

Detection and Characterization of Coordinated Online Behavior: A Survey

LORENZO MANNOCCI, University of Pisa, Italy and Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy

MICHELE MAZZA, Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy

ANNA MONREALE, University of Pisa, Italy

MAURIZIO TESCONI, Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy

STEFANO CRESCI, Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy

Coordination is a fundamental aspect of life. The advent of social media has made it integral also to online human interactions, such as those that characterize thriving online communities and social movements. At the same time, coordination is also core to effective disinformation, manipulation, and hate campaigns. This survey collects, categorizes, and critically discusses the body of work produced as a result of the growing interest on coordinated online behavior. We reconcile industry and academic definitions, propose a comprehensive framework to study coordinated online behavior, and review and critically discuss the existing detection and characterization methods. Our analysis identifies open challenges and promising directions of research, serving as a guide for scholars, practitioners, and policymakers in understanding and addressing the complexities inherent to online coordination.

Additional Key Words and Phrases: Coordinated behavior, coordination, synchronization, online social networks, network science

ACM Reference Format:

Lorenzo Mannocci, Michele Mazza, Anna Monreale, Maurizio Tesconi, and Stefano Cresci. 2024. Detection and Characterization of Coordinated Online Behavior: A Survey. 1, 1 (August 2024), 36 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Coordination, the process in which multiple connected actors are involved to pursue goals [79], is a fundamental aspect in the existence of various life forms, including human beings. From flocks of birds engaging in synchronized flight to insects working together in colonies, coordination enhances efficiency, safety, and resource utilization [15]. The ability to coordinate actions boosts the chances to overcome environmental challenges, fostering not only individual survival but also the resilience and success of entire communities. For these reasons, human coordination has been extensively scrutinized in multiple scientific disciplines interested in the dynamics of our offline interactions [114, 127].

With the advent of social media platforms, coordination has also become a fundamental component of *online* interactions. Social media users are now provided with a broad array of tools to coordinate with each other, such as hashtags that enable them to collectively discuss specific topics [12, 121]. Online platforms have become a suitable

Authors' addresses: **Lorenzo Mannocci**, University of Pisa, Italy and Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy, lorenzo.mannocci@phd.unipi.it, lorenzo.mannocci@iit.cnr.it; **Michele Mazza**, Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy, michele.mazza@iit.cnr.it; **Anna Monreale**, University of Pisa, Italy, anna.monreale@unipi.it; **Maurizio Tesconi**, Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy, maurizio.tesconi@iit.cnr.it; **Stefano Cresci**, Institute for Informatics and Telematics, National Research Council (IIT-CNR), Italy, stefano.cresci@iit.cnr.it.

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.

© 2024 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

1

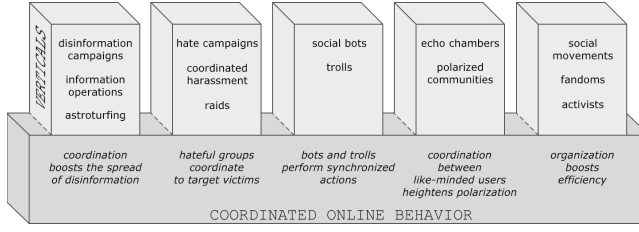


Fig. 1. Coordination is a fundamental aspect of online human interactions and the study of coordinated online behavior can complement the analyses of many other online phenomena.

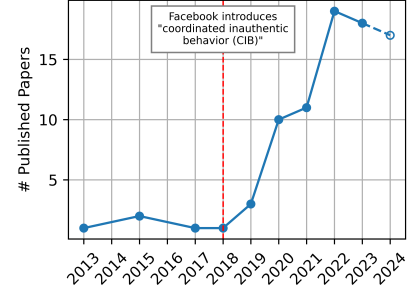


Fig. 2. Number of articles published yearly on coordinated online behavior. A steep rise is observed after Facebook introduced CIB in 2018.

environment for organizing social and political movements worldwide, giving rise to phenomena such as online activism [89, 98], boycotts [9, 70], and protests [74, 118]. The 2011 Arab Springs are a notable example, being largely organized through social media [54]. At the same time however, scholars found evidence of online coordination being exploited by nefarious actors for all sorts of malicious purposes. For instance, disinformation campaigns often leverage actors that coordinate their actions to maximize the outreach of their false narratives [62, 113, 117, 130]. Similarly, coordination is employed within information operations [18, 95, 130] and astroturfing, which involves creating the false appearance of grassroots support for a target cause, product, or person [62]. Also social bots and trolls exploit online coordination to amplify messages, manipulate trends, or spread disinformation [71, 82, 85, 141]. Finally, coordination also results from, and contributes to, the formation of echo chambers and online polarization [131]. Figure 1 highlights the complex, multifaceted, and intertwined nature of online coordination, which represents the underlying substrate of many diverse yet interconnected phenomena that permeate online social environments.

Recognizing the profound impact of online coordination on social media as well as its consequences on the offline world, both researchers [98, 100, 130, 138] and industry stakeholders [35, 36, 39] devoted a great deal of efforts to study its dynamics and to develop effective strategies to detect and mitigate its malicious instances. Research sped up significantly after 2018, when Facebook introduced the concept of *coordinated inauthentic behavior* (CIB) [35], marking a milestone in the development of the field. As shown in Figure 2, subsequent years saw a surge of interest on the subject, testified by the steadily growing number of published papers. This survey is motivated by this thriving interest on online coordination, which resulted in the availability of a large body of work. However, while Facebook’s interest towards online coordination was constrained to inauthentic behaviors as a response to the threat of orchestrated campaigns [35], here we embrace a more holistic and unbiased view by focusing on the broader concept of *coordinated behavior*. This inclusive approach allows for the analysis of a broader spectrum of works, including those focused on legitimate collective actions, offering a more comprehensive understanding of the coordination dynamics that shape digital spaces and fostering nuanced perspectives that go beyond mere threat detection. In spite of the many efforts, the existing literature still reflects the complexities and ambiguities surrounding this phenomenon. Among them is the limited agreement on a shared definition, which also hinders operationalization. Complexity also emerges from the diversity of methods proposed for detecting and characterizing coordinated online behavior, which impairs comparisons between different works and limits the generalizability of findings. Finally, the use of the same coordination technique by actors with disparate motivations poses challenges to estimating the impact and effects of online coordination.

This survey offers an extensive overview and critical analysis on coordinated behavior, starting from reconciling existing definitions from industry and academia, and proposing a general and comprehensive conceptual framework. We systematically analyze and categorize existing approaches for detecting and characterizing online coordination, elucidating open challenges and delineating promising directions for future research. Our work provides a roadmap for scholars, practitioners, and policymakers navigating the evolving complexities of coordinated online behavior.

Significance. Comprehensively modeling coordinated online behavior has far-reaching implications. On a theoretical level, it reconciles diverse definitions and provides a foundational framework for future research. On a technical level, it critically evaluates and categorizes existing detection and characterization methods, informing the development of the next generation of robust and adaptive tools for studying both malicious and neutral instances of coordinated behavior. By shedding light on online coordination—a fundamental dynamic of computer-mediated human behavior—this survey contributes to safeguarding online integrity and to fostering positive interactions in digital spaces. It offers valuable insights for shaping future methodologies, platforms, and policies, and it also contributes to enriching the interdisciplinary research occurring at the intersection between computer science and social dynamics.

Scope. Coordinated online behavior is orthogonal to many of the research topics tackled in fields such as Web and social media analysis, online social networks security, as well as social computing and computational social science, as shown in Figure 1. Therefore, a large number of works from these scientific communities implicitly or explicitly deal with online coordination. However, this survey is constrained to those papers that address online coordination explicitly and that provide relevant contributions to its detection, characterization, or understanding. In practice, we identified an initial set of candidate papers by selecting from Google Scholar and Scopus all those papers including the term “behavior” and at least one term among “coordinated”, “inauthentic”, “collaborative”. Each candidate paper was then evaluated by one of the authors of this survey to filter out unrelated works. Furthermore, we also manually checked all references from each of the related works retained after the previous filtering step, in order to identify possible additional related works. The final set of works categorized and critically discussed in this survey includes 84 papers published between 2014 and 2024, as shown in Figure 2.

Organization. This survey is structured as follows. Section 2 presents the theoretical foundations of coordinated online behavior, proposing a new general definition and laying out a comprehensive conceptual framework. Section 3 bridges the theoretical and methodological parts by defining the detection and characterization tasks. Section 4 discusses the existing literature on the detection task, while Section 5 focuses on the characterization task. Section 6 summarizes the main outstanding challenges and suggests promising directions of research. Finally, Section 7 concludes the survey.

2 CONCEPTUAL FRAMEWORK

The study of online coordination has its roots in the earlier studies of offline coordination. Recently, this study was advanced both by commercial platforms and academia, with a large array of different proposals. This section examines previously proposed definitions of coordination and discusses their advantages and limitations. Based on the results of this analysis, we then propose a general definition and a comprehensive conceptual framework.

2.1 Offline coordination

Coordination has already been extensively studied well before the emergence of social media across disciplines such as computer science, organization theory, management science, economics, and psychology [79, 114, 127]. Although

the meaning of coordination is intuitive, researchers suggested many definitions to frame the concept. A concise and precise definition was given in [77]:

Definition 2.1. Coordination (1988): *The additional information processing performed when multiple, connected actors pursue goals that a single actor pursuing the same goals would not perform [77].*

Definition 2.1 denotes coordination as the organizational overhead that multiple actors incur into when pursuing goals together. We note that this and similar definitions [5, 6, 76, 79] implicitly leverage the fundamental components of coordination, which we explicitly define as follows:

Definition 2.2. Coordination components: *A set of two or more actors who perform activities in order to achieve goals.*

Definition 2.2 introduces the fundamental components of coordination: *actors*, *activities*, and *goals*. Being coordination a nuanced concept, the theoretical modeling of these components can have major implications on downstream analyses and results. For example, in the case of communities or groups of users, each user in the group can be treated as a standalone actor, or alternatively the entire group may be considered as a single actor. Similar choices must be made when modeling the activities that allow actors to coordinate. Each actor typically performs multiple activities during any given time frame, and each of these activities might contribute differently to the overall coordination. Therefore, the choice of activities to model during an analysis directly impacts the resulting observed coordination [78]. Finally, the analyst is often interested in knowing the goal that the actors pursue when performing the activities. However, the actors may not all have the same goal, or even have any explicit goal at all [77]. These reflections on the components of coordination surface some of the challenges that early scholars faced since the 80s when studying offline coordination. Interestingly, many of these challenges carry over to the study of online coordination, informing the development of a new comprehensive definition and conceptual framework.

2.2 Concepts and definitions by online platforms

Beginning around 2016, mounting societal pressure impelled major social media platforms to confront pervasive challenges such as the organized dissemination of mis- and disinformation [135]. In consequence of this pressure, each platform adopted disclosure practices to communicate their results at exposing orchestrated deceptive activities perpetrated by organized actors. Given the importance of coordination for the success of large-scale disinformation campaigns [98, 101], within these public disclosure initiatives each platform addressed some instances of malicious online coordination. This section explores the concepts and definitions introduced by major social media platforms that are related to online coordination, discussing both their merits and limitations.

2.2.1 Facebook/Meta. After the public disclosure that the Internet Research Agency (IRA) had strategically exploited the platform to influence the 2016 US presidential election [99], Facebook began publishing reports detailing how their services were abused and the actions taken in response. In unveiling further actions against the IRA, in July 2018 Facebook introduced the concept of *coordinated inauthentic behavior* (CIB) [36]. A few months later, they supplied it with a first definition, and in October 2019 with a second one.

Definition 2.3. Coordinated inauthentic behavior (2018): *Groups of pages or people working together to mislead others on who they are or what they are doing [35].*

Definition 2.4. Coordinated inauthentic behavior (2019): *The use of multiple Facebook or Instagram assets (accounts, pages, groups, or events), working in concert to engage in inauthentic behavior, i.e., to mislead people or Facebook, where the use of fake accounts is central to the operation [38].*

Definition 2.3 underscores the collaborative nature of disinformation campaigns [117], emphasizing the objective of misleading others about the purported identity of the involved actors. To this end, it introduces the concept of *inauthenticity* of the actors, which is ever since often used in conjunction with the notion of coordination. Definition 2.4 explicitly articulates the concept of inauthenticity and elucidates that the act of deceiving others involves the extensive use of fake accounts [37]. A widespread critique of these initial definitions is that they exclusively address CIB, overlooking other types of malicious and possibly harmful coordination, let alone the neutral or benign ones [19, 47, 52, 98]. In September 2021, Facebook provided additional definitions focusing on harmfulness rather than inauthenticity.

Definition 2.5. Coordinated social harm (2021): *Networks of primarily authentic users who organize to systematically violate policies to cause harm on or off the platform [41].*

Definition 2.6. Coordinated mass harassment (2021): *Coordinated efforts of mass harassment that target individuals at heightened risk of offline harm [40].*

Definitions 2.5 and 2.6 adopt the concept of harmfulness in place of inauthenticity, thereby broadening the scope to also encompass authentic yet coordinated actors. These, in fact, hold the potential to cause negative consequences both on and off the platforms, as underscored in Definition 2.5.

2.2.2 *Twitter/X*. In October 2018, the platform released a public archive containing data about identified *information operations* (IOs).¹ Although not explicitly stated in Twitter’s definition at the time, a certain degree of coordination is necessary for the success of an IO [18, 130]. However, it was not until January 2021 that Twitter adopted a similar approach to Facebook and released a definition that explicitly addresses instances where coordination is leveraged to cause harm both online and offline.

Definition 2.7. Information operation (2018): *People directly involved in manipulation that can be reliably attributed to a government or state-linked actor [128].*

Definition 2.8. Coordinated harmful activity (2021): *Groups, movements, or campaigns that are engaged in coordinated activity resulting in harm on and off of Twitter [129].*

2.2.3 *YouTube/Google*. In its reports, primarily concerning abuses that occurred on YouTube, Google makes reference to *coordinated influence operations* (CIO) [49]. Even though Google did not provide a definition for CIOs, they nonetheless highlighted the importance of coordination in these online manipulations.

2.2.4 *Reddit*. In contrast to other platforms, Reddit embraced the broad concept of *content manipulation* to characterize the campaigns that violate its rules [104]. The platform shared only a small number of such campaigns, leaving it unclear whether these represent the entirety of the identified cases, or only a selection of them. Despite the absence of an explicit reference to coordination, large-scale content manipulations gain advantage from coordinated activities, similarly to what we discussed earlier about IOs.

Definition 2.9. Content manipulation (2022): *Things like spam, community interference, vote manipulation, and other attempts to artificially promote content [104].*

2.2.5 *Ambiguities and limitations*. This brief survey of the main concepts and definitions proposed by social media shows that online platforms are at the forefront in the analysis of coordinated online behavior. However, their conceptualizations are driven primarily by pressing practical regulation needs and by immediate contingencies, rather than by methodological rigor and theoretical soundness [29]. Faced with specific instances of malicious coordination,

¹<https://transparency.twitter.com/en/reports/moderation-research.html> (accessed: 31/07/2024)

Table 1. Examples of operational definitions used in recent academic literature. No definition is general enough to comprehensively describe coordinated online behavior. However, each definition grasps one or more relevant properties (highlighted in bold) that we leverage in our framework.

reference	definition
Nizzoli et al. [98]	Unexpected, suspicious, or exceptional similarity between a number of users
Cinelli et al. [19]	The number of times two accounts behaved similarly , such as when they repeatedly retweet the same post
Giglietto et al. [47]	The act of making people and/or things be involved in an organized cooperation
Magelinski and Carley [73]	Many instances of a tweet-behavior, i.e. tweeted hashtag [...] within a small predetermined time window
Weber and Neumann [137]	Anomalous levels of coincidental behavior
Magelinski et al. [74]	Users [that] take the same actions within minutes of one another
Zhang et al. [142]	Accounts that co-appear , or are synchronized in time
Pacheco et al. [101]	[Users exhibiting a] surprising lack of independence
Hristakieva et al. [55]	Coordination between users implies a shared intent
Keller et al. [62]	A group of people who want to convey specific information to an audience

online platforms adopt different and at times contrasting definitions, adding confusion and ambiguity to the already challenging task of defining an inherently nuanced and complex phenomenon. The current tangled landscape of platform definitions means that certain coordinated efforts may be categorized as such by some platforms but not by others, contingent upon the concepts they adhere to. What may be identified as coordinated behavior on one platform could be overlooked or dismissed on another, leading to inconsistencies in detection and mitigation. As a notable example, in June 2020 many teenagers organized on TikTok to reserve tickets for a Donald Trump rally to be held in Tulsa, OK (US). By mass-reserving and later cancelling their participation, they prevented others from making reservations and artificially inflated expected attendance numbers, ultimately causing an embarrassing number of empty seats at the rally [30]. Facebook’s head of security commented that, while tactical and sophisticated, they would have not acted upon this campaign as an instance of CIB (as per Definition 2.4), since it did not make use of fake accounts nor aimed to mislead Facebook users [30]. A similar case is the Chinese Spamouflage campaign, a cross-platform propaganda effort against the Hong Kong pro-democracy movement [97]. Various platforms reported takedowns of Spamouflage instances, with Google and Twitter labeling it as CIO and IO (as per Definition 2.7), respectively. However, Reddit interpreted Spamouflage differently, recognizing its low-quality and one-sided content, but refraining from considering these activities as rule violations.² Evidently, platforms address instances of coordinated online behavior disparately, and the criteria to establish the legitimacy of user behaviors lack objectivity. These discrepancies underscore the need for a unified understanding and standardized approach to defining and addressing online coordination, one that transcends platform-specific idiosyncrasies and fosters a more cohesive response to this challenge.

