

# VIRUS

29

## THE DIGITAL AND THE SOUTH: QUESTIONINGS VOL. 2

PORUGUÊS-ESPAÑOL | ENGLISH  
REVISTA . JOURNAL  
ISSN 2175-974X  
CC-BY-NC-SA

UNIVERSITY OF SAO PAULO  
INSTITUTE OF ARCHITECTURE AND URBANISM  
NOMADS.USP  
REVISTAS.USP.BR/VIRUS  
DECEMBER 2024



# VI20

THE DIGITAL AND THE SOUTH: QUESTIONINGS VOL. 2

O DIGITAL E O SUL: TENSIONAMENTOS VOL. 2

LO DIGITAL Y EL SUR: CUESTIONAMIENTOS VOL. 2

## EDITORIAL

001 THE DIGITAL AND THE SOUTH: QUESTIONINGS VOL. 2

O DIGITAL E O SUL: TENSIONAMENTOS VOL. 2

LO DIGITAL Y EL SUR: CUESTIONAMIENTOS VOL. 2

MARCELO TRAMONTANO, JULIANO PITA, PEDRO TEIXEIRA, CAIO NUNES, ISABELLA CAVALCANTI, RENAN TEIXEIRA, ALINE LOPES

## INTERVIEW

004 THE TECHNOCENE AND THE REESTABLISHMENT OF A HORIZON OF URGENCY

O TECNOCENO E O RESTABELECIMENTO DE UM HORIZONTE DE URGÊNCIA

EL TECNOCENO Y EL RESTABLECIMIENTO DE UN HORIZONTE DE URGENCIA

HENRIQUE PARRA, PEDRO TEIXEIRA, MARIO VALLEJO

## AGORA

015 DYSPHORIA AS THE POTENCY OF CONTRADICTIONS: A BET BY PAUL B. PRECIADO

DA DISFORIA COMO POTÊNCIA DAS CONTRADIÇÕES: UMA APOSTA DE PAUL B. PRECIADO

MARCOS BECCARI

023 DIGITAL FRAMEWORKS / MODERN URBAN FRAMES

ESTRUTURAS DIGITAIS / ESTRUTURAS URBANAS MODERNAS

CARLOS FEFERMAN

033 GLOBAL SOUTH ADRIFT: DIGITAL REGULATION IN THE EUROPEAN UNION AND BRAZIL

SUL GLOBAL À DERIVA: REGULAÇÃO DIGITAL NA UNIÃO EUROPEIA E NO BRASIL

MAGNO MEDEIROS

042 DIGITAL ACTIVISM AND PLATFORM (DE)REGULATION IN ELECTORAL CONTEXT

ATIVISMO DIGITAL E (DES)REGULAÇÃO DE PLATAFORMAS NO CONTEXTO ELEITORAL

ARNALDO DE SANTANA SILVA, MILENA CRAMAR LÔNDERO, VITÓRIA SANTOS

- 052 COSMOPLATFORMIZATION: DIGITAL PLATFORMS FROM THE GLOBAL SOUTH  
COSMOPLATAFORMIZAÇÃO: PLATAFORMAS DIGITAIS A PARTIR DO SUL GLOBAL  
ELI BORGES JUNIOR, EVANDRO LAIA, BRUNO MADUREIRA
- 060 SOCIAL ROBOTS: A SOCIO-TECHNICAL CONTROVERSY  
*BOTS SOCIAIS: UMA CONTROVÉRSIA SOCIOTÉCNICA*  
RAMON FERNANDES LOURENÇO
- 069 LAND, FREEDOM, AND DIVERSITY: METAPHORS TO THE DIGITAL WORLD?  
TERRA, LIBERDADE E DIVERSIDADE: METÁFORAS PARA O MUNDO DIGITAL?  
LUCCA AMARAL TORI
- 079 BETWEEN PHYSICAL AND VIRTUAL WINDOWS: OPENINGS OF LIVING IN THE PANDEMIC  
ENTRE JANELAS FÍSICAS E VIRTUAIS: ABERTURAS DO MORAR NA PANDEMIA  
PAULA LEMOS VILAÇA FARIA

