of the cumulative rounding error of an algorithm as a Taylor expansion of the local rounding errors, *Master's Thesis, Univ. Helsinki* (1970).

[8] Ivakhnenko, A.G., Polynomial Theory of Complex Systems, *IEEE Transactions on Systems, Man and Cybernetics*, **4**, 364 - 378 (1971).

[9] Fukushima, K., Neural network model for a mechanism of pattern recognition unaffected by shift in position - Neocognitron, *Trans. IECE*, **J62-A:10**, 658 - 665 (1979).

[10] Fukushima, K., Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position, *Biological Cybernetics*, **36:4**, 193 - 202 (1980).

[11] LeCun, Y. et. al., Jackel: Backpropagation Applied to Handwritten Zip Code Recognition, *Neural Computation*, **1:4**, 541-551 (1989).

[12] Krizhevsky, A., Sutskever, I. and Hinton, G., ImageNet Classification with Deep Convolutional Neural Networks, *Advances in Neural Information Processing Systems*, **25** (2012).

[13] Perez, S., Alexa learns to talk like a human with whispers, pauses & emotion, *Tech Crunch*, 28 April 2017.

[14] Weedon, J., Nuland, W. and Stamos, A., Information Operations and Facebook, https://fbnewsroomus.files.wordpress.com/2017/04/facebook-and-information-operations-v1.pdf, 27 April 2017.

[15] Lomas, N., Lyrebird is a voice mimic for the fake news era, *Tech Crunch* 25 April 2017.

[16] van den Oord, A., WaveNet: A Generative Model for Raw Audio, https://arxiv.org/pdf/1609.03499.pdf (2016).

[17] Wu, Y. et. al., Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation, https://arxiv.org/abs/1609.08144 8 Oct 2016.

[18] Inglis, J.C., Lumpkin, M.D., Waltzman, R. and Watts, C., Cyber-enabled Information Operations, *Senate Committee on Armed Services, subcommittee on cybersecurity* 27 April 2017.

[19] Subrahmanian, V.S. et. al., The Twitter Bot Challenge, https://arxiv.org/abs/1601.05140 (2016).

[20] Goodfellow, I.J. et. al., Generative Adversarial Networks, https://arxiv.org/abs/1406.2661 10 June 2014.

[21] RADFORD, A., METZ, L. and CHINTALA, S., Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks, https://arxiv.org/abs/1511.06434 (2016).

[22] HO, J. and ERMON, S., Generative Adversarial Imitation Learning, *Advances in Neural Information Processing Systems* (2016).

[23] LEE, C.P., Framework for Botnet Emulation and Analysis, *Doctoral Dissertation, Georgia Institute of Technology* (2009).

[24] BOSHMAF, Y., Security Analysis of Malicious Social Bots, *Doctoral Dissertation, The University of British Columbia*

[25] SCHILLER, C. and BINKLEY, J.R., Botnets: The Killer Web Applications, *Syngress* (2011).

[26] FERRARA, E. et. al., The Rise of Social Bots, *Communications of the ACM*, **59:7** 96 - 104 (2016).

[27] DAVIS, C.A., et. al., BotOrNot: A System to Evaluate Social Bots, https://arxiv.org/abs/1602.00975 2 Feb 2016.

[28] http://botnerds.com/types-of-bots/

[29] GELFENBEYN, I. et. al., https://api.ai/, *Google, formerly Speaktoit* (2017).

[30] SLACK API, Slackbot, *Slack Technologies* https://api.slack.com/bot-users

[31] BOT FRAMEWORK, Open Source Bot Builder SDKs, *Microsoft* https://dev.botframework.com

[32] NON-PROFIT AI RESEARCH COMPANY, OpenAI, *Sponsors Sam Altman et. al.* https://openai.com/

[33] HTTPS://BOTWIKI.ORG/ABOUT/TEAM/, botwiki.org, *open catalog of friendly, useful, artistic online bots* https://openai.com/

[34] HOSTED BY SLACK, Botmakers, *https://botwiki.org/about/team/* https://botmakers.org

[35] ALEXANDER, O. et. al., The Digital Emily Project: Achieving a Photoreal Digital Actor, *USC Institute for Creative Technologies*, http://gl.ict.usc.edu/Research/digitalemily/ (2017).

[36] HOCHREITER, S. and SCHMIDHUBER, J., Long Short-Term Memory, *Neural Computation*, **9:8** 1735 - 1780 (1997).

[37] DONAHUE, JEFF, HENDRICKS, LISA ANNE, GUADARRAMA, SERGIO and ROHRBACH, MARCUS, Long term recurrent convolutional networks for visual recognition and description, http://arxiv.org/pdf/1411.4389v3.pdf, **v3** (2015).

[38] GOODFELLOW, I., BENGIO, Y. and COURVILLE, A., Deep Learning, http://www.deeplearningbook.org/, *MIT Press* (2016).