2.3 Concepts and definitions in academic literature

As shown in Figure 2, scientific interest in coordinated online behavior has drastically risen in the last few years. However, akin to social media platforms, scholars have generally refrained from proposing theoretically-grounded and general definitions. The majority of the existing studies adopt Facebook’s Definition 2.3 of CIB [84, 86]. Instead, those who propose their own conceptualization mainly provide *operational* definitions useful for the development of coordination detection methods [74, 98, 101]. Table 1 reports some examples of operational definitions proposed recently. As shown, no definition is comprehensive and general enough to adequately describe the multifaceted phenomenon of coordinated online behavior. The existing conceptualizations of the phenomenon appear to be influenced by the specific technique employed for its detection. Consequently, the criteria used to define coordination vary, encompassing aspects

²https://www.reddit.com/r/redditsecurity/comments/dp9nbg/reddit_security_report_october_30_2019/ (accessed: 31/07/2024)

ranging from similar behavior [19, 73, 74, 98, 101, 137] to synchronicity [73, 74, 142] and common intent [55, 62]. As practical examples of the above, some definitions and the resulting detection methods revolve around the anomalous use of hashtags or URLs by multiple users [47, 74], or repeated screen name swapping [101]. However, while no single definition comprehensively describes online coordination, each grasps one or more relevant properties, as highlighted in table. In turn, these properties can inform a general definition coordinated online behavior.

2.4 General definition and fundamental components of coordinated online behavior

We propose a new general definition of coordinated online behavior that overcomes the ambiguities and limitations of the existing definitions. The new definition leverages the components of offline coordination introduced in Section 2.1 and is informed by the various operational definitions proposed by social media platforms and by the academic literature, respectively discussed in Sections 2.2 and 2.3.

Definition 2.10. Coordinated online behavior: *A group of users who perform synergic actions in pursuit of an intent.*

actors
actions
intent

Definition 2.10 delineates coordinated online behavior based on three fundamental components—*actors*, *actions*, and *intent*—that are similar to the components of offline coordination outlined in Definition 2.2. Definition 2.10 and its three fundamental components enable the comprehensive mapping of all instances of online coordination, as discussed in the following.

2.4.1 Actors. Actors refer to the individuals or entities that are engaged in coordinated behavior. The attributes of the actors contribute to characterizing instances of coordination. For example, the way in which the actors represent themselves to those not involved in the coordinated behavior determines whether the coordination is authentic or otherwise. All instances of online coordination where the actors misrepresent themselves, as in the case of social bots [50, 55, 98, 101] and state-backed trolls [88], are cases of inauthentic coordination. Conversely, coordination among actors who accurately self-portray is deemed authentic [50]. In addition, the relationships between the actors determine whether the coordination is spontaneous, grassroots, or emergent (i.e., bottom-up) [98] or whether it is structured and well-organized (i.e., top-down) [130]. Finally, the number of the involved actors determines the scale of the coordination.

2.4.2 Actions. Actions represent the practical means that allow actors to coordinate. In coordinated behavior the actions are synergic, in that they are mutually reinforcing and potentially capable of producing a larger effect than that obtainable by individual actions alone. While actors and intent can be misrepresented or concealed, actions are typically visible and non-falsifiable. In other words, the actions with which the actors coordinate represent the digital breadcrumbs of the coordination. For this reason, the actions are the component based on which the majority of coordination detection methods are developed. Furthermore, the types and attributes of the actions also provide information towards characterizing instances of coordination. For example, the timing and synchronization of the actions among actors indicates the degree of planning and organization involved [74, 142]. The consistency of the actions and the actions performed in response to external events are further characteristics of coordinated behaviors. Finally, the types of actions and their content provides insights into the intent and goals of the actors [19].

2.4.3 Intent. The intent is the objective that the actors pursue when they coordinate. When studying coordinated online behavior, the intent of the actors is typically unknown, if not deliberately concealed. For example, actors involved in malicious or harmful coordination conceal their intent to avoid being stopped in their endeavor [55]. However,

also actors involved in neutral or benign forms of coordination might not openly state their goals and intent. For this reason, the observer often tries to infer the intent based on the visible actions performed by the actors. Furthermore, intent can be either shared and explicit among the actors, or implicit. For instance, the perpetrators of a disinformation campaign share the explicit goal of disseminating certain pieces of false information [62]. Conversely, fans of a sports team or public character may spontaneously engage in coordinated actions to support their idol, without having agreed on a specific objective or course of action [98]. The previous examples highlight that the characteristics of the intent are related to the degree of organization of the actors. Importantly, the intent also contributes to determining the harmfulness of the coordinated behavior. However, while there are cases in which it is relatively straightforward to categorize a coordinated behavior as harmful or otherwise—think for example of coordinated hate attacks [83] or state-backed disinformation campaigns [18]—there also exist situations where the harmfulness of the intent is inherently subjective. Coordinated efforts to promote a controversial political ideology may be perceived as harmful by some, while others may view them as legitimate expressions of free speech [100]. Similarly, coordinated campaigns to boycott a company or criticize a public figure may be seen as harmful by those targeted, but supporters may genuinely view them as justified forms of activism [93].

2.5 Defining dimensions of coordinated online behavior

As discussed in the preceding sections, coordinated online behavior constitutes a complex and multifaceted phenomenon whose instances are contingent upon the actions and intent of the involved actors. Here we leverage our discussion about the fundamental components of online coordination to introduce four defining dimensions of this phenomenon: authenticity, harmfulness, orchestration, and time-variance.

2.5.1 Authenticity. Authenticity refers to the degree of genuineness and transparency that the actors exhibit in their actions and overall online presence. Coordinated authentic behavior is executed by genuine actors and typically emerges organically within a community of users who share common interests or beliefs. Examples of authentic coordination are activists, social movements, and mutual support groups, which are driven by motivations such as a desire for social or political change or the cultivation of a sense of community and belonging [98, 118]. While authentic forms of coordination are also harmless in the majority of cases, as in the previous examples, there also exist less frequent cases of authentic yet harmful behaviors. For instance, certain coordinated hate groups openly encourage racist, xenophobic, or supremacist ideologies [90, 93, 138]. Conversely, coordinated inauthentic behavior entails the use of fake accounts, such as social bots, trolls, and fake personas [23, 108]. These are typically employed for purposes such as spreading disinformation, sowing confusion, or eroding trust in democratic institutions. As such, inauthentic coordination is often characterized by its deceptive nature and aim to manipulate unaware users [117]. Nonetheless, there exist cases of inauthentic yet harmless coordination. For example, online participants in the Arab Spring movements concealed their identities to avoid government surveillance and potential reprisals [54]. While their coordinated efforts were inauthentic in terms of individual identity disclosure, they remained largely harmless in intent, aiming to promote democratic ideals, social justice, and human rights. These examples highlight the difference between authenticity and harmfulness, which represent two orthogonal dimensions of coordinated online behavior.

2.5.2 Harmfulness. Harmfulness refers to the negative impact, consequences, or outcomes—both online and offline—resulting from the coordinated actions of the actors. As discussed in Section 2.4.3, harmfulness depends both on the shared intent and actions of the actors engaged in coordination, and on the viewpoint of the observer, constituting a much more conceptually intricate dimension of online coordination than authenticity. However, in spite of the inherent

subjectivity, there exist many clear cut cases of harmful and harmless coordination. For example, coordinated actors involved in the spread of disinformation, hate speech, and online harassment, represent straightforward cases of harmful coordination [18, 95]. In contrast, users who coordinate to collect and share information and other resources, such as in the aftermath of mass emergencies, represent cases of harmless coordination [74, 98, 100].

2.5.3 *Orchestration.* Orchestration represents the degree of planning and organization between the coordinated actors. This dimension is closely linked to the intent of the actors, in that highly orchestrated campaigns typically imply shared intent and goals between the participants. The orchestration of a coordinated campaign can be centralized or distributed. In centralized orchestration, a single actor or entity exercises control and coordination over the actions of all other actors involved in the coordinated behavior. This centralized authority dictates the timing, content, and strategy of the coordinated actions, allowing for tight coordination and synchronization. An example of strong and centralized orchestration is the coordinated behavior exhibited by social botnets, where large groups of automated accounts quasi-simultaneously perform predefined actions depending on the command of a botmaster entity [85]. In decentralized orchestration, coordination and control are distributed among multiple actors within a network, without a single central authority dictating the actions of all participants. Actors may self-organize, collaborate, or communicate autonomously, often guided by shared goals, interests, or ideologies. For instance, in January 2021, retail investors coordinated on Reddit to target short-selling activity by hedge funds on GameStop shares, causing a surge in the share price and triggering significant losses for the funds involved [70]. Instead, non-orchestrated coordinated behavior occurs when the actions of multiple actors spontaneously converge around a given topic, narrative, or activity. Certain viral social media trends are an example of non-orchestrated coordination emerging from the widespread adoption of a particular hashtag or activity that occurs organically as many users observe and emulate others' behavior [118].

2.5.4 *Time-variance.* Time-variance refers to the temporal characteristics and the dynamic nature of coordinated online behavior. It grasps possible changes in the types, timing, frequency, and intensity of the actions, which in turn may reflect changes in the intent of the actors, as well as adaptations or responses to external stimuli. Examples of largely static coordinated behavior are the activities of some spammers and bots, who repeatedly perform the same actions adhering to a fixed pattern without much adaptation or variation [23, 100]. Conversely, many information operations are dynamic and time-varying, presenting different characteristics at different points in time. Among the changing characteristics are the types of actors involved in the coordination (e.g., whether automated or human-operated) or the topics of discussion [117, 130]. Time-variance also strongly depends on the duration of the coordinated behavior itself. Actors involved in certain state-sponsored disinformation campaigns operate on online platforms for extended periods, spanning years or even decades [64]. Over such lengthy time frames, the actors adapt their tactics, narratives, and targets in response to shifts in intent, changes in technology and platforms, or advancements in countermeasures. This extended duration implies a relatively gradual and nuanced evolution of the coordinated behavior. Conversely, other forms of online coordination rely on expendable or disposable accounts created for short-lived and fast-paced activities [8]. These actors are employed for specific tasks and then discarded or deactivated once their purpose is fulfilled or they are detected. As a result, these ephemeral instances of coordination are rapid, intense, and short-lived.

2.6 Taxonomy and final remarks

Our conceptual framework of coordinated online behavior encompasses the three fundamental components presented in Section 2.4 and the four defining dimensions discussed in Section 2.5, providing a general and flexible scheme for studying, categorizing, and comprehensively mapping the multiple instances of online coordination.

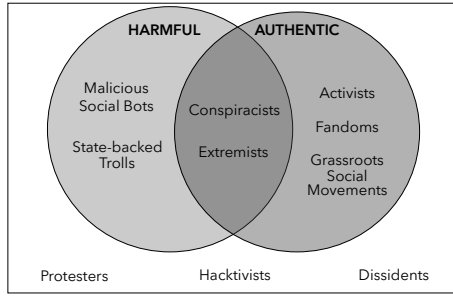


Fig. 3. Taxonomy of coordinated online behavior obtained by considering the dimensions of *harmfulness* and *authenticity* of our conceptual framework. The framework conveniently allows the mapping of disparate instances of online coordination.

often initiated by malicious social bots and state-sponsored trolls [1, 18, 130]. Other problematic phenomena are those featuring harmful yet authentic behaviors, such as the activity of conspiracy theorists and hateful extremists [83, 90, 93]. Conversely, the activity of grassroots social movements, fandoms, and other activists represent instances of harmless and authentic coordination [118]. Finally, harmless yet inauthentic behaviors lay outside of the partially overlapping sets, and are exemplified by hacktivists, dissidents, and other anonymous protesters [54]. Alternative taxonomies can be obtained by leveraging the dimensions of orchestration and time-variance, which would highlight additional phenomena to those shown in Figure 3.

2.6.2 Final remarks. The previous instantiation of our conceptual framework in a taxonomy based on the dimensions of harmfulness and authenticity concludes the theoretical part of the survey. The following section bridges the theoretical and methodological parts by presenting the problem definition. Subsequently, we systematically review the proposed methodologies for detecting and characterizing coordinated online behavior, showing how these tasks stem from the conceptual modeling of the phenomenon presented in this section.

3 PROBLEM DEFINITION

The problem of identifying and investigating different types of coordinated online behavior involves defining two functions $f(\cdot)$ and $g(\cdot)$ that respectively implement the tasks of coordinated behavior *detection* and *characterization*, as outlined in Figure 4. Given a set of users and their actions on one or more online platforms, $f(\cdot)$ identifies possible coordinated groups of users. Instead, $g(\cdot)$ extracts additional information for each detected group, thus contributing to determining the nature, intent, and the overall characteristics of the involved actors (e.g., whether they are inauthentic, harmful, etc.). The detection and characterization tasks are related to the Definition 2.10 of coordinated online behavior and its components in that the function $f(\cdot)$ implementing the detection task does so via the analysis of user *actions*, while the function $g(\cdot)$ implementing the characterization task provides information about the *actors* and their *intent*. A comprehensive overview of the entire process is depicted in Figure 4. The input is represented by the set of users U to analyze and their activities H . The detection task differentiates coordinated users from non-coordinated ones. Depending on the detection method, the distinction between the two can be expressed as binary labels assigned to the users, as two or more sets (e.g., clusters) of either coordinated or non-coordinated users, or as two or more coordinated or non-coordinated communities (i.e., nodes and edges) from a network. These are subsequently scrutinized during the

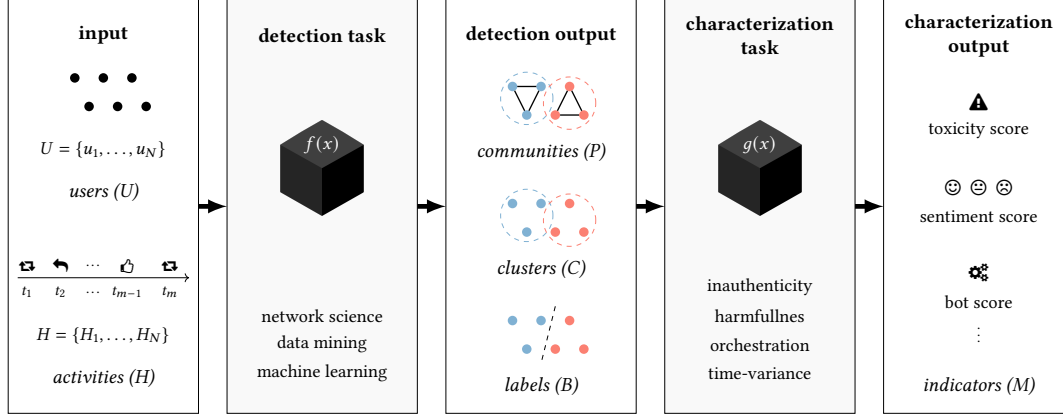


Fig. 4. The analytical process of studying coordinated online behavior, involving the *detection* and *characterization* tasks. The input to the overall process is a set of users U and their activities H on one or more platforms. The output of the detection task is either a set of binary labels B , clusters C , or network communities G that differentiate coordinated and non-coordinated users. The characterization task receives these in input and outputs a set of indicators M .

characterization task, which computes a set of indicators for each coordinated user, set, or community. The indicators are selected so as to provide information about the characteristics of the coordinated actors and their behavior. For example, computing bot scores is a common method to estimate the inauthenticity of coordinated users.