## PROJECT

- 087 ECOLOGICAL ENSEMBLE  
CONJUNTO ECOLÓGICO  
ANA CECILIA PARRODI ANAYA

## **SOCIAL ROBOTS: A SOCIO-TECHNICAL CONTROVERSY** **BOTS SOCIAIS: UMA CONTROVÉRSIA SOCIOTÉCNICA**

RAMON FERNANDES LOURENÇO

Ramon Fernandes Lourenço holds a Bachelor's degree in Social Communication – Public Relations and a Master's degree in Information Science. He is a PhD candidate in the Postgraduate Program in Contemporary Integration of Latin America at the Federal University for Latin American Integration, Brazil. His research topics are related to Digital Communication, Political Communication, International Relations, Technological Mediation, Social Networks, Actor-Network Theory, Sociotechnical Networks, and Research Methods in the Digital Environment. uel.ramon@gmail.com

<http://lattes.cnpq.br/8171408485283759>

60

ARTICLE SUBMITTED ON AUGUST 4, 2024

Lourenço, R. F. (2024). Social Robots: A Socio-technical Controversy. *VIRUS*, (29). The Digital and the South: Questionings Vol. 2. 60-68  
<https://doi.org/10.11606/2175-974x.virus.v29.229585>

## Abstract

The rise of far-right groups in Latin American and Caribbean countries reveals the use of strategies to manipulate public discussions, including the massive use of fake profiles on social media. Therefore, using case study methodology, this article aims to analyze the concept of social bots by describing the main initiatives for detecting automated profiles on X, formerly Twitter. By using the Actor-Network Theory, it was possible to uncover the complexity involved in defining a computerized profile, pointing to the need to establish an umbrella concept that encompasses practices such as automated profiles, robots, hybrid profiles, sockpuppets, and meatpuppets. Ultimately, we identify that the Botometer X, Pegabot, Bot Sentinel, and Bot Slayer platforms are based on automated monitoring methodologies, with data on user behavior patterns as the essential elements that indicate whether users are humans or robots.

**Keywords:** Social bot, Social media, Sockpuppet, Meatpuppet, Actor-Network Theory

## 1 Introduction

Following the main political movements in Latin America in recent years, a growing tension has been directly linked to using strategies to manipulate digital conversational dynamics. The advance of far-right groups, which gained strength with the election of Jair Bolsonaro in Brazil in 2018 and more recently with the election of Javier Milei in Argentina in 2022, highlights the need to expand studies on how manipulating public discussions on social media has been damaging to democracies, especially in the Global South. Questioning public institutions is among the traditional strategies of far-right groups, such as allegations of electoral fraud (Yáñez, 2022), the production and massive sharing of fake news (Esquivel, 2022), the use of robots and other mechanisms to manipulate public discussions (Azevedo Júnior & Lourenço, 2023), resulting in increased polarization and threaten electoral processes.

The importance of analyzing these controversies grows as the speed of Internet connections evolves, increasing the number of agents connected to the World Wide Web and becoming one of the most striking features of digitizing politics. In addition to identifying and understanding this content as a whole, it is also necessary to map the circulation flows, identifying the agents and their roles in these networks. That is why we need to focus on analyzing the agents who participate in discussions on the Internet, trying to understand who they are, how they organize themselves, and how they influence the production and sharing of information. One specific profile of an agent deserves much attention: social bots. This term can be translated as social robots and describes automated accounts for sharing and interacting with social media content.

Automated accounts feature prominently in discussions about digital controversies, mainly due to a negative view of their behavior in manipulating public discussion. Such discussions raise the need for greater control of these accounts on social media, seeking their identification, banning, and accountability for managing these social robots. Based on this debate, this article aims to analyze the concept of social bots by describing the main initiatives for detecting automated profiles on X, formerly Twitter. Its specific objectives are to describe the behaviors analyzed by these initiatives, to explain how human and non-human profiles are differentiated, to delimit the differences between the categories of robots and cyborgs, and, finally, to explain the most common behaviors of automated accounts.