3.1 Detection of coordinated online behavior

3.1.1 Input. Let $I = \langle U, H \rangle$ be the problem input, where $U = \{u_1, \dots, u_N\}$ denotes the set of users and $H = [H^{u_1}, \dots, H^{u_N}]$ represents an ordered vector of activities performed by those users. We define the activity of a user u_j as $H^{u_j} = [h_1^{u_j}, \dots, h_T^{u_j}]$ representing the vector of chronologically ordered actions performed by u_j . An action is defined by the quadruple $h = \langle \text{type}, \text{target}, \text{content}, \text{timestamp} \rangle$ which describes the *type* of action executed by a user on a specific *target* or *content*, at a given *timestamp*. Users can execute actions of different *type* such as posting, resharing, befriending, and more. A *target* is another user of the platform who is affected by the action. For example, in the case of a retweet action on the platform Twitter/X, the target is the author of the retweeted tweet. For some actions the target is undefined, as in the case of the posting action. The *content* of an action is a post (e.g., a tweet, comment, submission, and more, depending on the platform). Posts contain one or more elements of content, such as text, image, URL, mention, hashtag, and more. In case the content contains multiple elements, the corresponding action is called compound action [74]. Similarly to the target, also the content can be undefined depending on the type of action, as in the case of a befriending or following action. To wrap up, the type of action and its timestamp are always defined, while one of content and target might be optional, depending on the type of action.

3.1.2 Methods and output. As presented in Section 4, most of the existing literature on coordinated behavior detection analyzes both the set of users and their actions. Some studies only leverage the content of the actions, without considering their type [25, 83, 133]. Furthermore, certain works do not take into consideration the timings of the actions [14, 17, 19, 83], while some others only consider the timings [8]. The task of detecting coordinated online behavior is modelled by the function $f(U, H)$, which can provide three different outputs depending on the adopted

method, corresponding to different levels of detail and information on the coordinated users:

$$f(U, H) = \begin{cases} P = \{P_1, \dots, P_i, \dots, P_k\}, & P_i = (V_i, E_i), V_i \subseteq U & \text{communities} \\ C = \{C_1, \dots, C_i, \dots, C_k\}, & C_i \subseteq U & \text{clusters} \\ B = \{B_c, B_u\}, & B_c \cup B_u = U & \text{binary labels} \end{cases} \quad (1)$$

In the most general case, the output of $f(\cdot)$ is a set P of *communities* of coordinated users. Coordinated communities P_i are sub-networks where the nodes are users from U , and the edges—with their weights—encode the level of coordination among the users. Communities are typically outputted by those methods that adopt an internal network representation, which is then analyzed with community detection algorithms. Coordinated communities are an information rich representation, given that the presence and weight of links between the coordinated users facilitates subsequent analyses, such as those needed for the characterization task. Another possible output consists of a set of *clusters* of users. The clusters C_i are produced by methods that adopt tabular representations of the users, which are then analyzed with clustering algorithms. These methods typically ignore the relationships between the users but are able to identify multiple groups of coordinated users. Finally, the least informative output is given by those methods based on classification algorithms. These methods assign *binary labels*, partitioning the initial set of users U in two labeled groups of coordinated (B_c) and non-coordinated (B_u) users. These labelled groups do not provide information about neither the relationships between the users nor the existence of multiple coordinated groups of users in U .

3.2 Characterization of coordinated online behavior

3.2.1 Input. The characterization task is modelled by the function $g(Y, H) = M$ whose inputs are the groups of coordinated users resulting from the detection task $f(U, H) = Y \in \{P, C, B\}$, defined in Eq. (1), with their activities H .

3.2.2 Methods and output. As discussed in Section 5, the characterization task aims at computing a set of quantitative indicators M to measure distinctive properties of the detected coordinated behaviors in terms of the defining dimensions that we presented in Section 2.5: authenticity, harmfulness, orchestration, and time-variance. The indicators that can be used in the characterization task partly depend on the methods and outputs of the detection task. For example, assortativity measures the extent to which nodes with a high degree in a network are connected to other nodes with a high degree, and vice versa. This indicator was used to gain insights into the inner structure and organization of certain coordinated communities [98]. However, assortativity can be computed only if the coordination detection method outputs communities, rather than clusters or binary labels. On the contrary, other indicators can be computed independently of the detection method, such as the aforementioned bot scores that are commonly used as an estimator of the inauthenticity of the coordinated users [50, 55, 98, 101]. The utility of the characterization task is not limited to shedding light on the nature of the detected coordinated behaviors nor to distinguishing between different instances of the phenomenon. In fact, the output of characterization task can also be leveraged to validate the output of the detection task, as in those frequent cases when a ground-truth of coordinated users is unavailable.

4 DETECTION OF COORDINATED BEHAVIOR

Coordination detection methods can be classified into two main categories depending on their underlying approach: network science or machine learning. The following sections discuss the existing solutions in each category.

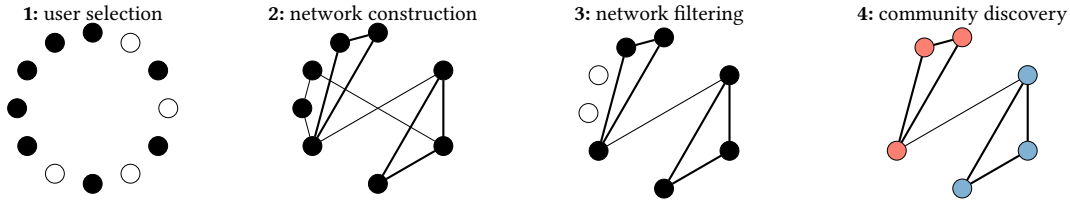


Fig. 5. Main steps of the network science methods for the detection of coordinated online behavior. 1: The selected users become nodes in a network. 2: User similarities are computed with a similarity function and assigned to the edge weights of the network. 3: The network is filtered so as to retain only similarities with given properties. 4: Community discovery is performed to detect groups of strongly coordinated users.

4.1 Network science methods

Network science coordination detection methods build a network of users or posts, where the links between the nodes in the network represent the presence, and possibly also the extent [98], of coordination. In spite of the existing differences, all methods in this category carry out the sequence of steps shown in Figure 5, namely: (i) user selection, (ii) coordination network construction, (iii) network filtering, and (iv) community discovery. We now discuss the objective and the implementation options available for each step.

4.1.1 User selection. This step selects an initial subset $U' \subseteq U$ of users according to some criteria that depend on the purpose of the analysis. This initial selection is motivated by the observation that a small fraction of users accounts for the majority of actions on a social network [102], especially those associated with harmful behaviors [106]. Selecting a subset of users also has the positive consequence of reducing the computational cost of the subsequent steps, which can easily become time- and computation-intensive for large networks [20].

Implementation. Multiple choices can be made to select a subset of relevant users. A common choice is to select the most active users as they produce the largest share of actions and content. Most active users can be defined as those who publish a large number of original posts (*super-producers*) [66, 73, 101], or as those having a large number of re-shares (*super-spreaders*) [19, 28, 55, 69, 87, 98, 120]. Other than activity, network centrality or influence, suspicious behavior, location, and timings are used for user selection [10, 72]. Furthermore, many other general criteria can also be adopted, such as selecting all users who posted certain keywords, or all followers of a given user. Combining multiple criteria allows for even more fine-grained user selections.

4.1.2 Coordination network construction. This step builds a coordination network between the users $U' \subseteq U$ previously selected.³ A coordination network is a type of network where links exist only between coordinated nodes. As per Definition 2.10, coordination implies synergic actions between users. In network science methods, this concept is operationalized with *co-actions*. A co-action represents two users performing the same action on the same target or content. For instance, two users who comment, like, or re-share the same post are generating a co-action. As shown in Figure 6, the two coordinated users need not be directly connected in the social or interaction network. This characteristic makes coordination networks particularly suitable for surfacing coordination between seemingly unrelated users, such as those involved in inauthentic or harmful behavior, so much so that some scholars specifically refer to *latent* coordination networks [138]. In the most general case, a coordination network is a multiplex network $G(V, E, W, \mathcal{L})$. In order to build G , one must define the types of co-actions C_a to consider (e.g., re-shares, mentions,

³A few works follow the general approach of network science methods, but build networks of *content* rather than *users*. These are discussed in Section 4.1.5.

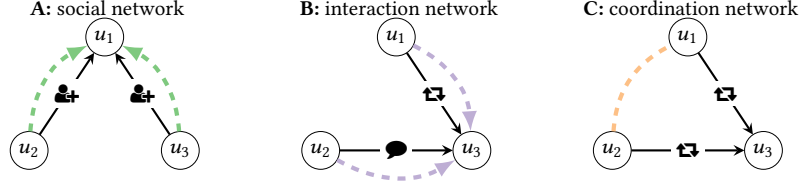


Fig. 6. Differences between social (A), interaction (B), and coordination (C) networks. Solid black edges represent actions on the online platform, while dashed colored edges show how actions are translated into edges in the corresponding type of network. Coordination networks are typically undirected and link users performing similar actions at around the same time. Differently to social and interaction networks, coordination networks allow connecting users even if they never directly interact with one another.

follows, etc.). When multiple co-actions are used, G is a multiplex network with L layers, where $\mathcal{L} = \{1, \dots, L\}$ is the set of layers, each corresponding to a co-action $i \in C_a$. However, the majority of existing works leverage a single co-action. In this case, the number of layers is $L = 1$ and G is a single layer network. Each layer in G is an undirected weighted graph $G^i(V^i, E^i, W^i)$, where $V^i \subseteq U'$, E^i and W^i respectively denote the set of nodes, edges, and weights of layer i . We highlight that $V = \bigcup_{i \in \mathcal{L}} V^i$, $E = \bigcup_{i \in \mathcal{L}} E^i$, and $W = \bigcup_{i \in \mathcal{L}} W^i$. Given a layer i , an edge e_{jk}^i is created if there exists a co-action of type i between two users (u_j, u_k) . The edge weight $w_{jk}^i = \text{sim}_i(u_j, u_k)$ is obtained via a similarity function that computes pairwise user similarities in U' . Different similarity functions can be used for different co-actions.

Implementation. Building the coordination network G requires defining the types of co-actions and the corresponding similarity functions. The vast majority of works in literature rely on a single co-action and the resulting networks are single-layered. Table 2 reports the main implementation details for the works that built single layer coordination networks. In table, when multiple co-actions are listed for the same author or work, this means that multiple single layer networks were built, rather than a multiplex network resulting from the simultaneous analysis of multiple co-actions. The few existing works that built multiplex coordination networks are instead described in Table 5.

As shown in Table 2, the most common type of co-action is *co-sharing* (e.g., the *co-retweet* action on Twitter/X) [17–19, 27, 28, 31, 50, 55, 61, 65, 67, 69, 72, 98, 101, 108, 109, 120, 121, 130, 136–138]. Other frequently used co-actions are *co-reply/co-comment* [63, 137, 138] and *co-like* [57], which occur when two users comment or leave a reaction to the same post. The previous co-actions are based on the type of interaction between users and content in an online platform. Other co-actions are instead based on the content of user posts. For example, *co-post* (e.g., *co-tweet* on Twitter/X and *co-parley* on Parler) [17, 26, 31, 51, 62, 109, 130, 132, 133], *co-image*, and *co-video* [140] represent the publishing of posts with the same text, image, or video by multiple users. More specific co-actions are also possible, such as *co-text-image* that occurs when multiple users post images that contain the same text [115]. Similarly, *co-mention* [1, 74, 87, 89, 90, 130, 137, 138], *co-URL* [1, 11, 14, 27, 31, 44–47, 52, 72–74, 89–91, 93, 130, 137, 138], and *co-hashtag* [1, 13, 27, 31, 72–74, 88–90, 101, 130, 134, 137, 138] represent two users publishing a post with the same user mention, URL, or hashtag. Regarding the latter, the majority of works consider publishing a post with a *single* common hashtag as a valid co-action [1, 31, 73, 74, 88–90, 130, 134, 137, 138]. However, others argued that using the same *set* or *sequence* of hashtags represents a stronger and more reliable signal of coordination [13, 27, 72, 101]. When a co-action is defined in such a way that multiple atomic actions are required for two users to be considered as coordinated, that co-action is a *compound* action [74]. More generally, compound actions refer to actions that involve the simultaneous occurrence or combination of multiple individual actions or sub-events. Therefore, a co-action requiring the posting of the same set or sequence of hashtags is a compound action composed of multiple homogeneous elements: $\langle \text{hashtag}, \text{hashtag}, \dots \rangle$. However, compound co-actions can also be defined with heterogeneous elements, such as $\langle \text{hashtag},$

Table 2. Network science methods for detecting coordinated behavior based on *single layer* user networks. For each group of works we report the considered co-actions, similarity functions, filtering criteria, and community detection methods.

reference	action	similarity	filters [†]	community detection
[18]	retweet	cardinality	threshold, ADJ	modularity clustering
[50]	retweet	cardinality	EDO	Louvain
[65]	retweet	cardinality	threshold, ADO	Louvain
[67]	retweet	cardinality	backbone, ADJ	Louvain
[108]	retweet	cardinality	EDO	
[19, 28, 55, 69, 98, 120]	retweet	cosine similarity TF-IDF	backbone	Louvain
[121]	retweet	cosine similarity TF-IDF	backbone, EDO	Leiden
[136]	retweet	cardinality	threshold, ADJ	
[17]	retweet, tweet	cardinality	threshold, EDO	Louvain, connected components
[62]	retweet, tweet	cardinality	threshold, EDO	Louvain
[109]	retweet, tweet	cardinality	threshold, EDO	
[26]	tweet	cardinality	ADO	
[66]	tweet	cardinality	threshold	cohesive campaign
[100]	tweet	text similarity	threshold, ADO	
[92, 94]	parley	cardinality	threshold, kNN graph	Leiden
[132]	text	cardinality	backbone	Louvain
[91]	tweet, URL	cardinality, cosine similarity	threshold, EDO	Louvain
[93]	tweet, parley, URL, username	text similarity, cardinality, cosine similarity	threshold, kNN graph	Louvain
[72]	retweet, tweet, URL, hashtag	cosine similarity TF-IDF, text similarity	threshold, ADO	
[137, 138]	retweet, URL, hashtag, mention, reply	cardinality	ADJ	focal structures
[130]	retweet, tweet, URL, hashtag, mention	cosine similarity TF-IDF	ADJ	Leiden
[101]	retweet, hashtag, image, handle change, synchronization	Jaccard coefficient, cardinality, cosine similarity	threshold, EDO	
[14]	URL	cardinality	kNN graph	Louvain
[11, 44, 46, 47, 52, 105]	URL	cardinality	threshold, ADO	connected components
[1]	URL, hashtag, mention	unweighted	EDO	Leiden
[90]	URL, hashtag, mention	cardinality	threshold, EDO	Louvain
[45, 115]	URL, text-image	cardinality	threshold, ADO	connected components
[95]	image	cardinality	threshold, kNN graph	
[140]	image, video	cardinality	threshold, ADO	connected components
[13]	hashtag	cardinality	threshold	connected components
[88]	hashtag	cardinality	threshold, backbone	Louvain
[134]	hashtag	cardinality		
[57]	like	cardinality	threshold	
[63]	comment	cardinality	threshold	k-means, hierarchical clustering
[87]	mention	cosine similarity TF-IDF	threshold	Louvain

[†] ADJ: adjacent time window, EDO: evenly distributed overlapping time window, ADO: action-driven overlapping time window

mention) or *(URL, mention)* [74]. When building coordination networks, the use of compound actions increases the confidence of labelling a group of users as coordinated, since compound actions are less likely to occur by chance. However, this approach risks neglecting simpler, milder, or looser forms of coordination. Lastly, *co-handle changes* refer to multiple users using the same handle or username at different points in time [101].

After selecting one or more co-actions to identify similar user behaviors, it is necessary to define the corresponding similarity functions. A similarity function computes the weight w_{jk} of the edges connecting two coordinated users u_j and u_k based on how similar their behaviors are according to the chosen co-action. Independently of the type of co-action, the majority of works use similarity functions based on the *cardinality* of the co-action between the

two users [11, 13, 14, 17, 18, 26, 31, 44–47, 50–52, 58, 62, 63, 65, 67, 91–95, 105, 108, 109, 115, 120, 132, 134, 136–138, 140]. For example, when using *co-hashtags*, the similarity function may simply count the number of times the two users used the same hashtags. Another widely adopted function is the *cosine similarity* between the two user vectors [19, 27, 27, 28, 55, 69, 72, 87, 91, 93, 98, 101, 121, 130]. These can be binary vectors, frequency vectors, or TF-IDF weighted vectors. The latter allows discounting the importance of popular or viral content, boosting instead the relevance of unpopular items [19, 27, 28, 55, 69, 72, 87, 98, 121, 130]. Depending on the choice of co-action, the previous similarity functions must be preceded by additional processing steps. For example, the use of co-actions such as *co-post*, *co-image*, and *co-video* involve counting the number of times two users shared the same content, which severely limits the possibility to detect certain coordinated behaviors. For this reason, some scholars relaxed this requirement by also considering the posting of *similar*—as opposed to *equal*—content. This solution requires defining additional similarity functions for texts, images, and videos. In literature, text similarity was computed via correlation [66], cosine similarity of document embeddings [27, 31, 72, 91, 92, 94], Ratcliff/Obershelp algorithm [100], or Jaccard similarity [51]. All these works first computed similarities between the text of two users’ posts. Then, they set a threshold to select highly similar texts. Finally, they computed user similarities based on the number of similar texts [31, 51, 91, 92, 94], or as the average of the similarities between the highly similar texts [27, 72]. Analogously, *co-image* involves a pre-processing step for representing the images with their embeddings [95] or with RGB color histograms [101]. Then, image similarity is computed via measures such as the Euclidean distance [95]. Finally, a similarity threshold is applied and user similarity is computed as the cardinality [95] or the Jaccard coefficient [101] of the sets of similar images.