The methodology used was a case study, which enabled a detailed analysis and description of the initiatives mapped (Eisenhardt, 1989; Yin, 2009). This mapping was carried out between May 2020 and October 2022, making up Twitter's move to X. The focus on social media X is relevant because it is more open to monitoring initiatives of this nature and is also widely used by important institutions and world leaders. Teixeira (2018) reinforces the importance of social media for building public opinion and anticipates the challenge of ensuring that these social robots are not used en masse as an instrument of colonization by anti-democratic groups. Indeed, it is precisely these groups that lead the use of manipulation strategies with automated profiles around the world, especially in Brazil (Ruediger et al., 2017) and the United States (Bessi & Ferrara, 2016), but with a growing presence in the Global South, in other Latin American and Caribbean countries (Tricontinental: Institute for Social Research, 2021). Before detailing the initiatives, a theoretical contribution is appropriate by pointing out the difficulties of defining a social robot and its role in social media. In identifying whether a human or an automated algorithm manages an account, any attempt at simplification based on a few specific characteristics of these accounts has proved unproductive. More is needed to

analyze the pattern of profile pictures or the names and addresses of accounts. It is necessary to follow in the footsteps of these actors. To advance this understanding, it is worth mentioning the principle of symmetry and technical mediation (Latour, 2012; Law, 1992; Callon, 2004), which are fundamental elements for this discussion.

Actor-network theory (ANT) has made significant progress in analyzing non-human elements' role in bringing about transformations, a process based on the principle of symmetry. This principle is anchored in the denial of a natural primacy of man over things, or, as Santaella and Cardoso point out, "(...) Latour rejects both a determinism of technique over the human (materialism) and the determinism of the human over technique (anthropocentrism)." (Santaella & Cardoso, 2015, p. 169). ANT starts from the symmetry between the agents in a network, with their nature not mattering in principle but rather the actions they undertake.

The same applies to the controversy over whether a social media account is managed by a human, as descriptive characteristics of the account are not decisive for analysis. Lemos (2013) reinforces this understanding by stating that "(...) entities have their attributes acquired as a result of the relationship with other entities and not by their inherent qualities" (Lemos, 2013, pp. 64-65). In this sense, an automated profile is the sum of the social media, the algorithm programmed to automate a specific action, and the actions of the programmer who created the social robot. Since actions play a leading role from the perspective of ANT, it is crucial to go into greater detail. Thus, technical mediation is built on four pillars: interference, composition, reversible obscuration, and delegation (Latour, 2012). Interference is the action program carried out by the agent in the network in which it is inserted. It is the action of interfering in ongoing flows, generating a transformation. For this article, it is the visualization of the actions of social robots that interfere in public discussions.

The second pillar understands that in a network, every action generates a series of other actions, or, in other words, a series of articulated actions. This pillar points to the characteristic that every action can be broken down into micro-actions, thus revealing other agents. In short, we verify that in the action of automated accounts, it is not possible to impute responsibility only to the social robot but also to the programmer who carried it out and the person who hired the programmer and financed the whole process of sharing false information. In this sum of responsibilities, the action's meaning is ascertained, and the network shows itself as a fluid and ever-changing space. Composition articulates actions resulting from a first movement: "It is the multiplication of sub-programs that results in composition." (Melo, 2011, p. 10).

Thus, it is possible to see how intricate the notion of technical mediation is and its relevance to this challenge of understanding and identifying automated profiles on social media. Reversible obscuration is the third pillar of technical mediation, revealed as the actions are composed.

Whenever a network acts as a single block, then it disappears, replaced by the action itself and the apparently unique author of this action. At the same time, the form in which the effect is produced is also erased: in the circumstances, it is neither visible nor relevant. It then happens that something much simpler emerges - a (working) television, a well-run bank, or a healthy body - for a while, to cover the networks that produced it. (Law, 1992, p. 385)

In normality, the action's complexity is then obscured and simplified into its principal agent or main effect. In the current context, crises have revealed the complexity of the actions of these agents on social media. It took the interference in critical electoral processes to intensify the analysis of the role and effects of these agents. When analyzing the role, we see the last pillar of technical mediation: delegation. Since social robots are lines of code programmed by programmers who respond to a need pointed out by a third party, the principle of delegation becomes easier to understand. It is the ability to delegate an action program to an actor in the network.

Mediation is built through this complex correlation between the pillars presented, and the solution for identifying automated social media accounts follows the path of analyzing mediation. This should make it clear that

[...] the idea of mediation is being related here to a sharing of responsibility for the action between various actors, respecting the action of all those involved in the technique in question. This is what the author means by composition since only the sum of all the agents involved can give meaning to mediation. (Santaella & Cardoso, 2015, p. 171)

After this brief contextualization, the next section will present some definitions of social robots and also describe how the main initiatives to combat these practices are moving towards, following the traces left by these actors.