4.1.3 Network filtering. This optional, yet crucial, step allows to filter out nodes and edges from the coordination network so as to only retain network structures that convey meaningful and reliable coordination. Independently of the method used to achieve this objective, performing network filtering also has the advantage of reducing the size of the final coordination network, which speeds-up subsequent analyses.



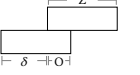
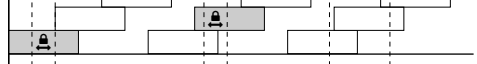
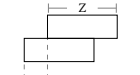
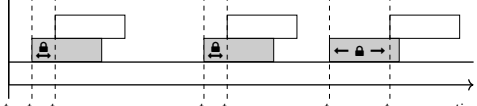
Implementation. There are three main approaches for coordination network filtering: (i) fixed thresholds, (ii) statistical validation, and (iii) the timings of the co-actions. Methods based on fixed similarity thresholds discard all edges in the network whose weight $w < w_{th}$, and all the resulting disconnected nodes [11, 13, 17, 18, 23, 31, 44–47, 52, 57, 62, 63, 65, 72, 87, 88, 90–92, 94, 95, 100, 101, 109, 115, 133, 136, 140]. The threshold w_{th} is chosen in such a way to retain only strongly coordinated users, as they are implicitly considered to be more relevant. Relevant nodes can also be identified via eigenvector centrality, by pruning those nodes that do not exceed a certain threshold [27, 72]. The similarity and centrality thresholds are typically selected arbitrarily, without a strong underlying theoretical motivation [98]. Moreover, coordination networks resulting from the analysis of different datasets, or from the use of different types of co-actions, inevitably result in different edge weight and node centrality distributions. This variability makes the repeated use of “standard” thresholds unsuitable and mandates in-depth case-by-case analyses. To alleviate this burden some works leverage *k*-nearest neighbors graphs (*k*-NNG), where two users u_j and u_k are connected only if u_j is among the *k*-nearest neighbors of u_k [14, 92, 94, 95]. This filtering operation retains only the *k* strongest neighborhoods in the coordination network, thus reducing the emphasis on the edge weight. However, determining the best value for *k* is also challenging since *k* strongly depends on the characteristics of the different networks and their layers. Other methods for identifying relevant network structures are those that retain statistically-meaningful edges, independently of their weight [43, 111, 126]. This approach, used by several recent works [19, 28, 55, 67, 69, 88, 98, 120, 121, 132], is not biased towards fixed arbitrary levels of similarity or coordination, but instead erases network structures that convey limited information, allowing to focus on meaningful expressions of coordination.

Table 3. Types and characteristics of the time windows and the coordination networks used by network science methods.

reference	time window		network	
	type	size	type	layer(s)
[136–138]	adjacent	15 min, 1 hour, 6 hour, 1 day	user	single
[130]	adjacent	1 day, 1 week	user	single
[18, 67]	adjacent	1 week	user	single
[50]	evenly distributed overlapping	1 sec	user	single
[17]	evenly distributed overlapping	from 1 sec to 250 sec	content	single
[62, 108, 109]	evenly distributed overlapping	1 min	user	single
[90]	evenly distributed overlapping	from 1 min to 30 min	user	single
[1]	evenly distributed overlapping	5 min	user	single
[73, 91]	evenly distributed overlapping	5 min	user	multiple (L=2)
[51, 74, 89]	evenly distributed overlapping	5 min	user	multiple (L=3)
[101]	evenly distributed overlapping	30 min	user	single
[121]	evenly distributed overlapping	1 week	user	single
[65]	action-driven overlapping	from 1 sec to 1 hour	user	single
[26]	action-driven overlapping	from 1 sec to 11 day	user	single
[100]	action-driven overlapping	10 sec	user	single
[27, 72]	action-driven overlapping	10 sec	user	single, multiple (L=4)
[44–47, 52, 105, 115]	action-driven overlapping	from 10 sec to 1 min	user	single
[11]	action-driven overlapping	25 sec	user	single
[140]	action-driven overlapping	1 min	user	single
[31]	action-driven overlapping	1 min, 1 hour, 1 day	user	multiple (L=4)
[133]	action-driven overlapping	10 tweets	content	single

Filtering methods based on the timings of the co-actions use a sequence of time windows of equal size $Z = t_{\text{end}} - t_{\text{start}}$. The filtering typically occurs when computing user similarities, by only retaining the co-actions that occur within the same time window. In other words, this filter corresponds to adding a further constraint to the actions that determine the coordination, in that such actions must be temporally close to one another. Table 3 displays the type and length of the time windows used in literature, along with the type of coordination network built. As shown in Tables 3 and 4, time windows can be either *adjacent* [1, 18, 67, 130, 137, 138] or *overlapping*. Furthermore, overlapping time windows can be *evenly distributed* in time [1, 17, 50, 51, 62, 73, 74, 89–91, 101, 108, 109, 121], or *action-driven*—that is, positioned based on the timings of the actions [11, 26, 27, 31, 44, 44–47, 52, 65, 72, 100, 105, 115, 133, 140]. As shown in Table 4, in the latter case each time window starts when a relevant action occurs. In the former case instead, the additional parameter step $\delta < Z$ defines the temporal offset between two consecutive time windows. The amount of overlap O between two consecutive overlapping time windows is thus $O = Z - \delta$. Given a sequence of actions by two or more users, the choice of time windows and their parameters influence the number of such actions that are valid co-actions. These are indicated with the 🗄️ lock icon in the example shown in the rightmost column of Table 4. Let $d = t_2 - t_1$ be the delay with which user u_2 performs an action at time t_2 , with respect to the same action that u_1 performed at time t_1 . Independently of the type of time window, actions whose $d > Z$ are never considered as co-actions. Then, actions with $d \leq Z$ are always co-actions when using action-driven overlapping time windows. Conversely, actions with $d \leq O$ are always co-actions when using evenly distributed overlapping time windows, and can possibly be co-actions when $O < d \leq Z$. Instead, with adjacent time windows, also actions that are close in time (i.e., $d \ll Z$) are occasionally not considered valid co-actions, when they occur across the boundary between two windows, as in the case of the actions occurred at t_3 and t_4 in the example of Table 4. Therefore, using overlapping rather than adjacent time windows ensures that no actions close in time are missed. However, it increases the number of windows required to cover the same time frame, leading to more computational demands. Another critical consideration is the length Z of the time window. A larger Z includes more actions as co-actions, while a shorter one restricts valid co-actions to highly synchronized events, while also increasing the number of windows and the resulting computations. The choice of Z also depends on

Table 4. Time window types, their parameters, and their effect on valid co-actions. To the right, a sequence of actions occurs at times t_1, \dots, t_6 . The same sequence results in different valid co-actions, marked by the lock icon, depending on the time window type.

type	parameters	action sequence and valid co-actions
adjacent		
evenly distributed overlapping		
action driven overlapping		

the goal of the analysis. Literature indicates that inauthentic or harmful coordinated behaviors are highly synchronized and can often be detected with short windows [74, 101], whereas emergent human behaviors that are typically less orchestrated and loosely synchronized, require longer windows to capture their relaxed temporal dynamics [98].

A few works departed from the traditional time windows filters presented above. For example, [121, 137, 138] built a distinct coordination network for each time window. Then, [137, 138] aggregated all networks, each corresponding to a different adjacent time window, by computing the pairwise sums of the network weights. The summation includes a temporal decay weighting scheme that emphasizes the contribution of recent time windows over older ones. Instead, Tardelli et al. [121] built a multiplex temporal network where the layers are obtained from a sequence of evenly distributed overlapping time windows. Differently from all other approaches, the multiplex temporal network is then analyzed as a whole. A further innovation is introduced in [90], which solves an optimization task to select the best window size Z in an overlapping time windows setting. Finally, we remark that the filtering methods discussed in this section can be, and oftentimes are, used in combination for greater efficiency and to further reduce the network size when analyzing very large datasets.

4.1.4 Community discovery. This step aims to detect a set of coordinated communities $P = \{P_1, \dots, P_n\}$ via community discovery on the coordination network. Multiplex networks are either flattened before performing community discovery [56, 137, 138] or an algorithm suitable for multiplex networks is used [75, 121].

Implementation. Most of the works dealing with single layer networks carry out community discovery with LOUVAIN [14, 17, 25, 50, 62, 65, 67, 87, 88, 90, 91, 93, 95, 132], or more broadly via modularity clustering [18]. Other recent works [1, 92, 94, 121, 130] relied on LEIDEN [123]. Among the advantages of these approaches is their scalability, which makes them suitable for the analysis of large networks. In addition, modularity clustering provides a hierarchical community structure, allowing for analyses at different levels of granularity. The authors of [98, 120] combined community discovery via LOUVAIN with the iterative application of a progressively increasing edge weight filtering threshold on the coordination network. As a result of the moving threshold, they studied how the structure and the properties of the coordinated communities change across the whole spectrum of coordination. Moreover, the moving threshold implicitly defines a measure for the extent of coordination observed at each iteration, for each coordinated community, thus providing a continuous score of coordination rather than a binary label. Others proposed a variation of focal structures analysis [110] to identify influential sets of nodes in a coordination network [137, 138]. An alternative to community discovery is the application of a particularly restrictive set of filters that results in splitting the coordination network in

Table 5. Network science methods for detecting coordinated behavior based on *multiplex* user networks, where each layer corresponds to a different co-action. For each group of works we report the considered co-actions, similarity functions, filtering criteria, and the optional flattening step applied before the community detection method.

reference	action	similarity	filters [†]	flattening	community detection
[31]	share, message, URL, hashtag	cardinality	threshold, ADO		Louvain, IPVC [‡]
[27, 72]	retweet, tweet, URL, hashtag	cosine similarity TF-IDF	threshold, ADO	unweighted edge union	
[51]	retweet, text, URL, reply	cardinality	EDO	multigraph	connected components
[73]	URL, hashtag	cardinality	EDO		multi-view clustering
[74]	URL, hashtag, mention	cardinality	EDO		multi-view clustering
[89]	URL, hashtag, mention	cardinality	EDO	sum cardinality	
[91]	URL, tweet	cardinality, cosine similarity	threshold, EDO	unweighted edge union	Louvain

[†] EDO: evenly distributed overlapping time window, ADO: action-driven overlapping time window. [‡] IPVC: iterative probabilistic voting consensus

multiple connected components [11, 13, 44–47, 51, 52, 105, 115, 140]. Since each component is disconnected from the others, this process essentially produces an output that is similar to that of community discovery, in that it identifies groups of connected and highly coordinated users. However, the application of very restrictive filters also discards much of the coordination network [98].

The works presented in Table 5 performed community discovery on a multiplex coordination network. Some authors reduced the complexity of dealing with multiplex networks by flattening them into single layer networks where nodes and edges are the union of the nodes and edges of each layer. Flattened networks can be either weighted [89] or unweighted [27, 72, 91], with the latter resulting in the loss of some information. Alternative approaches involve representing the multiplex network as a single-layer multigraph with labeled parallel edges [51], leveraging multi-view clustering to combine different network layers [73, 74], or applying LOUVAIN on each layer and an iterative probabilistic voting consensus algorithm to achieve consensus clustering [31].

4.1.5 Content networks. A few works built and analyzed networks of *content* rather than *users*. After selecting the initial set of users, these methods build a fully connected network where the nodes are contents posted by the users (e.g., posts) and edge weights are proportional to the similarity between the linked contents. The community discovery process on a content network yields clusters of highly similar contents, regarded as proxies for coordination. User and content networks are interchangeable because each node in a user network can be mapped to the content they published, and each node in a content network can be mapped to its respective author, allowing for a dual representation of the same underlying actions. The main characteristics of the works based on content networks are reported in Table 6. Some works built and analyzed content networks where the nodes were Twitter cashtags [22, 23] or hashtags [25, 138]. Edge weights were given by the number of co-occurrences of the cashtags and hashtags in the same tweets. The analyses revealed suspicious clusters of contents that were later linked to online manipulations by coordinated actors. Lee et al. [66] created a network of similar texts published by multiple users, identifying clusters through connected components and maximal clique analysis. Similarly, [133] constructed two content networks: one for cospasta tweets, clustering coordinated actors via hierarchical clustering, and another to analyze the temporal evolution of co-occurring hashtags. Finally, [95] built an image similarity network by computing image embeddings and using Euclidean distance for comparison, then applied LOUVAIN to detect groups of similar images.

4.2 Data mining and machine learning methods

Methods in this category follow the typical knowledge discovery in databases (KDD) analytical process, involving data collection and preparation, application of data mining and machine learning techniques, validation of the discovered

Table 6. Network science methods for detecting coordinated behavior based on *content* networks, where nodes are posted contents and edge weights encode the similarity between the linked contents. For each group of works we report the types of content, similarity functions, filtering criteria, and the community detection methods.

reference	nodes	similarity	filters [†]	community detection
[66]	texts	text similarity scores	threshold	loose strict campaign, cohesive campaign
[133]	texts, hashtags	cosine similarity, cardinality	EDO, threshold	hierarchical clustering
[25, 138]	hashtags	cardinality	threshold	Louvain
[22, 23]	cashtags	cardinality	threshold	
[95]	images	Euclidean distance	kNN graph	Louvain

[†] EDO: evenly distributed overlapping time window

patterns, and extraction of insights by interpretation of the results. The types of analyzed data include texts [3, 4, 7, 103, 107, 116], images [116], audio [83], interactions [112, 118, 142, 143], and temporal data [8, 42, 61]. In terms of machine learning techniques, some works leverage unsupervised approaches and focus on identifying groups of users exhibiting similar behaviors, while others rely on the availability of labeled data and apply supervised techniques to classify each user as either coordinated or non-coordinated. A comprehensive overview of these methods is presented in Table 7.

4.2.1 Unsupervised. Unsupervised methods work with unlabelled data and, therefore, do not exploit any knowledge on the membership of users to coordinated groups. Multiple unsupervised methods [3, 4, 103, 116] are based on the TEXTCLUST [16] stream clustering algorithm. TEXTCLUST is designed to group similar textual documents into micro-clusters that represent the topics recently discussed in the stream. The detection of rapidly growing clusters of documents is used to identify inorganic or orchestrated campaigns by coordinated users. This approach is conceptually similar to the analysis of a content network of similar documents that includes a time-based filter. A similar approach is also used in [7], where text and node embeddings are clustered via DBSCAN. Keller et al. [61] modeled user daily tweeting activity with a binary matrix whose cells $x_{jt} = 1$ represent a user u_j who tweeted at least once on day t , while cells $x_{jt} = 0$ indicate no tweeting on that day. Then, groups of users with similar tweeting behaviors are found via expectation-maximization. Here, highly similar tweeting behaviors are considered as a proxy for coordination. Others modeled the sequence of user activities as temporal point processes. Sharma et al. [112] proposed the AMDN-HAGE generative model to jointly account for user activities and hidden group behaviours. The ATTENTIVE MIXTURE DENSITY NETWORK (AMDN) component models observed activity traces as a temporal point process, while the HIDDEN ACCOUNT GROUP ESTIMATION (HAGE) component models user groups as mixtures of multiple distributions. The synchronized groups of accounts detected by AMDN-HAGE are deemed coordinated, and the work makes the assumption that all detected coordinated groups are malicious. Similarly, Zhang et al. [142] jointly learned a distribution of user-group assignments based on how consistent each assignment is to the user embedding space and to some prior knowledge such as temporal logic. Finally, they used expectation-maximization to cluster the users according to the different distributions.