### 3 The Definitions of Social Robots

Social bots are computer programs designed to mimic human behavior on social media, such as retweeting messages from a profile or constantly tweeting the same message (Davis et al., 2016; Ferrara et al., 2016; Ruediger et al., 2017). These programmed activities aim to influence public opinion in favor of specific groups. These automated agents can be seen in full action in contemporary electoral processes, especially with the rapid growth of polarization in Latin America and the Caribbean due to the advance of the extreme right. Ferrara and other authors highlight the central objective of these agents: "A robot is an algorithm that produces content and interacts with people on social networks, emulating and even altering their behavior." (Ferrara et al., 2016, p. 96). The attempt to imitate a human user and alter their behavior reveals the nefarious link between these tools and the main controversies in the political field, whose ultimate goal is to influence public opinion (Ruediger et al., 2017).

The diversity of behaviors performed by these agents reinforces the need to avoid the simplifying description of a single category and to explore the complexities of these practices. With a growing understanding of these agents, monitoring and combating this practice is gaining momentum. It is therefore necessary to understand the main types of technological agents surrounding the most significant public discussions on social media today, which can be classified into human and non-human. Accounts legitimately managed by humans have specific behavior patterns. They have a great diversity of content, interact with a network of users, and invest time in consuming information on other profiles. Some may even have a high volume of posts on the same day, differing from the standard average. A typical presence in Internet discussions who exhibits this type of behavior is the troll, "an individual who seeks to interfere with the progress of a discussion in a particular online community by posting nasty or out-of-context comments." (Zago, 2012, p. 151). Trolls can carry out these actions with their own or fake profiles, but a human still manages the profile.

Social robots, on the other hand, have a more limited set of behaviors compared to humans. However, with the growing development of this technology, they become increasingly complex when imitating a human profile, making the detection process increasingly complicated. Because of these advances, it is necessary to make an initial distinction when working with the concept of robots. In principle, there is a need to establish an umbrella category to describe any social media profile with a certain level of automation: social bots. Then, those we call just bots are operated entirely by a computer program. At the same level are hybrid profiles, operated part of the time by an algorithm and part of the time by humans, called cyborgs (Duarte et al., 2016). Cyborgs are the latest strategies to give profiles more credibility and circumvent robot detection mechanisms.

Among the practices that define cyborgs, those that have gained popularity are sockpuppets and meatpuppets. Liu and other authors (2016) define sockpuppets as multiple accounts controlled by the same individual, while meatpuppets are multiple accounts controlled by a group of people, usually from the same organization (Liu et al., 2016). It is interesting to note the coexistence of these strategies for a single purpose - to meet the objectives of the organization they work for - whether they are automaton profiles (robots), those managed by the same human operator (sockpuppets), or a group of humans operating fake social media profiles (meatpuppets). Solorio, T., Hasan, R. & Mizan, M. (2013) explore some of the standard actions of this type of organization, demonstrating the growing challenge of media detection initiatives:

(...) *smart sockpuppet* can therefore avoid detection by using multiple IP addresses, modifying writing style, and altering behavior. In addition, a malicious user can create dormant accounts that perform benign edits from time to time but are used as puppets when necessary. Identifying these accounts as puppets is not obvious, as these accounts can have a long and diverse editing history. (Solorio et al., 2013, p. 59)

Finally, the growing challenge of trying to locate and identify accounts managed by robots and cyborgs is clear. It is now time to present some of the initiatives and demonstrate the main information used as a parameter to identify these accounts.

### 4 The Initiatives and their Methodologies

Before starting the presentation of the platforms, it should be noted that the initiatives mapped for this case study were identified before Twitter's move to X, which implied severe changes to the data usage permissions of these tools.

### a. Botometer X

*Botometer X* (<https://botometer.iuni.iu.edu>) is available on a website, created in 2014 under the old name of *BotOrNot*. According to the descriptions on the page, the tool uses data from a user's activities on X and returns a score based on the likelihood of this account being a social robot. The higher the score, the more likely the account is to be characterized as a social robot. The classification system devised by the platform's authors is based on six main classes: a) network patterns, b) user characteristics, c) friends, d) temporality, e) content, and f) sentiment. Network characteristics analyze information dissemination patterns by analyzing networks of retweets, mentions, and co-occurrence of hashtags (Davis et al, 2016). These analyses are carried out using their statistical characteristics, which reveal distribution patterns and relationships between their elements. User characteristics identify metadata including language, geographical locations, and the date and time of account creation. Information on friends includes an analysis of followers, profiles followed by the account, and posts made, among other information (Davis et al., 2016).

In addition to the networks, information on the user and their friends, the characteristics of temporality, content, and sentiment of the content are used in *Botometer*. With regard to temporal characteristics, information is captured on posting patterns and information consumption on the platform. Regarding the content of posts, linguistic information is analyzed using natural language processing and sentiment analysis, seeking to identify the main emotions captured from each post (Davis et al., 2016, p. 274). *Botometer X* currently operates in archive mode, providing historical data collected until May 31, 2023.

### b. Pegabot

*Pegabot* (<https://pegabot.com.br/>) is a Brazilian initiative that was launched in 2018. According to its creators, its aim is to contribute to the fight against disinformation in Brazil, targeting journalists, experts, and civil society organizations. Its dynamic follows the logic of *Botometer*, assigning a score to an analyzed profile. According to the information available on the project's website, the tool analyzes the profile's posting history in search of patterns in three main categories: profile analysis, network analysis, and sentiment analysis. Profile analysis takes into account the name, number of profiles followed and followers, description text, number of posts, and favorites. This information is processed using metrics such as the character count of the name, evaluation of the age of the profile, number of tweets, and existence of a profile photo, among others.

Each of these elements has a direct influence on the profile's score. Interaction dynamics are carried out by collecting a sample of the user's timeline and identifying hashtags and mentions of the profile. To this end, it identifies the distribution of hashtags and mentions, seeking to understand whether the user is forwarding spam messages. It also performs sentiment analysis on a sample of the 100 most recent posts. With this, it seeks to identify whether a specific emotion is prevalent, whether negative or positive. The more neutral the profile, the lower the score and the more likely it is to be a social robot. The website currently has errors, precisely because of the change in the policy for accessing X's data.

### c. Bot Sentinel

*Bot Sentinel* (<https://botsentinel.com>) was created in 2018 and worked in two ways, on a website and also integrated with Twitter. However, in 2022 the platform entered into a dispute with Twitter over the allegation that it was violating the company's policies, losing its integration with social media. After that, *Bot Sentinel* continues to work only on the website with historical data and some recent information, due to the limitations imposed by X's new policy. On the website there is still an information-rich dashboard where you can follow the monitoring carried out automatically on X. Previously there was a function on Twitter called *Check user*, present on users' profiles. By clicking on this option, *Bot Sentinel* was triggered to check whether the profile was a robot, helping to combat disinformation.

According to the description on the *Bot Sentinel* website, it is based on a machine learning model that uses Twitter's own rules as a standard, unlike other platforms that create models based on researchers' interpretation of collected data. There is no detail on what types of information are used and how Twitter's rules were applied in the analysis. The classification system is also not detailed but is limited to

describing that accounts are classified in a system that goes from zero to one hundred percent, with the higher the percentage, the greater the possibility of the account being involved in harassment, trolling, or deceptive tactics.

#### d. Bot Slayer

From the same creators as *Botometer*, *Bot Slayer* (<https://osome.iuni.iu.edu/tools/botlayer/>) differs from previous platforms by focusing on the flow of sharing malicious information on Twitter. Throughout this process, the platform also analyzed the characteristics of the users involved in sharing certain content, indicating whether or not the profile might be an automaton. To arrive at this classification, Hui and other authors (2019) indicate that *Bot Slayer* extracts four characteristics from each profile: volume, trendiness, diversity, and botness. For volume information, the tool counts the number of posts involving the user over a given period of time. The Trendiness characteristic "is calculated as the ratio between the entity's volume in two consecutive time windows" (Hui et al., 2019: p. 3). Diversity is a value calculated from the ratio between the number of unique users and the number of posts, and botness is the measure of an account's classification as a possible social robot.