In addition to the previous works, others proposed simpler methods or indicators as signals of possible coordinated behavior. Bellutta and Carley [8] analyzed account creation times, given that accounts involved in coordinated malicious activities are typically created in bursts [119]. They computed the daily histogram of account creations to which they applied a burst detection algorithm, identifying spikes in account creations by comparing the number of accounts created in a given day against the average number of daily accounts created in a reference time window. Another simple unsupervised technique is proposed in [118], where coordination is measured as the extent to which users converge—spontaneously or in an organized fashion—on the use of certain hashtags. The Gini coefficient, an indicator

Table 7. Data mining and machine learning methods for detecting coordinated behavior. For each group of works we report the input types, the machine learning approach, the learning paradigm, and whether the method takes time into account.

reference	input	machine learning approach	learning	time
[3, 4, 103]	text streams	text stream clustering	unsupervised	✓
[116]	text streams, image captions	text stream clustering	unsupervised	✓
[7]	text, user-content network	clustering of text and node embeddings	unsupervised	✓
[61]	daily tweeting activity	expectation-maximization	unsupervised	✓
[118]	hashtags	peak detection	unsupervised	✓
[8]	account creation timestamps	burst detection	unsupervised	✓
[80]	user activities	contrast pattern mining	unsupervised	✓
[81]	user activities	convergent cross mapping	unsupervised	✓
[32, 33]	URL, hashtag, image, mention	networked Markov chains	unsupervised	✓
[112]	user activities	temporal point processes, gaussian mixture models	unsupervised	✓
[142]	user activities	temporal point processes, expectation-maximization	unsupervised	✓
[143]	user activities	representation learning, conditional embedding, neural encoding	supervised	✓
[42]	network, temporal, semantic features	outlier detection	supervised	✓
[83]	metadata, audio transcripts, thumbnails	ensemble classification	supervised	
[107]	text	peak detection, multiclass classification	supervised	✓

of inequality, is applied to the distribution of used hashtags to create time series where saddles close to 0—indicative of low inequality—correspond to no coordination whereas peaks close to 1—indicative of a situation where all users use the same few hashtags—correspond to strong coordination. Two alternative techniques are proposed in [80] and [81]. The former leverages contrast pattern mining to extract anomalous behavior, while the latter uses convergent cross mapping to discover cause-and-effect relationship and to construct an influence network. Similarly, [32, 33] use a discrete-time stochastic model to analyze coordinated activity, representing the user behaviors as interacting Markov chains.

4.2.2 Supervised. Supervised methods require labeled input data that they use to learn a function for predicting whether users belong to a coordinated group. Zhang et al. [143] tackles a binary classification task to distinguish between coordinated and non-coordinated users that participate in cross-platform campaigns. They leverage information from an *aid platform* where coordinated users are known, to detect unknown coordinated users on a *target platform*. Input data consists of a known coordinated activity set on the aid platform, plus unknown user activities on the target platform. The relationship between the two activities are modeled with neural time series encoders before being fed to a multi-layer perceptron for prediction. Mariconti et al. [83] developed a system to identify YouTube videos that are likely to be raided (i.e., targeted by coordinated hate attacks). They used information coming from metadata, audio transcripts, and video thumbnails to compute machine learning features for each video. The prediction is given by an ensemble classifier that determines the likelihood that a new video will be raided based on a ground-truth of previously raided videos. A subsequent study focused on attributing coordinated hate attacks to the communities that organized them [107]. The system uses a peak detector to identify an abrupt rise in the comment activity of a YouTube video, a signal of a coordinated attack. Then, it leverages a trained classifier based on linguistic features from the comments to the video and a set of Reddit and 4chan communities, identifying the community responsible for each attack based on linguistic patterns and similarities. A simpler approach is proposed in [42] where a small ground-truth of 20 organic and inorganic campaigns is built by some subject-matter experts. Each campaign—the ground-truth ones and others to analyze—is characterized with indicators (i.e., features) belonging to three dimensions of participant behavior: network, time, and semantics. Then, [42] automatically labels as coordinated and inorganic those unknown campaigns whose

indicators lay within the 95% confidence interval of ground-truth inorganic campaigns and outside the 95% confidence interval of ground-truth organic or non-coordinated campaigns.

4.3 Discussion

4.3.1 Modeling complex coordination. In Section 2 we highlighted that coordinated online behavior is complex and multifaceted, often involving various actions across multiple platforms. Consequently, methods that analyze only a single type of action—particularly single-layer network approaches—risk missing significant coordination activities. To address this limitation, using *multiplex* networks and *compound* actions can offer more comprehensive insights. The network science framework presented in Section 4.1 supports both single- and multi-layer (i.e., multiplex) networks. However, few studies considered multiple co-actions and built multiplex networks, as summarized in Table 5. Moreover, to fully leverage the benefits of multiplex networks, these should not be flattened before running the community detection algorithm, which must be specifically designed for multiplex networks so as to utilize the multiple layers effectively [75]. Unfortunately, only a few studies possess these characteristics [27, 31, 72, 73, 89], making the analysis of coordinated behavior across multiple actions a largely unexplored area of research. Compound co-actions, combining simpler actions (e.g., a post with both a hashtag and user mention), can enhance detection confidence because they are less likely to occur by chance [74]. However, this approach is also almost completely unexplored. In conclusion, too little research has been conducted so far on using multiple co-actions for detecting coordinated online behavior, leaving the actual advantages of these methods over single-action analysis unclear. Moreover, any potential benefits must be weighed against the increased complexity and computational costs that they introduce.

4.3.2 Network science vs. machine learning. Most methods for detecting coordinated behavior rely on network science, offering greater generality and expressiveness than machine learning approaches. They effectively model complex interactions and relationships, capturing varying degrees of coordination [98], tracking temporal changes [121], and providing detailed characterizations of coordinated groups [120]. The disadvantage is the computational cost of analyzing large networks and the requirement for human analysts to interpret results, making them less suitable for being directly used in automatic decision making systems. A further drawback is their sensitivity to specific parameters, such as those defining the way in which the timings of user actions are modeled [136], which can significantly impact the obtained results and for which few guidelines currently exist [90, 121]. Instead, data mining and machine learning methods often rely on oversimplified assumptions, such as the existence of a sharp binary distinction between coordinated and non-coordinated actors, which overlooks the nuanced nature of online coordination. This binary approach can undermine the theoretical and practical reliability of results. Furthermore, these methods face limitations due to the lack of comprehensive labeled datasets needed for training effective models. Despite these challenges, machine learning methods excel in encoding diverse actions and characteristics of users, and offer outputs that are directly applicable for automatic decision-making systems, unlike the more complex network science methods. In conclusion, network science methods are particularly suitable when the goal is an in depth understanding of the studied phenomenon, while machine learning methods can be used—with due caution—for quick decisions and when scalability is a concern.

5 CHARACTERIZATION OF COORDINATED BEHAVIOR

When the coordination detection method is not integrated into an automated decision-making system, its output may initiate the characterization task. The aim of this task is to describe each coordinated user, group, or community along one or more of the defining dimensions outlined in Section 2.5. Characterizing the detected instances of coordination

Table 8. Indicators used for characterizing coordinated behavior. For each group of indicators we report the works that used them, the high-level concept implemented, and the defining dimensions of coordination for which the indicators provide ●, or could provide ○, information. Table rows are grouped based on the type of information (i.e., user, content, network) leveraged by the indicator.

	reference	concept	indicators	defining dimensions				
				auth.	harm.	orch.	time	other
user	[4, 8, 17, 25, 26, 50, 55, 87, 89–91, 98, 100, 101, 120, 138]	automation	bot scores	●			○	
	[26, 52, 55, 66, 67, 87, 90, 98, 108, 112, 120, 140, 142]	moderation	suspended users		●		○	
	[90]	username diversity	entropy	●			○	
	[108]	activity	account creation burstiness	●			●	
content	[13, 23, 25, 44, 59, 61, 62, 66, 105, 108, 109, 118, 130, 132, 134, 138]	activity	number of posts, retweets, ...				●	●
	[105, 132]	engagement	number of views, likes, ...				●	●
	[11, 17, 59, 100]	timings	action time interval	○		○	●	●
	[8, 13, 23, 25, 27, 61, 81, 87, 88, 90, 93, 98, 107, 112, 132, 134, 137, 138, 142, 142]	socio-linguistics	attitudes, concerns, emotions, stances, topics, words		●		○	●
	[59, 88, 137, 138]	repetitiveness	text similarity scores		●	○	○	●
	[121, 133]	socio-linguistics	topic and hashtag evolution		●		●	●
	[14, 47]	news diversity	entropy, Gini coefficient			○	○	●
	[8, 46, 134]	news reliability	suspended and blacklisted URLs		●		○	
	[7, 8, 105, 120]	news reliability	NewsGuard and MBFC [†] scores		●		○	
	[55]	manipulation	propaganda scores		●		○	
	[69, 107]	offensiveness	toxicity scores		●		○	
	[65, 69, 98, 120, 121]	political bias	MBFC [†] scores, hashtag bias				●	●
network	[17, 18, 55, 65, 69, 87–90, 98, 120]	coordination	edge weights			●	○	●
	[19, 26, 28, 47, 51, 63, 65, 87, 88, 91, 100, 132]	influence	centrality scores			●	○	
	[18, 23, 87, 98]	homophily	assortativity			●	○	
	[18, 26, 47, 51, 63, 65, 87, 88, 91, 95, 98]	cohesiveness	clustering coefficient, modularity, density			●	○	
	[67, 121]	stationarity	user influx and outflux, user evolution				●	

†: Media Bias/Fact Check

can also provide valuable information for validating detection results, especially in the absence of ground truth data. Characterizing coordinated online behavior is a semi-automatic task, as it involves some degree of manual analyses and observations by human analysts supported by some automatically-computed indicators that provide information on the coordinated actors. Table 8 summarizes the main indicators used for characterizing coordinated actors, whose discussion is presented in the remainder of this section. Indicators in table are grouped based on the information they leverage, which can be related to (i) users, (ii) content, or (iii) networks. The table also highlights the dimensions of coordination that each indicator addresses. Full dots indicate dimensions where the indicator has already been applied, while empty dots denote dimensions where it has potential use that has not yet been explored.

5.1 Authenticity

Authenticity refers to the extent to which users, groups, or communities correctly represent themselves to others on a platform. As a consequence, inauthenticity is found when an actor misrepresents itself, such as to mislead others on who they are or what they do. In other words, authenticity is a property of the coordinated actors. In literature, the most common proxy for inauthenticity is an account’s degree of *automation*, obtained via a *bot score* [4, 8, 17, 25, 26, 50, 55, 87, 89–91, 98, 100, 101, 120, 138]. By definition, social bots are accounts that make use of some or full automation [21]. Thus, accounts with a high bot score are likely to be mostly automated. Unfortunately however, while useful, the use of bot scores as indicators of inauthenticity faces some limitations. First, computing bot scores is a challenging task

per se, and it is prone to errors [24]. Second, not all bots are inauthentic, as there exist some types of self-declared bots that operate for neutral or benign purposes [21]. Finally and most importantly, there exist multiple types of users that are inauthentic but that do not make use of automation (e.g., trolls, dissidents). Therefore, using bot scores as indicators of inauthenticity can cause both false positive and false negative errors. In addition to bot scores, some also used username similarity [90] and bursts of account creations [108] as proxies for inauthenticity. Table 8 also shows that the few existing indicators of inauthenticity are all based on user information. While this is expected since authenticity is a property of the actors rather than their actions, we also note that extremely short time intervals between user actions [11, 17, 59, 100] could be used as a red flag of automation and, by extension, also of inauthenticity.

5.2 Harmfulness

Harmfulness measures the extent to which the actions of the coordinated actors can potentially cause negative consequences. Therefore, the analysis of the actions can provide valuable information about the intent, or potential, for harm of the coordinated actors. The production or re-sharing of content (e.g., text, images, links) is the most information-rich type of action and, in fact, nearly all indicators of harmfulness in Table 8 are derived from the analysis of content. A first line of work involves the use of traditional NLP methods to analyze attitudes, concerns, emotions, stances, topics, hashtags, and words, as these provide useful insights into the intent and aims of coordinated actors [8, 13, 23, 25, 27, 61, 81, 87, 88, 90, 93, 98, 107, 112, 121, 132–134, 137, 138, 142, 142]. For instance, certain words, negative emotions, and stances can amplify social discord and polarize public opinion. The prevalence of these indicators is largely due to the abundance of readily available tools and their ease of application. However, these indicators are often quite generic and consequently lack substantial power and informativeness. Text similarity scores were also used to identify harmful campaigns that aim to create false consensus through repeated sharing of similar posts [59, 88, 137, 138]. Other straightforward indicators of harmfulness are those based on the presence of toxic (i.e., hateful or offensive) [69, 107] or propagandistic [55] content. A different line of work measured harmfulness in terms of the unreliability of the news shared by coordinated actors, which was measured via NewsGuard’s⁴ reliability ratings and Media Bias/Fact Check⁵ factual reporting and credibility ratings [7, 8, 105, 120]. Others measured unreliability in terms of suspended and blacklisted URLs [8, 46, 134]. Finally, the fraction of coordinated users that were suspended or banned by a social platform has been extensively used as an indicator of harmfulness [26, 52, 55, 66, 67, 87, 90, 98, 108, 112, 120, 140, 142]. This stands out as the only indicator of harmfulness based on user, rather than content, information. However, we note that harmfulness indicators based on content moderation lack predictive capabilities and can only be used retrospectively, since moderation actions typically occur some time after the account have engaged in the harmful activities [125].

5.3 Orchestration

Orchestration expresses the extent to which the coordination is well-organized, rather than spontaneous and emergent. Since orchestration is related to the level of organization of a group of coordinated actors, all orchestration indicators are based on the analysis of coordination networks, as shown in Table 8. As such, these indicators are typically computed after the application of a network science detection method. A large coordination score among a group of users was used as a proxy for orchestration, based on the idea that tightly synchronized groups are likely to be well-organized [17, 18, 55, 65, 69, 87–90, 98, 120]. Centrality scores measure the importance of nodes within a network or community. Hence, nodes with high centrality wield significant influence over the dissemination of information.

⁴<https://www.newsguardtech.com/> (accessed: 31/07/2024)

⁵<https://mediabiasfactcheck.com/> (accessed: 31/07/2024)

As such, the presence of nodes with high centrality scores in a coordinated community could indicate a structured hierarchy [19, 26, 28, 47, 51, 63, 65, 87, 88, 91, 100, 132]. Conversely, assortativity is a measure of homophily in a network. Measuring assortativity can help determine whether coordinated behavior is orchestrated or spontaneous by assessing the extent to which highly connected nodes tend to connect with other highly connected nodes [18, 23, 87, 98]. For example, high assortativity between nodes with large degree might indicate a centralized orchestrated structure. Instead, high assortativity between nodes with low degree might indicate decentralized orchestration, while low assortativity or disassortativity might indicate spontaneous behaviors. Likewise, low modularity, high density, or high clustering coefficient indicates a tendency for nodes to form tightly-knit groups, suggesting well-organized and potentially orchestrated behavior [18, 26, 47, 51, 63, 65, 87, 88, 91, 95, 98]. Finally, the time intervals between the actions of coordinated actors can provide information on whether the coordinated behavior is orchestrated or spontaneous because tightly synchronized actions suggest a higher level of premeditated organization, whereas more irregular intervals might indicate spontaneous, less structured coordination. The study of the internal structure and organization of coordinated communities is also relevant for investigating the diverse strategies of online manipulation. Organized influence campaigns often unfold in distinct ways: some aim to mobilize pre-existing organic communities, leveraging existing social ties and trust, while others seek to establish new communities controlled by a central entity [117]. Additionally, organic communities can emerge from grassroots efforts, where users are genuinely motivated by a common goal without external orchestration. These different strategies leave “network footprints” that can be captured by orchestration indicators [18].

5.4 Time-variance

Time-variance indicators aim at quantifying temporal changes in a wide array of characteristics of the coordinated actors. These include the number, intent, and behavior of the actors, which may result in changes in the types, timing, frequency, and intensity of their actions. Many measures can serve as time-variance indicators, including the indicators previously discussed for the dimensions of authenticity, harmfulness, and orchestration. Indeed, as shown in Table 8, all indicators for authenticity, harmfulness, and orchestration can potentially be analyzed over time, even if they have not been already studied in this way yet. However, utilizing a measure as a time-variance indicator requires repeated and extended monitoring to detect possible changes over time. This requirement for prolonged monitoring poses an additional challenge compared to other types of indicators, which explains the scarcity of studies that employed time-variance indicators or that investigated the temporal dynamics of coordinated behavior [121]. The few existing detailed temporal analyses examined user flows between coordinated communities or the variation of users between different time windows [67, 121]. Other temporal analyses of user activity can involve examining bursts of account creations [108], abnormal post or retweet volumes [13, 23, 25, 44, 59, 61, 62, 66, 105, 108, 109, 118, 130, 132, 134, 138], or inflated engagement metrics, such as the number of views [105, 132]. The analysis of topics and hashtags over time can reveal the evolution of the narratives promoted or discussed by the coordinated actors [121, 133]. The timing of the actions, such as the intervals between pairs of actions that result in a co-action, can provide insights into the type of coordination and strategies employed, which may evolve over time [11, 17, 100].