From these four platforms, it is possible to see the diversity of methodologies used to identify anomalous behavior on digital social media. These methodologies range from methods that use less information, such as Pegabot, to those that are more complete and fully described in scientific publications, such as *Botometer*. There are also major similarities in some of the cases, such as the proximity of the methodologies and techniques used by *Botometer X* and Pegabot. Now, in order to delve deeper into this analysis, it is necessary to look at the main methodologies used by these platforms to collect and analyze large volumes of data, as this is a major challenge in implementing tools to monitor the flows of disinformation circulating on social media.

### 5 The Challenges of the Hunt

The processes of interaction and information sharing on social media generate large volumes of information, which has a high impact on the costs of collecting, processing and analyzing this data. And the exponential growth of these flows is perhaps one of the biggest challenges facing initiatives that seek to analyze the anomalous behaviors that influence public discussions. So, in order to deal with this large volume of data, some techniques make use of great computational potential and are making these initiatives viable. Embracing the challenge of analyzing the main methodologies for detecting bots on social media, Alothal, E., Zaki, N., Mohamed, E. A., & Alashwal, H. (2018) observed that the platforms so far use the behavioral patterns of each account as a fundamental element. In other words, as can be seen in the description of the four platforms mentioned in this study, the various variables analyzed are directly related to the behavior of profiles, how they interact with other users, and how they share information, among other aspects. However, the lack of consensus on which characteristics best represent a social robot on social media is still a challenge.

To help in this process, Ferrara and other authors (2016) systematized in their study the main information used by detection initiatives. The authors highlight the relevance of elements such as the number of posts and reposts, replies, mentions, number of shares made by the analyzed account, its age at creation, and the length of the user's name. As a result, they define that a social robot has a high number of reposts, an account with a more recent creation date, a low number of posts, and a username with many characters (Ferrara et. al., 2016, p. 102). It is essentially an artificially created account, with random names, with the sole aim of replicating content as much as possible. Its life cycle is different from a human user account, as it is usually created for a specific task, acting systematically for a short period of time, before being identified and taken down by robot detection mechanisms. However, identifying the fixed characteristics of social robot accounts faces a number of challenges when taking into account the performance of cyborg profiles, precisely because of their high capacity for adaptation. This makes it difficult to implement mass detection strategies for these profiles.

Another challenge linked to this issue is the main methodologies applied in monitoring platforms. Some authors classify these platforms into three main groups: graph-based, crowdsourcing and machine learning (Alothal et al., 2018; Ferrara et al., 2016). The initiatives that fall into the first classification are those that use the social graph, or social connections, as the main element of the analysis (Alothal et al., 2018). In this process, relational information is highlighted, such as connections between accounts, posts and reposts, mentions and the use of common hashtags. In other words, everything that can demonstrate a connection between users and content. The other line of development

of these systems is based on *crowdsourcing*, i.e. involving the collaborative work of several human users in the task of identifying social robots.

This method is a hybrid between humans and non-humans, as it uses human analysis skills combined with computational strategies to standardize information at scale. On this point, Ferrara and other authors (2016) state that "(...) robot detection is a simple task for humans, whose ability to evaluate conversational nuances such as sarcasm or persuasive language, or to observe emerging patterns and anomalies, is as yet unparalleled in machines" (Ferrara et al., 2016, p. 101). The last group of methodologies is based on machine learning, a method that uses powerful computational resources to identify anomalous behavior (Alothali et al., 2018). The focus of this method is the processing of large volumes of data, facilitated by the choice of the type of information processed. According to Ferrara and other authors, approaches that use machine learning focus on behavioral pattern information, stating that such patterns can be coded and assimilated by machines to distinguish between humans and social robots (Ferrara et al., 2016).

The four detection platforms analyzed are mainly graph-based and machine-learning, focusing on developing automated detection applications with little dependence on human action. This feature is due to the cost of maintaining entire teams to analyze the high volume of information produced on social media. That is why hybrid methods, which guarantee more satisfactory results and are adaptable to the updating processes of social robots, are outside the leading platforms implemented. Beyond this difference, a common element between the different methods is the concern of increasing data analysis capacity while reducing computational processing. This challenge is presented by the continuous growth in information sharing on social media.

## 6 Final Considerations

Faced with this scenario where machines learn human behaviors and mimic their virtual presence, debates on social media are increasingly susceptible to massive manipulation processes. In this environment, threats to democracies in Latin America and the Caribbean are growing, ingeniously coordinated by extreme right-wing groups whose control dynamics are based on political polarization. Therefore, understanding the main strategies that are at work in these media is essential if we are to find new solutions to these problems.

Following this path, we chose to inspect the black box of the term increasingly used in the fields of Social Sciences, social bots or social robots. Thus, with the help of ANT, it was possible to understand that this is a complex category encompassing various types of practices involving people and algorithms with varied functions, created to adapt to every new platform targeting to detect manipulation. By understanding technical mediation and its four pillars, it became possible to observe the complex network of agents and agencies that circulate in these processes, noting that the term social robot can no longer contain all the meanings needed to describe the strategies implemented. Instead, it should be understood as an umbrella concept, encompassing everything from completely autonomous profiles to cyborgs that function in a hybrid way. Among these hybrid profiles, sockpuppets and meatpuppets result from the most current strategies for evading detection platforms.

Therefore, the analysis of the methods of the four platforms mentioned in this study revealed that these systems focus on the behavioral patterns of suspicious profiles, monitoring the frequency of posts, profile information, and the network of relationships built up by these agents. However, pinpointing a rigid standard that defines an automated account is challenging despite similarities, especially given the hybrid strategies of cyborg profiles. Platforms mainly use graph-based and machine-learning methods to deal with the growing volume of information on social media, relying on automated computing power to process large volumes of data. These methods, however, find it more challenging to deal with hybrid profiles due to their capacity for adaptation coordinated by human actions.

Finally, in the current context of tensions, where extremist political groups and large corporations threaten the institutions and sovereignty of southern countries, digital dynamics on social media are a complex battlefield with plenty of network actors that must be analyzed from an interdisciplinary perspective. The following ways range from the urgent need to regulate social media to holding the large corporations that own these platforms accountable. If the scenario remains as it is, the manipulation ecosystem may not be threatened, as detection initiatives are divided between those operated by groups that implement transparent and scientifically rigorous tools, which have a small impact. Others are linked directly to large corporations but are used as a subterfuge to respond to external pressures for control and the fight against disinformation.