5.5 Other general indicators

In addition to the indicators that provide information about the four defining dimensions of coordinated behavior, other general-purpose indicators have also been used. These come in handy to provide additional context on the coordinated behavior and may be particularly relevant depending on the application context, as in the case of the

indicators of political polarization or bias that are used to characterize coordinated communities involved in online electoral debates [65, 69, 98, 120, 121]. Political bias was computed based on MBFC scores or by leveraging hashtag polarity. Analyses of the level of activity of coordinated actors on a platform [13, 23, 25, 44, 59, 61, 62, 66, 105, 108, 109, 118, 130, 132, 134, 138], of the engagement they obtain [105, 132], or of the content they produce, can provide additional contextual information. Regarding the latter, standard socio-linguistic analyses were used to draw insights into the narratives and stances of coordinated actors [8, 13, 23, 25, 61, 87, 88, 90, 93, 98, 107, 112, 121, 132–134, 137, 138, 142, 142]. Similarly, others considered the diversity of the news shared by the coordinated actors [14, 47] or assessed the presence of patterns in the timings of their actions [11, 17, 59, 100]. Finally, acknowledging that coordination is a nuanced and non-binary concept, some studies computed coordination scores, for example based on the edge weights of the coordination network [17, 18, 55, 65, 69, 87–90, 98, 120]. Analyzing the degree of coordination can contribute to shedding light on the extent of collaboration and influence within the group, and on the effectiveness of their collective actions.

5.6 Compound indicators

The previous discussion involved indicators used individually. However, some works also investigated the combined use of multiple indicators to derive even deeper insights, possibly revealing complex interactions and enhancing the overall understanding of coordinated behavior. For example, authors of [18, 55, 87, 98, 120] measured trends and correlations between indicators of coordination, propaganda, moderation, cohesiveness, and automation, finding that certain coordinated communities exhibit markedly different behaviors. They also found that assessments of the nature of online coordination—such as determining if it is harmful or inauthentic—can be unreliable when based on a single indicator, highlighting the importance of combining multiple indicators for more accurate results [55]. Similarly, others considered the interplay between coordination, toxicity, and political bias [69]. Finally, Tardelli et al. [120] proposed a comprehensive approach to characterize coordinated communities by analyzing multiple indicators. They examined coordination, automation, suspensions, news unreliability, and political bias, visualizing each community’s scores on these dimensions with a radar chart. They observed that larger scores across these dimensions indicate suspicious behavior. Therefore, communities with wider radar charts were deemed more suspicious than those with smaller areas. The area of a radar chart thus serves as a compound indicator, derived from combining the underlying dimensions. The suspiciousness of a coordinated community can thus be quantitatively assessed by calculating the radar chart’s area, or any other measure that correlates positively with the underlying dimensions. This approach exemplifies how multiple individual indicators can be combined to derive more detailed insights. Considering that there are only a few instances of the use of compound indicators in literature, this represents a largely unexplored area of research.

5.7 Discussion

Table 8 shows that many and diverse indicators are used to measure harmfulness. These include indicators of hate speech, toxicity, and propaganda, as well as indicators based on platform moderation decisions. Conversely, inauthenticity has only been examined through the lens of automation, despite being the focus of many works. From a practical standpoint, understanding the authenticity of an account is inherently more challenging than assessing its harmfulness, as authenticity pertains to the nature of the actors rather than their actions. Many online platforms afford users a significant degree of anonymity, making it difficult to verify their true identities. In contrast, many actions performed online leave unforgeable traces, which can be more easily analyzed to determine harmful behavior. Despite the challenges however, future work should focus on developing additional indicators of authenticity. Among the untapped sources of information are the completeness and credibility of the user profile, the regularity of activity patterns, the duplication or

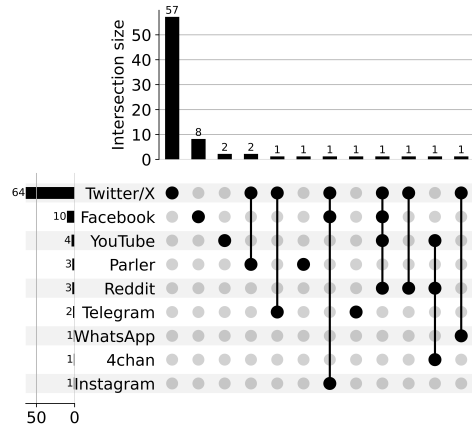


Fig. 7. Distribution of works based on the combination of analyzed platforms. The vast majority of works analyzed a single platform, mainly Twitter/X. Only a few works performed cross-platform analyses, represented by multiple connected dots.

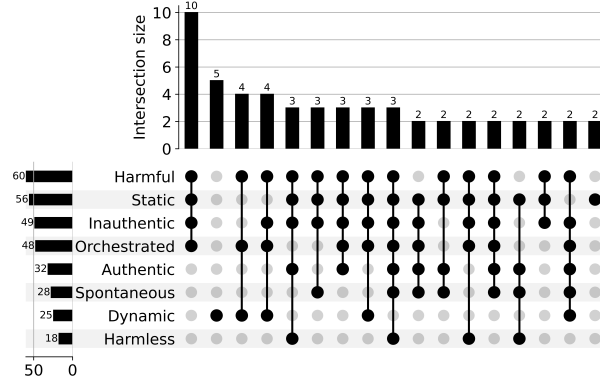


Fig. 8. Distribution of works focused on different types of coordinated behavior, categorized based on the defining dimensions of authenticity, harmfulness, orchestration, and time-variance. Vertical bars show the number of works for each combination of dimensions, while horizontal bars show the number of works for each individual dimension. Only combinations covered by 2+ works are shown.

repetitiveness of posted content, the lack of consistency in language or interests, the authenticity of posted multimedia content, the cross-platform consistency of profile and behavior, the similarity of profile and behavior to that of previously moderated accounts, and the verification status, which however should be evaluated differently depending on whether such status can be purchased on the analyzed platform.⁶

6 OPEN CHALLENGES AND FUTURE RESEARCH DIRECTIONS

Multiplatform-ness. Some coordinated campaigns unfold simultaneously on multiple platforms resulting in *multiplatform* coordinated behavior [139]. These campaigns exhibit overall similar behaviors, such as the use of campaign-specific hashtags, with minor variations between platforms. In contrast, *cross-platform* campaigns involve different platforms playing distinct roles, with coordinated behavior varying significantly across the involved platforms. For example, organizers of a targeted hate attack might use private groups on a messaging platform to plan their actions before moving to the target platform to flood it with hateful comments [59]. Detecting coordination on the planning platform might involve analyzing group joinings, while on the target platform it could involve examining synchronized commenting. The example highlights the challenges of multi- and cross-platform analyses, which require extensive data collection, careful selection of actions or machine learning features for each platform, and identification of the same groups of users across multiple platforms. While the latter can be mitigated via the analysis of *content* rather than *user* networks—as discussed in Section 4.1.5—the remaining outstanding challenges result in a scarcity of multi- and cross-platform analyses. Figure 7 shows that most studies analyzed coordinated behavior on Twitter/X, with only a few multiplatform analyses [7, 31, 59, 93, 94, 105, 107, 143], highlighting that this area still necessitates much investigation.

Temporal variability. Most studies on coordinated behavior incorporate time to some degree, as shown in Tables 3, 7, and 8. However, the majority of such works only perform superficial temporal analyses. Figure 8 highlights that only a few works have deeply investigated the *temporal dynamics* of coordinated behavior, using methods like multiplex temporal

⁶link.gale.com/apps/doc/A746844400/AONE (accessed: 31/07/2024)

networks [121] and temporal point processes [112, 142]. Thus, analyzing temporal dynamics remains a promising yet relatively unexplored research avenue. Temporal analyses offer numerous advantages, including examining (i) the temporal stability of coordinated groups, (ii) changes in membership, structure, and actions, and (iii) adaptation to countermeasures or other stimuli [48]. Further, dynamic analyses seem to yield more accurate results compared to static ones [121]. However, this approach faces several challenges, such as setting appropriate temporal parameters and addressing the computational demands resulting from the time-intensive nature of dynamic analyses.

Heterogeneity. The great degree of heterogeneity inherent of coordinated behavior does not only manifests through temporal variations, but also through the multitude of characteristics that the different instances of coordinated behaviors exhibit. Figure 8 shows the distribution of studies that addressed specific types of online coordination, based on the defining dimensions introduced in Section 2.5. As shown by the vertical bars, coordinated *harmful inauthentic* and *orchestrated* behavior is the type of coordination that was studied the most. As discussed in Section 2, this is due to both practical reasons related to the urgency of contrasting online manipulations, as well as to the large body of work that adopted Facebook’s initial Definition 2.3 of coordinated inauthentic behavior. Almost all the other remaining combinations involve *harmful* behavior and *static* analyses. Figure 8 also shows that *authentic spontaneous* and *harmless* coordinated behavior is completely unstudied. While harmful, inauthentic, and orchestrated behaviors are of particular practical relevance for content moderation purposes, the study of harmless, authentic, and spontaneous coordination should not be overlooked as it provides valuable insights into fundamental aspects of online human interactions. Therefore, the analysis of the less problematic forms of online coordination should be increased in the near future.

Multimodality. Online coordination can involve multiple content modalities, such as text, images, and videos. This diversity of content types adds complexity to detecting and analyzing coordinated behavior, as each modality has unique characteristics requiring different analytical techniques. For instance, text-based coordination might involve analyzing text similarities or hashtags use [26, 132, 134]. Image-based coordination could require image recognition tasks or computing image similarities [45, 101, 140], while video-based coordination demands techniques like video content analysis and transcript examination. As shown in Tables 2, 6, and 7, however, only a handful of works investigated modalities other than text [45, 83, 95, 101, 115, 140]. Accounting for and integrating different modalities offers the potential to capture a more comprehensive picture of coordinated behavior, but also poses new challenges. These include developing features and methodologies to assess similarities based on multimedia content and addressing the computational demands of multimodal analyses. Nonetheless, given the rising trend of online platforms centered around multimedia content, multimodal analyses appear promising and well motivated.

Generative AI. The rise of generative AI presents both challenges and opportunities. One potential challenge is the increased difficulty in detecting coordinated actors that use such techniques to mask their activities. AI is already capable of producing human-like text with given properties, and increasingly also images, audio, and video [96]. Moreover, it has already been used to simulate authentic online human behaviors in such a way as to evade detection by existing systems [53]. Progress in generative AI could thus make it harder to detect future instances of coordinated behavior or to assess the authenticity of coordinated actors. However, it is still unknown what the effect of these techniques will be on the landscape of online coordination. At the same time, generative AI (e.g., LLMs) can be used to simulate human behaviors within agent-based models [122], enabling the creation of realistic simulations of coordinated activities. In future, these models could be leveraged to train or evaluate detection methods, providing a controlled environment to probe methods’ capabilities at detecting different instances and types of coordinated behavior.

Scalability. Large-scale coordinated campaigns often involve tens of thousands of users and millions of interactions, creating a heavy computational burden. This issue is especially pronounced for methods that rely on computing pairwise similarities between users, such as the majority of network science methods. Consequently, some studies analyze only a tiny fraction of users—sometimes as little as 1% [28, 98]—which limits the scope and reliability of coordination detection results. Future efforts should focus on increasing the scalability of current detection methods to handle the vast amounts of data generated by large-scale campaigns. Possible solutions to this issue can include the use of *approximation algorithms* to obtain near-accurate results with reduced computational complexity, or *sampling techniques* to select representative data subsets, making computations more manageable. Additionally, *online algorithms* could be used to update results continuously as new data arrives rather than recomputing from scratch, and *graph summarization* techniques could be used to condense large networks while preserving key structural properties for faster analysis.

Data availability. The scarcity of labeled authoritative data limits the development of new coordination detection methods. This constraint hampers the creation of supervised methods, leading to the predominance of unsupervised approaches that we documented in Section 4. Moreover, it complicates the fine-tuning and evaluation of all methods, now primarily assessed through post-hoc characterizations rather than via quantitative metrics. Unfortunately, the challenge of data availability is expected to intensify in the future, as access to platforms data continues to diminish in the post-API age [60, 124]. Potential solutions include using certain public data as ground truth or running simulations. Examples of the former approach are those works that built balanced ground-truths based on authoritative and publicly available repositories,⁷ obtaining datasets that are suitable for training [34, 72] and evaluation [121] purposes. Existing datasets are however constrained to only one type of coordinated behavior on a single platform, which is insufficient for training detectors capable of generalizing to the many instances of online coordination, such as those sketched in Figure 3. This calls for further endeavors to create more comprehensive ground-truths. An alternative solution is to evaluate detection methods in simulated environments, as done in recent studies that employed agent-based models to mimic coordinated behavior [58, 86, 103]. In any case, ensuring access to reliable benchmarks—either real or synthetic—will remain a critical issue that demand significant attention in the coming years, particularly given the dynamic and complex nature of coordinated online behavior, which necessitate diverse and frequently updated data.

Subjectivity. The existence of many indicators of harmfulness—reported in Table 8—is advantageous from a practical standpoint. However, this multitude of indicators reflects great variability in how harmfulness has been practically defined. More broadly, the subjective nature of harmfulness—discussed in Sections 2.4.3 and 2.5.2—presents several challenges to its assessment. The first of such challenges is the variability with which different stakeholders define what constitutes harmful behavior, which hinder reaching a consensus on definitions and indicators [30]. Additionally, perceptions of harmfulness are influenced by cultural and social norms, which differ across individuals, regions, and communities. Consequently, what is seen as harmful by some might not be perceived the same way by others. Finally, assessing harmfulness involves evaluating both the intent behind the actions and their impact. While the impact can sometimes be measured—as with misinformation that leads to real-world violence [133] or that reduces vaccine uptake [68]—the intent is often hidden or ambiguous, making it difficult to measure accurately. Given that a certain degree of subjectivity in the assessment of online coordination is inevitable and bound to persist, future works must prioritize transparency in defining and implementing it, favoring nuanced analyses based on multiple indicators.

⁷<https://transparency.x.com/en/reports/moderation-research.html> (accessed: 31/07/2024)

Attribution. Detecting coordination, especially in cases where malicious actors aim to remain hidden, is already challenging. Attributing these instances to the responsible entities is even more difficult. Attribution does not merely involve identifying the accounts that take part in the coordination, but most importantly uncovering the entities behind them—such as movements, organizations, or states—and understanding their goals. Several obstacles hinder this process, including the anonymity afforded by online platforms and the sophisticated techniques used to obfuscate identities, activities, and intentions. Potential solutions include strengthening the collaboration with platforms, which have access to more data than what is publicly visible or available via APIs or scraping. Additionally, advanced forensic techniques and multidisciplinary approaches that combine technical, social, and behavioral analysis are essential. Despite these solutions, the task will remain daunting and require continuous research and experimentation.

Multidisciplinarity. The study of coordinated behavior is intrinsically multidisciplinary due to the complex interplay of technical, social, and regulatory factors involved. This multidisciplinary nature is evident in the diversity of the studies on the subject, which employ a wide array of methods from multiple scientific communities. For instance, computer scientists and complex systems scholars develop computational methods to detect instances of coordinated behavior [74, 120, 138]. Social scientists study the emergent behaviors and interactions that lead to massive online coordination [25, 65], while political scientists analyze the impact of coordinated campaigns on real-world events, such as elections [17, 98, 140]. The involvement of media and communication experts in understanding the dissemination of (mis)information further highlights the multidisciplinary scope [47, 115]. In this context, multidisciplinarity is both a challenge and an opportunity. Challenges arise from the need to coordinate efforts among various scientific communities and to organize and keep track of the extensive body of work on the subject—a task for which this survey aims to provide a contribution. In spite of the challenges however, the comprehensive understanding of the phenomenon and the effective mitigation of its nefarious instances can only be achieved through the collaboration between diverse stakeholders, such as scholars, platform operators, policymakers, and civil society organizations.

Ethical dilemmas. Research on coordinated online behavior also faces some ethical dilemmas. The process of collecting and analyzing online data involves monitoring user interactions and behaviors, which can raise concerns about surveillance and consent. The inherent subjectivity of coordinated behavior leads to additional challenges, as certain coordinated actions may be punished on some platforms while being tolerated on others [30], resulting in imbalanced interventions and potentially unfair decisions. Characterizing coordinated actors further complicates matters, as this requires assessing the authenticity and harmfulness of actors who might have strong motivations—such as social and political reasons [2]—for remaining anonymous, raising privacy concerns. Additionally, developing methods to detect coordinated actors presents risks as these tools could be misused to silence certain groups or minorities, thereby selectively threatening freedom of expression. Ultimately, these issues demand that great attention to ethical and normative considerations should be devoted in future research.

7 CONCLUSIONS

Coordination is a fundamental dynamic of human behavior and its study holds great theoretical significance and numerous practical implications. Theoretically, understanding coordinated online behavior contributes to shedding light on social dynamics, group formation, and collective action. Practically, this research is pivotal in combating online manipulations, such as disinformation campaigns and orchestrated hate attacks, and fostering more inclusive online spaces that promote genuine cooperation and constructive discourse. However, the field faces several significant challenges, including the complexity of multimodal and cross-platform analyses, the opportunities and perils of

generative AI, the scarcity of labeled data, the subjectivity and ethical dilemmas inherent to the evaluation of coordinated behavior. These diverse challenges, coupled with the broad interest in the topic, highlight the need for multidisciplinary efforts from diverse research communities. In the coming years, these efforts should aim to develop a comprehensive array of datasets, methods, and studies to better understand and address this complex and dynamic phenomenon.

ACKNOWLEDGMENTS

This work was supported in part by project SERICS (PE00000014) under the NRRP MUR program funded by the EU – NGEU; in part by PNRR-M4C2 “FAIR-Future Artificial Intelligence Research”-Spoke 1 “Human-Centered AI,” funded under the NextGeneration EU Program under Grant PE00000013; “SoBigData.it—Strengthening the Italian RI for Social Mining and Big Data Analytics”—Prot. IR0000013. H2020 Research Infrastructures (Grant No. 654024).

AUTHOR CONTRIBUTIONS

Lorenzo Mannocci, Michele Mazza, Anna Monreale, Maurizio Tesconi and Stefano Cresci designed the research; Lorenzo Mannocci and Michele Mazza performed the review; Lorenzo Mannocci, Michele Mazza, Anna Monreale, Maurizio Tesconi and Stefano Cresci wrote the paper; Anna Monreale, Maurizio Tesconi and Stefano Cresci supervised the work; Maurizio Tesconi provided funding.

REFERENCES

- [1] Iuliia Alieva, Lynnette H. X. Ng, and Kathleen M. Carley. 2022. Investigating the spread of Russian disinformation about biolabs in Ukraine on Twitter using social network analysis. In *IEEE BigData*. 1770–1775.
- [2] Hans Asenbaum. 2018. Cyborg activism: Exploring the reconfigurations of democratic subjectivity in Anonymous. *New Media & Society* 20, 4 (2018).
- [3] Dennis Assenmacher, Lena Adam, Heike Trautmann, and Christian Grimme. 2020. Towards real-time and unsupervised campaign detection in social media. In *AAAI FLAIRS*. 303–307.
- [4] Dennis Assenmacher, Lena Clever, Janina S. Pohl, Heike Trautmann, and Christian Grimme. 2020. A two-phase framework for detecting manipulation campaigns in social media. In *HCII*. 201–214.
- [5] Helmy H. Baligh. 1986. Decision rules and transactions, organizations and markets. *Management Science* 32, 11 (1986), 1480–1491.
- [6] Helmy H. Baligh and Richard M. Burton. 1981. Describing and designing organizational structures and processes. *International Journal of Policy Analysis and Information Systems* 5, 4 (1981), 251–266.
- [7] Fabio Barbero, Sander op den Camp, et al. 2023. Multi-modal embeddings for isolating cross-platform coordinated information campaigns on social media. In *MISDOOM*. 14–28.
- [8] Daniele Bellutta and Kathleen M. Carley. 2023. Investigating coordinated account creation using burst detection and network analysis. *Journal of Big Data* 10, 1 (2023), 1–17.
- [9] W. Lance Bennett and Alexandra Segerberg. 2012. The logic of connective action. *Information, Communication & Society* 15, 5 (2012), 739–768.
- [10] Stephen P. Borgatti. 2005. Centrality and network flow. *Social Networks* 27, 1 (2005), 55–71.
- [11] D. A. Broniatowski. 2021. *Towards statistical foundations for detecting coordinated inauthentic behavior on Facebook*. Technical Report. Institute for Data, Democracy and Politics, The George Washington University.
- [12] Axel Bruns, Tim Highfield, and Jean Burgess. 2013. The Arab Spring and social media audiences: English and Arabic Twitter users and their networks. *American Behavioral Scientist* 57, 7 (2013), 871–898.
- [13] Keith Burghardt, Ashwin Rao, Georgios Chochlakis, Baruah Sabyasachee, Siyi Guo, Zihao He, Andrew Rojecki, Shrikanth Narayanan, and Kristina Lerman. 2024. Socio-linguistic characteristics of coordinated inauthentic accounts. In *AAAI ICWSM*, Vol. 18. 164–176.
- [14] Cheng Cao, James Caverlee, Kyumin Lee, Hancheng Ge, and Jin-Wook Chung. 2015. Organic or organized?: Exploring URL sharing behavior. In *ACM CIKM*. 513–522.
- [15] Valerio Capraro and Matjaž Perc. 2024. In search of the most cooperative network. *Nature Computational Science* 4, 4 (2024), 257–258.
- [16] Matthias Carnein, Dennis Assenmacher, and Heike Trautmann. 2017. Stream clustering of chat messages with applications to Twitch streams. In *ER*. 79–88.
- [17] Victor Chomel, Maziyar Panahi, and David Chavalarias. 2023. Manipulation during the French presidential campaign: Coordinated inauthentic behaviors and astroturfing analysis on text and images. In *CNA*. 121–134.

- [18] Lorenzo Cima, Lorenzo Mannocci, Marco Avvenuti, Maurizio Tesconi, and Stefano Cresci. 2024. Coordinated behavior in information operations on Twitter. *IEEE Access* 11 (2024), 61568–61585.
- [19] Matteo Cinelli, Stefano Cresci, Walter Quattrociocchi, Maurizio Tesconi, and Paola Zola. 2022. Coordinated inauthentic behavior and information spreading on Twitter. *Decision Support Systems* 160 (2022), 113819.
- [20] Aaron Clauset, Mark EJ Newman, and Cristopher Moore. 2004. Finding community structure in very large networks. *Physical Review E* 70, 6 (2004), 066111.
- [21] Stefano Cresci. 2020. A decade of social bot detection. *Commun. ACM* 63, 10 (2020), 72–83.
- [22] Stefano Cresci, Fabrizio Lillo, Daniele Regoli, Serena Tardelli, and Maurizio Tesconi. 2018. \$FAKE: Evidence of spam and bot activity in stock microblogs on Twitter. In *AAAI ICWSM*. 580–583.
- [23] Stefano Cresci, Fabrizio Lillo, Daniele Regoli, Serena Tardelli, and Maurizio Tesconi. 2019. Cashtag Piggybacking: Uncovering spam and bot activity in stock microblogs on Twitter. *ACM Transactions on the Web* 13, 2 (2019), 11:1–11:27.
- [24] Stefano Cresci, Kai-Cheng Yang, Angelo Spognardi, Roberto Di Pietro, Filippo Menczer, and Marinella Petrocchi. 2024. Demystifying misconceptions in social bots research. [arXiv:2303.17251](https://arxiv.org/abs/2303.17251)
- [25] Adya Danaditya, Lynnette H. X. Ng, and Kathleen M. Carley. 2022. From curious hashtags to polarized effect: Profiling coordinated actions in Indonesian Twitter discourse. *Social Network Analysis and Mining* 12, 1 (2022), 105.
- [26] Bart De Clerck, Juan Carlos Fernandez Toledano, Filip Van Utterbeeck, and Luis EC Rocha. 2024. Detecting coordinated and bot-like behavior in Twitter: the Jürgen Conings case. *EPJ Data Science* 13, 1 (2024), 40.
- [27] Priyanka Dey, Luca Luceri, and Emilio Ferrara. 2024. Coordinated activity modulates the behavior and emotions of organic users: A case study on tweets about the Gaza conflict. In *ACM WWW*. 682–685.
- [28] N. Di Marco, Sara Brunetti, Matteo Cinelli, and Walter Quattrociocchi. 2024. Post-hoc evaluation of nodes influence in information cascades: The case of coordinated accounts. [arXiv:2401.01684](https://arxiv.org/abs/2401.01684)
- [29] Carlotta Dotto, Seb Cubbon, Stefano Cresci, Serena Tardelli, and Leonardo Nizzoli. 2021. How to improve our analysis of ‘coordinated inauthentic behavior’. <https://medium.com/1st-draft/how-to-improve-our-analysis-of-coordinated-inauthentic-behavior-a4ec62ce9bfb> (accessed: 31/07/2024).
- [30] Evelyn Douek. 2020. What does “coordinated inauthentic behaviour” actually mean? <https://slate.com/technology/2020/07/coordinated-inauthentic-behavior-facebook-twitter.html>. (accessed: 31/07/2024).
- [31] Auriant Emeric and Chomel Victor. 2023. Interpretable cross-platform coordination detection on social networks. In *CNA*. 143–155.
- [32] Keeley Erhardt and Dina Albassam. 2023. Detecting the hidden dynamics of networked actors using temporal correlations. In *ACM WWW*. 1214–1217.
- [33] Keeley Erhardt and Alex Pentland. 2024. Hidden messages: mapping nations’ media campaigns. *Computational and Mathematical Organization Theory* 30, 2 (2024), 161–172.
- [34] Fatima Ezzeddine, Omran Ayoub, Silvia Giordano, Gianluca Nogara, Ihab Sbeity, Emilio Ferrara, and Luca Luceri. 2023. Exposing influence campaigns in the age of LLMs: A behavioral-based AI approach to detecting state-sponsored trolls. *EPJ Data Science* 12, 1 (2023), 46.
- [35] Facebook. 2018. Coordinated inauthentic behavior explained. <https://about.fb.com/news/2018/12/inside-feed-coordinated-inauthentic-behavior/>. (accessed: 31/07/2024).
- [36] Facebook. 2018. Removing bad actors from Facebook. <https://about.fb.com/news/2018/07/removing-bad-actors-on-facebook/>. (accessed: 31/07/2024).
- [37] Facebook. 2019. How we respond to inauthentic behavior on our platforms: Policy update. <https://about.fb.com/news/2019/10/inauthentic-behavior-policy-update/>. (accessed: 31/07/2024).
- [38] Facebook. 2019. Inauthentic behavior. <https://transparency.fb.com/policies/community-standards/inauthentic-behavior/>. (accessed: 31/07/2024).
- [39] Facebook. 2019. Removing coordinated inauthentic behavior in UAE, Egypt and Saudi Arabia. <https://about.fb.com/news/2019/08/cib-uae-egypt-saudi-arabia/>. (accessed: 31/07/2024).
- [40] Facebook. 2021. Advancing our policies on online bullying and harassment. <https://about.fb.com/news/2021/10/advancing-online-bullying-harassment-policies/>. (accessed: 31/07/2024).
- [41] Facebook. 2021. Removing new types of harmful networks. <https://about.fb.com/news/2019/10/inauthentic-behavior-policy-update/>. (accessed: 31/07/2024).
- [42] Camille Francois, Vladimir Barash, and John Kelly. 2023. Measuring coordinated versus spontaneous activity in online social movements. *New Media & Society* 25, 11 (2023), 3065–3092.
- [43] Diego Garlaschelli and Maria I. Loffredo. 2008. Maximum likelihood: Extracting unbiased information from complex networks. *Physical Review E* 78 (2008), 015101. Issue 1.
- [44] Piyush Ghasiya and Kazutoshi Sasahara. 2022. Rapid sharing of Islamophobic hate on Facebook: The case of the Tablighi Jamaat controversy. *Social Media + Society* 8, 4 (2022).
- [45] Fabio Giglietto, Giada Marino, Roberto Mincigrucci, and Anna Stanziano. 2023. A workflow to detect, monitor, and update lists of coordinated social media accounts across time: The case of the 2022 Italian election. *Social Media + Society* 9, 3 (2023).
- [46] Fabio Giglietto, Nicola Righetti, Luca Rossi, and Giada Marino. 2020. Coordinated link sharing behavior as a signal to surface sources of problematic information on Facebook. In *ACM SMSociety*. 85–91.

- [47] Fabio Giglietto, Nicola Righetti, Luca Rossi, and Giada Marino. 2020. It takes a village to manipulate the media: Coordinated link sharing behavior during 2018 and 2019 Italian elections. *Information, Communication & Society* 23, 6 (2020), 867–891.
- [48] Fabio Giglietto, Massimo Terenzi, Giada Marino, Nicola Righetti, and Luca Rossi. 2020. Adapting to mitigation efforts: Evolving strategies of coordinated link sharing on Facebook. SSRN:3775469
- [49] Google. 2020. Updates about government-backed hacking and disinformation. <https://blog.google/threat-analysis-group/updates-about-government-backed-hacking-and-disinformation/>. (accessed: 31/07/2024).
- [50] Timothy Graham, Axel Bruns, Guangnan Zhu, and Rod Campbell. 2020. *Like a virus: The coordinated spread of coronavirus disinformation*. Technical Report. Centre for Responsible Technology, The Australia Institute.
- [51] Timothy Graham, Sam Hames, and Elizabeth Alpert. 2024. The coordination network toolkit: A framework for detecting and analysing coordinated behaviour on social media. *Journal of Computational Social Science* (2024), 1–22.
- [52] Anatoliy Gruzd, Philip Mai, and Felipe B. Soares. 2022. How coordinated link sharing behavior and partisans’ narrative framing fan the spread of COVID-19 misinformation and conspiracy theories. *Social Network Analysis and Mining* 12, 1 (2022), 118.
- [53] Bing He, Mustaque Ahamad, and Srikan Kumar. 2021. PETGEN: Personalized text generation attack on deep sequence embedding-based classification models. In *ACM KDD*.
- [54] Philip N. Howard, Aiden Duffy, Deen Freelon, Muzammil M. Hussain, Will Mari, and Marwa Maziad. 2011. Opening closed regimes: What was the role of social media during the Arab Spring? SSRN:2595096
- [55] Kristina Hristakieva, Stefano Cresci, Giovanni Da San Martino, Mauro Conti, and Preslav Nakov. 2022. The spread of propaganda by coordinated communities on social media. In *ACM WebSci*. 191–201.
- [56] Roberto Interdonato, Matteo Magnani, Diego Perna, Andrea Tagarelli, and Davide Vega. 2020. Multilayer network simplification: Approaches, models and methods. *Computer Science Review* 36 (2020), 100246.
- [57] Laura Jahn and Rasmus K. Rendsvig. 2023. Towards detecting inauthentic coordination in Twitter likes data. arXiv:2305.07384
- [58] Laura Jahn, Rasmus K. Rendsvig, and Jacob Stærk-Østergaard. 2023. Detecting coordinated inauthentic behavior in likes on social media: Proof of concept. arXiv:2305.07350
- [59] Maurice Jakesch, Kiran Garimella, Dean Eckles, and Mor Naaman. 2021. Trend alert: A cross-platform organization manipulated Twitter trends in the Indian general election. In *ACM CSCW*. 1–19.
- [60] Julian Jaurisch, Jakob Ohme, and Ulrike Klinger. 2024. *Enabling research with publicly accessible platform data: Early DSA compliance issues and suggestions for improvement*. Technical Report. Weizenbaum Institute.
- [61] Franziska B. Keller, David Schoch, Sebastian Stier, and JungHwan Yang. 2017. How to manipulate social media: Analyzing political astroturfing using ground truth data from South Korea. In *AAAI ICWSM*. 564–567.
- [62] Franziska B. Keller, David Schoch, Sebastian Stier, and JungHwan Yang. 2020. Political astroturfing on Twitter: How to coordinate a disinformation campaign. *Political Communication* 37, 2 (2020), 256–280.
- [63] Baris Kirdemir, Oluwaseyi Adeliyi, and Nitin Agarwal. 2022. Towards characterizing coordinated inauthentic behaviors on YouTube. In *ROMCIR*. 100–116.
- [64] Charles Kriel and Alexa Pavliuc. 2019. Reverse engineering Russian Internet Research Agency tactics through network analysis. *Defence Strategic Communication* 6 (2019).
- [65] Aytalina Kulichkina, Nicola Righetti, and Annie Waldherr. 2024. Protest and repression on social media: Pro-Navalny and pro-government mobilization dynamics and coordination patterns on Russian Twitter. *New Media & Society* (2024).
- [66] Kyumin Lee, James Caverlee, Zhiyuan Cheng, and Daniel Z. Sui. 2013. Campaign extraction from social media. *ACM Transactions on Intelligent Systems and Technology* 5, 1 (2013), 9:1–9:28.
- [67] Renan S. Linhares, José M. Rosa, Carlos H. G. Ferreira, Fabricio Murai, Gabriel Nobre, and Jussara Almeida. 2022. Uncovering coordinated communities on Twitter during the 2020 US election. In *IEEE/ACM ASONAM*. 80–87.
- [68] Sahil Loomba, Alexandre De Figueiredo, Simon J Piatek, Kristen De Graaf, and Heidi J Larson. 2021. Measuring the impact of COVID-19 vaccine misinformation on vaccination intent in the UK and USA. *Nature Human Behaviour* 5, 3 (2021), 337–348.
- [69] Edoardo Loru, Matteo Cinelli, Maurizio Tesconi, and Walter Quattrociocchi. 2023. The influence of coordinated behavior on toxicity. arXiv:2310.01283
- [70] Lorenzo Lucchini, Luca M. Aiello, Laura Alessandretti, Gianmarco De Francisci Morales, Michele Starnini, and Andrea Baronchelli. 2022. From Reddit to Wall Street: The role of committed minorities in financial collective action. *Royal Society Open Science* 9, 4 (2022), 211488.
- [71] Luca Luceri, Silvia Giordano, and Emilio Ferrara. 2020. Detecting troll behavior via inverse reinforcement learning: A case study of Russian trolls in the 2016 US election. In *AAAI ICWSM*. 417–427.
- [72] Luca Luceri, Valeria Pantè, Keith Burghardt, and Emilio Ferrara. 2024. Unmasking the Web of deceit: Uncovering coordinated activity to expose information operations on Twitter. In *ACM WWW*.
- [73] Thomas Magelinski and Kathleen M. Carley. 2020. *Detecting coordinated behavior in the Twitter campaign to reopen America*. Technical Report. School of Computer Science, Carnegie Mellon University.
- [74] Thomas Magelinski, Lynnette H. X. Ng, and Kathleen Carley. 2022. A synchronized action framework for detection of coordination on social media. *Journal of Online Trust and Safety* 1, 2 (2022).
- [75] Matteo Magnani, Obaida Hanteer, Roberto Interdonato, Luca Rossi, and Andrea Tagarelli. 2021. Community detection in multiplex networks. *ACM Computing Surveys (CSUR)* 54, 3 (2021), 1–35.