## References

- Alothali, E., Zaki, N., Mohamed, E. A., & Alashwal, H. (2018). Detecting social bots on twitter: a literature review. *Proceedings of the International conference on innovations in information technology (IIT)*, 175–180. <https://ieeexplore.ieee.org/document/8605995>
- Azevedo Júnior, A. C., & Lourenço, R. F. (2023). Lideranças populistas, firehosing e a dinâmica algorítmica: um estudo dos posicionamentos de Jair Bolsonaro. *Más Poder Local*, (54), 96–123. <https://doi.org/10.56151/maspoderlocal.150>.
- Bessi, A., & Ferrara, E. (2016). Social bots distort the 2016 US Presidential election online discussion. *First Monday*, 21(11), 7–11. <https://doi.org/10.5210/fm.v21i11.7090>
- Callon, M. (2004). Por uma abordagem da ciência, da inovação e do mercado. O papel das redes sócio-técnicas. In A. Parente (Org.), *Tramas da Rede* (pp. 64–79). Porto Alegre: Sulina.
- Davis, C. A., Varol, O., Ferrara, E., Flammini, A., & Menczer, F. (2016). Botornot: A system to evaluate social bots. In J. Bourdeau, J. A. Hender & R. N. Nkambou (Eds.) *Proceedings of the 25th international conference companion on world wide web* (pp. 273–274). International World Wide Web Conferences Steering Committee. <https://doi.org/10.1145/2872518.2889302>
- Duarte F., J. I., Rodríguez G., G. E., Lares, J., & Sosa B., J. R. (2017). Venezolanos en Twitter: ¿Humanos, Bots o Ciborgs? Modelo de Clasificación. *Tekhné*, 1(19), 47–59. <https://doi.org/10.62876/tekhn.v1i19.3309>
- Eisenhardt, K. M. (1989). Building Theories from Case Study Research. *The Academy of Management Review*, 14(4), 532–550. <https://doi.org/10.2307/258557>
- Esquivel, E. (2022). La Manipulación en redes socio digitales. Una aproximación a sus estrategias. In A. C. Azevedo Júnior & L. Panke (Eds.), *Eleições, Propaganda e Desinformação* (pp. 85–98). Paraíba: EDUEPB.
- Ferrara, E., Varol, O., Davis, C., Menczer, F., & Flammini, A. (2016). The rise of social bots. *Communications of the ACM*, 59(7), 96–104. <https://doi.org/10.1145/2818717>
- Hui, P., Yang, K-C., Torres-Lugo, C., Monroe, Z., McCarty, M., Serrette, B., Pentchev, V. & Menczer, F. (2019). BotSlayer: real-time detection of bot amplification on Twitter. *Journal of Open Source Software*, 4(42), 1706. <https://doi.org/10.21105/joss.01706>
- Instituto Tricontinental de Pesquisa Social. (2021). Novas roupas, velhos fios: a perigosa ofensiva das direitas na América Latina. *Dossiê nº 47 do Instituto Tricontinental de Pesquisa Social*. <https://thetricontinental.org/pt-pt/dossie-47-ofensiva-da-direita-na-america-latinaamerica-latina/>
- Latour, B. (2012). *Reagregando o social: uma introdução à teoria do ator-rede*. Bahia: Edufba.
- Law, J. (1992) Notes on the theory of the actor-network: Ordering, strategy, and heterogeneity. *Systems Practice*, (5), 379–393. <https://doi.org/10.1007/BF01059830>.
- Lemos, A. (2013). Espaço, mídia locativa e teoria ator-rede. *Galáxia*, 13(25), 52–65. <https://revistas.pucsp.br/index.php/galaxia/article/view/13635/11399>
- Liu, D., Wu, Q., Han, W. & Zhou, B. (2016). Sockpuppet gang detection on social media sites. *Frontiers of Computer Science*, (10), 124–135. <https://doi.org/10.1007/s11704-015-4287-7>
- Melo, M. de F. A. de Q. e. (2011). A pipa e os quatro significados da mediação sociotécnica: articulações possíveis entre a Educação e a Psicologia para o estudo de um brinquedo. *Revista Brasileira de Pesquisa em Educação em Ciências*, 10(2). <https://periodicos.ufmg.br/index.php/rbpec/article/view/3982>

- Ruediger, M. A., Grassi, A., Freitas, A., Contarato, A. S., Silva, D. C., Beltrão, K., Calil, L., Silva, L. R., Barboza, P., & Bastos, R. (2017). *Robôs, redes sociais e política no Brasil: estudo sobre interferências ilegítimas no debate público na web, riscos à democracia e processo eleitoral de 2018* (v. 2). Rio de Janeiro: FGV DAPP. <https://hdl.handle.net/10438/24843>
- Santaella, L., & Cardoso, T. (2015). O desconcertante conceito de mediação técnica em Bruno Latour. *MATRIZes*, 9(1), 167–185. <https://doi.org/10.11606/issn.1982-8160.v9i1p167-185>
- Solorio, T., Hasan, R., & Mizan, M. (2013). A case study of sockpuppet detection in Wikipedia. *Proceedings of the Workshop on Language Analysis in Social Media*, 59–68. <https://aclanthology.org/W13-1107.pdf>
- Teixeira, V. C. (2018). Competição Eleitoral no Cenário Brasileiro Utilizando a Internet: Ágora ou Clientela. *Esferas*, 1(12), 9–18. <https://doi.org/10.31501/esf.v1i12.8267>
- Yañez, M. V. (2022). ¿Es un recurso el discurso de Fraude electoral? Elecciones en el continente americano 2019–2021. In A. C. Azevedo Júnior & L. Panke (Eds.), *Eleições, Propaganda e Desinformação* (pp. 153–180). Paraíba: EDUEPB.
- Yin, R. K. (2009). *Case Study Research Design and Methods* (5th ed.). Thousand Oaks, CA: Sage Publications. <https://utppublishing.com/doi/10.3138/cjpe.30.1.108>
- Zago, G. D. S. (2012). Trolls e jornalismo no Twitter. *Estudos em jornalismo e mídia*, 9(1). <https://doi.org/10.5007/1984-6924.2012v9n1p150>