- [76] Thomas W. Malone. 1987. Modeling coordination in organizations and markets. *Management Science* 33, 10 (1987), 1317–1332.
- [77] Thomas W. Malone. 1988. *What is coordination theory?* Technical Report 182. Massachusetts Institute of Technology (MIT), Sloan School of Management.
- [78] Thomas W. Malone and Kevin Crowston. 1990. What is coordination theory and how can it help design cooperative work systems?. In *ACM CSCW*. 357–370.
- [79] Thomas W. Malone and Stephen A. Smith. 1988. Modeling the performance of organizational structures. *Operations Research* 36, 3 (1988), 421–436.
- [80] Isura Manchanayaka, Zainab R. Zaidi, Shanika Karunasekera, and Christopher Leckie. 2024. Identifying coordinated activities on online social networks using contrast pattern mining. In *IEEE IJCNN*.
- [81] Isura Manchanayaka, Zainab R. Zaidi, Shanika Karunasekera, and Christopher Leckie. 2024. Using causality to infer coordinated attacks in social media. [arXiv:2407.11690](https://arxiv.org/abs/2407.11690)
- [82] Lorenzo Mannocci, Stefano Cresci, Anna Monreale, Athina Vakali, and Maurizio Tesconi. 2022. MulBot: Unsupervised bot detection based on multivariate time series. In *IEEE BigData*. 1485–1494.
- [83] Enrico Mariconti, Guillermo Suarez-Tangil, Jeremy Blackburn, Emiliano De Cristofaro, Nicolas Kourtellis, Ilias Leontiadis, Jordi L. Serrano, and Gianluca Stringhini. 2019. “You Know What to Do”: Proactive detection of YouTube videos targeted by coordinated hate attacks. In *ACM CSCW*. 1–21.
- [84] Michele Mazza, Guglielmo Cola, and Maurizio Tesconi. 2022. Ready-to-(ab) use: From fake account trafficking to coordinated inauthentic behavior on Twitter. *Online Social Networks and Media* 31 (2022), 100224.
- [85] Michele Mazza, Stefano Cresci, Marco Avvenuti, Walter Quattrociocchi, and Maurizio Tesconi. 2019. RTbust: Exploiting temporal patterns for botnet detection on Twitter. In *ACM WebSci*. 183–192.
- [86] Swapneel S. Mehta, Atilim G. Baydin, Bogdan State, Richard Bonneau, Jonathan Nagler, and Philip Torr. 2022. Estimating the impact of coordinated inauthentic behavior on content recommendations in social networks. In *AI4ABM*.
- [87] Kumari Neha, Vibhu Agrawal, Saurav Chhatani, Rajesh Sharma, Arun Balaji Buduru, and Ponnuram Kumaraguru. 2024. Understanding coordinated communities through the lens of protest-centric narratives: A case study on #CAA protest. In *AAAI ICWSM*, Vol. 18. 1123–1133.
- [88] Kin W. Ng and Adriana Iamnitchi. 2023. Coordinated information campaigns on social media: A multifaceted framework for detection and analysis. In *MISDOOM*. 103–118.
- [89] Lynnette H. X. Ng and Kathleen M. Carley. 2022. A combined synchronization index for grassroots activism on social media. [arXiv:2212.13221](https://arxiv.org/abs/2212.13221)
- [90] Lynnette H. X. Ng and Kathleen M. Carley. 2022. Online coordination: Methods and comparative case studies of coordinated groups across four events in the United States. In *ACM WebSci*. 12–21.
- [91] Lynnette H. X. Ng and Kathleen M. Carley. 2023. Do you hear the people sing? Comparison of synchronized URL and narrative themes in 2020 and 2023 French protests. *Frontiers in Big Data* 6 (2023), 1343108.
- [92] Lynnette H. X. Ng, Iain Cruickshank, and Kathleen M. Carley. 2021. Coordinating narratives and the Capitol riots on Parler. [arXiv:2109.00945](https://arxiv.org/abs/2109.00945)
- [93] Lynnette H. X. Ng, Iain J. Cruickshank, and Kathleen M. Carley. 2022. Cross-platform information spread during the January 6th Capitol riots. *Social Network Analysis and Mining* 12, 1 (2022), 133.
- [94] Lynnette H. X. Ng, Iain J. Cruickshank, and Kathleen M. Carley. 2023. Coordinating narratives framework for cross-platform analysis in the 2021 US Capitol riots. *Computational and Mathematical Organization Theory* 29, 3 (2023), 470–486.
- [95] Lynnette H. X. Ng, J. D. Moffitt, and Kathleen M. Carley. 2022. Coordinated through a Web of images: Analysis of image-based influence operations from China, Iran, Russia, and Venezuela. [arXiv:2206.03576](https://arxiv.org/abs/2206.03576)
- [96] Sophie J. Nightingale and Hany Farid. 2022. AI-synthesized faces are indistinguishable from real faces and more trustworthy. *Proceedings of the National Academy of Sciences* 119, 8 (2022).
- [97] Ben Nimmo, Ira Hubert, and Yang Cheng. 2021. Spamouflage breakout: Chinese spam network finally starts to gain some traction. <https://graphika.com/reports/spamouflage-breakout>. (accessed: 31/07/2024).
- [98] Leonardo Nizzoli, Serena Tardelli, Marco Avvenuti, Stefano Cresci, and Maurizio Tesconi. 2021. Coordinated behavior on social media in 2019 UK General Election. In *AAAI ICWSM*. 443–454.
- [99] United States. Office of the Director of National Intelligence. 2017. *Assessing Russian activities and intentions in recent US Elections*.
- [100] Diogo Pacheco, Alessandro Flammini, and Filippo Menczer. 2020. Unveiling coordinated groups behind White Helmets disinformation. In *ACM WWW*. 611–616.
- [101] Diogo Pacheco, Pik-Mai Hui, Christopher Torres-Lugo, Bao T. Truong, Alessandro Flammini, and Filippo Menczer. 2021. Uncovering coordinated networks on social media: Methods and case studies. In *AAAI ICWSM*. 455–466.
- [102] Sen Pei, Lev Muchnik, José S Andrade, Jr, Zhiming Zheng, and Hernán A Makse. 2014. Searching for superspreaders of information in real-world social media. *Scientific Reports* 4, 1 (2014), 5547.
- [103] Janina Pohl, Dennis Assenmacher, Moritz V. Seiler, Heike Trautmann, and Christian Grimme. 2022. Artificial social media campaign creation for benchmarking and challenging detection approaches. In *AAAI ICWSM*.
- [104] Reddit. 2022. 2022 Transparency report. <https://www.redditinc.com/policies/2022-transparency-report/>. (accessed: 31/07/2024).
- [105] Nicola Righetti, Fabio Giglietto, Azade E. Kakavand, Aytalina Kulichkina, Giada Marino, and Massimo Terenzi. 2022. Political advertisement and coordinated behavior on social media in the lead-up to the 2021 German federal elections. *Dusseldorf: Media Authority of North Rhine-Westphalia* (2022).

- [106] Ronald E Robertson. 2022. Uncommon yet consequential online harms. *Journal of Online Trust and Safety* 1, 3 (2022).
- [107] Mohammad H. Saeed, Kostantinos Papadamou, et al. 2024. TUBERAIDER: Attributing coordinated hate attacks on YouTube videos to their source communities. In *AAAI ICWSM*, Vol. 18. 1354–1366.
- [108] Marcel Schliebs, H. Bailey, J. Bright, and P. N. Howard. 2021. *China's inauthentic UK Twitter diplomacy: A coordinated network amplifying PRC diplomats*. Technical Report. Programme on Democracy and Technology, Oxford University.
- [109] David Schoch, Franziska B. Keller, Sebastian Stier, and JungHwan Yang. 2022. Coordination patterns reveal online political astroturfing across the world. *Scientific Reports* 12, 1 (2022), 4572.
- [110] Fatih Şen, Rolf T. Wigand, Nitin Agarwal, Serpil Yuce Tokdemir, and Rafal Kasprzyk. 2016. Focal structures analysis: Identifying influential sets of individuals in a social network. *Social Network Analysis and Mining* 6, 1 (2016), 17:1–17:22.
- [111] M. Ángeles Serrano, Marián Boguná, and Alessandro Vespignani. 2009. Extracting the multiscale backbone of complex weighted networks. *Proceedings of the National Academy of Sciences* 106, 16 (2009), 6483–6488.
- [112] Karishma Sharma, Yizhou Zhang, Emilio Ferrara, and Yan Liu. 2021. Identifying coordinated accounts on social media through hidden influence and group behaviours. In *ACM SIGKDD*. 1441–1451.
- [113] Kai Shu, Suhang Wang, Dongwon Lee, and Huan Liu. 2020. Mining disinformation and fake news: Concepts, methods, and recent advancements. In *Disinformation, misinformation, and fake news in social media*. 1–19.
- [114] Christophe Sibertin-Blanc, Frédéric Amblard, and Matthias Maillard. 2005. A coordination framework based on the sociology of organized action. In *AAMAS*. 3–17.
- [115] Felipe B. Soares. 2023. From sharing misinformation to debunking it: How coordinated image text sharing behaviour is used in political campaigns on Facebook. In *MISDOOM*. 45–59.
- [116] Lucas Stampe, Janina Pohl, and Christian Grimme. 2023. Towards multimodal campaign detection: Including image information in stream clustering to detect social media campaigns. In *MISDOOM*. 144–159.
- [117] Kate Starbird, Ahmer Arif, and Tom Wilson. 2019. Disinformation as collaborative work: Surfacing the participatory nature of strategic information operations. In *ACM CSCW*. 127:1–127:26.
- [118] Zachary C. Steinert-Threlkeld, Delia Mocanu, Alessandro Vespignani, and James H. Fowler. 2015. Online social networks and offline protest. *EPJ Data Science* 4, 1 (2015), 19.
- [119] Serena Tardelli, Marco Avvenuti, Maurizio Tesconi, and Stefano Cresci. 2020. Characterizing social bots spreading financial disinformation. In *HCI*. 376–392.
- [120] Serena Tardelli, Leonardo Nizzoli, Marco Avvenuti, Stefano Cresci, and Maurizio Tesconi. 2024. Multifaceted online coordinated behavior in the 2020 US presidential election. *EPJ Data Science* 13, 1 (2024), 1–27.
- [121] Serena Tardelli, Leonardo Nizzoli, Maurizio Tesconi, Mauro Conti, Preslav Nakov, Giovanni Da San Martino, and Stefano Cresci. 2024. Temporal dynamics of coordinated online behavior: Stability, archetypes, and influence. *Proceedings of the National Academy of Sciences* (2024).
- [122] Petter Törnberg, Diliara Valeeva, Justus Uitermark, and Christopher Bail. 2023. Simulating social media using large language models to evaluate alternative news feed algorithms. [arXiv:2310.05984](https://arxiv.org/abs/2310.05984)
- [123] Vincent A. Traag, Ludo Waltman, and Nees Jan Van Eck. 2019. From Louvain to Leiden: Guaranteeing well-connected communities. *Scientific Reports* 9, 1 (2019), 5233.
- [124] Rebekah Tromble. 2021. Where have all the data gone? A critical reflection on academic digital research in the post-API age. *Social Media + Society* 7, 1 (2021).
- [125] Amaury Trujillo, Tiziano Fagni, and Stefano Cresci. 2024. The DSA Transparency Database: Auditing self-reported moderation actions by social media. [arXiv:2312.10269](https://arxiv.org/abs/2312.10269)
- [126] Michele Tumminello, Salvatore Micciche, Fabrizio Lillo, Jyrki Piilo, and Rosario N. Mantegna. 2011. Statistically validated networks in bipartite complex systems. *PLoS One* 6, 3 (2011), 1–11.
- [127] Michael T. Turvey. 1990. Coordination. *American Psychologist* 45, 8 (1990), 938.
- [128] Twitter. 2018. Enabling further research of information operations on Twitter. https://blog.twitter.com/en_us/topics/company/2018/enabling-further-research-of-information-operations-on-twitter. (accessed: 31/07/2024).
- [129] Twitter. 2021. Coordinated harmful activity. <https://help.twitter.com/en/rules-and-policies/coordinated-harmful-activity>. (accessed: 31/07/2024).
- [130] Luis Vargas, Patrick Emami, and Patrick Traynor. 2020. On the detection of disinformation campaign activity with network analysis. In *ACM CCSW*. 133–146.
- [131] Vítor V. Vasconcelos, Sara M. Constantino, Astrid Dannenberg, Marcel Lumkowsky, Elke Weber, and Simon Levin. 2021. Segregation and clustering of preferences erode socially beneficial coordination. *Proceedings of the National Academy of Sciences* 118, 50 (2021).
- [132] Otavio R. Venâncio, Carlos H. G. Ferreira, Jussara M. Almeida, and Ana P. C. da Silva. 2024. Unraveling user coordination on Telegram: A comprehensive analysis of political mobilization during the 2022 Brazilian Presidential election. In *AAAI ICWSM*, Vol. 18. 1545–1556.
- [133] Padinjaredath Suresh Vishnuprasad, Gianluca Nogara, Felipe Cardoso, Stefano Cresci, Silvia Giordano, and Luca Luceri. 2024. Tracking fringe and coordinated activity on Twitter leading up to the US Capitol attack. In *AAAI ICWSM*, Vol. 18. 1557–1570.
- [134] Xinyu Wang, Jiayi Li, Eesha Srivatsavaya, and Sarah Rajtmajer. 2023. Evidence of inter-state coordination amongst state-backed information operations. *Scientific Reports* 13, 1 (2023), 7716.

- [135] Claire Wardle and Hossein Derakhshan. 2017. *Information Disorder: Toward an Interdisciplinary Framework for Research and Policymaking*. Council of Europe report DGI (2017)09. Council of Europe.
- [136] Derek Weber and Lucia Falzon. 2021. Temporal nuances of coordination networks. [arXiv:2107.02588](https://arxiv.org/abs/2107.02588)
- [137] Derek Weber and Frank Neumann. 2020. Who’s in the gang? Revealing coordinating communities in social media. In *IEEE/ACM ASONAM*. 89–93.
- [138] Derek Weber and Frank Neumann. 2021. Amplifying influence through coordinated behaviour in social networks. *Social Network Analysis and Mining* 11, 1 (2021), 111.
- [139] Tom Wilson and Kate Starbird. 2020. Cross-platform disinformation campaigns: Lessons learned and next steps. *HKS Misinformation Review* (2020).
- [140] William E. S. Yu. 2022. A framework for studying coordinated behaviour applied to the 2019 Philippine midterm elections. In *ICICT*. 721–731.
- [141] Savvas Zannettou, Tristan Caulfield, William Setzer, Michael Sirivianos, Gianluca Stringhini, and Jeremy Blackburn. 2019. Who let the trolls out?: Towards understanding state-sponsored trolls. In *ACM WebSci*. 353–362.
- [142] Yizhou Zhang, Karishma Sharma, and Yan Liu. 2021. VigDet: Knowledge informed neural temporal point process for coordination detection on social media. In *NeurIPS*. 3218–3231.
- [143] Yizhou Zhang, Karishma Sharma, and Yan Liu. 2023. Capturing cross-platform interaction for identifying coordinated accounts of misinformation campaigns. In *ECIR*. 694–702.