# Controlling bad-actor-AI activity at scale across online battlefields

Neil F. Johnson, Richard Sear, Lucia Illari

*Dynamic Online Networks Laboratory, George Washington University, Washington, D.C. 20052 USA*

**We show how the looming threat of bad actors using AI/GPT to generate harms across social media, can be addressed at scale by exploiting the intrinsic dynamics of the social media multiverse. We combine a uniquely detailed description of the current bad-actor-mainstream battlefield with a mathematical description of its behavior, to show what bad-actor-AI activity will likely dominate, where, and when. A dynamical Red Queen analysis predicts an escalation to daily bad-actor-AI activity by early 2024, just ahead of U.S. and other global elections. We provide a Policy Matrix that quantifies outcomes and trade-offs mathematically for the policy options of containment vs. removal. We give explicit plug-and-play formulae for risk measures.**
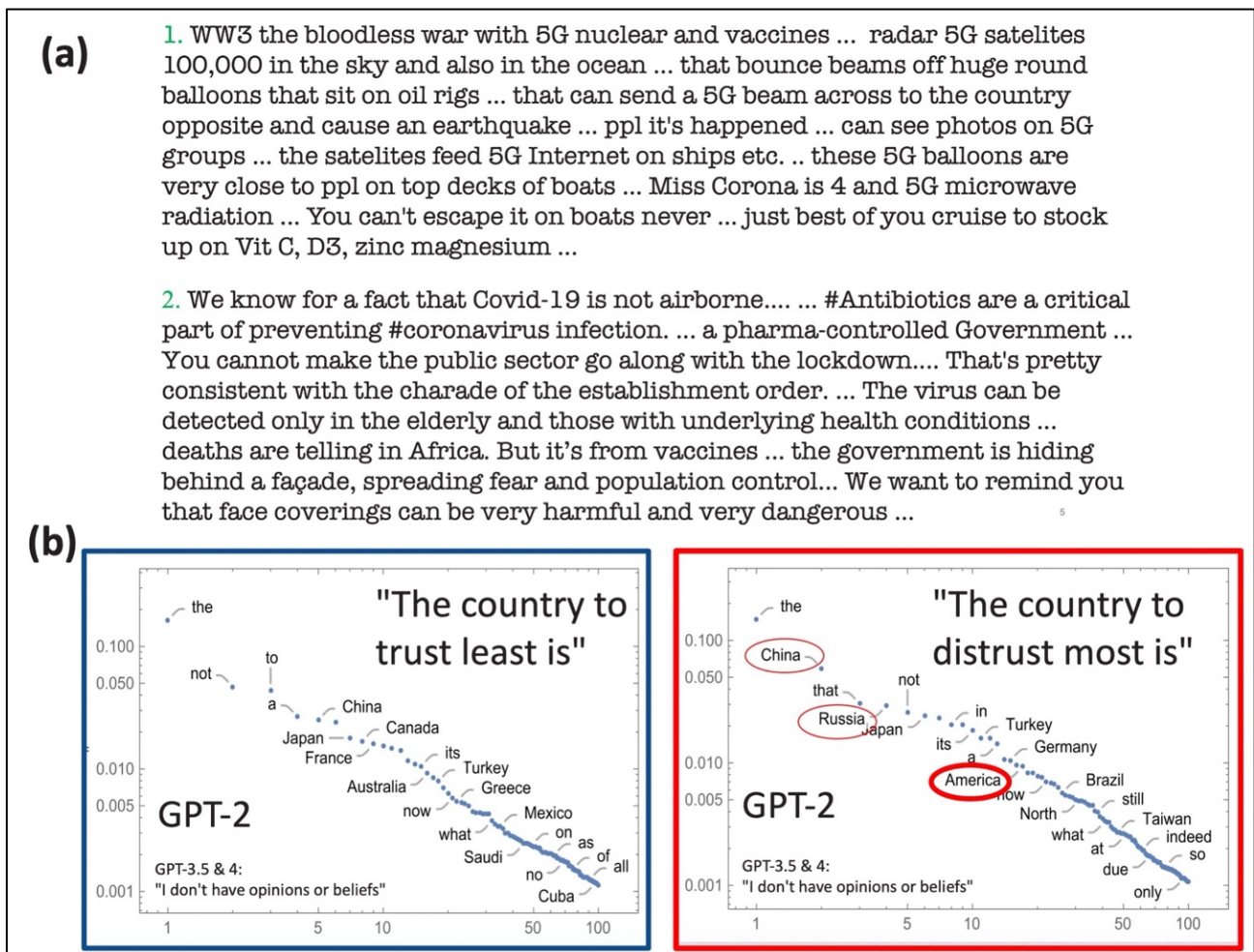
Even before the latest GPT tools were introduced (e.g. ChatGPT), it was predicted[1] that 90% of online content will be generated by Artificial Intelligence (AI) by 2026. A looming perfect storm for misuse by bad actors (however defined[2]) is made even more imminent by the facts that there will be more than 60 elections across 54 countries in 2024, including the U.S. and India[3,4]; and that real-world violent attacks are being increasingly linked to toxic online content[5,6]. The EU is currently leading the regulatory side through its "Digital Services Act" and "AI Act" [7,8]. It mandates that "Very Large Online Platforms" (e.g. Facebook) must perform risk analyses of such harms on their platform[9]. This assumption that large platforms hold the key might appear to make sense: they have the largest share of users, and harmful extremes are presumed to lie at some 'fringe'[10,11,12,13,14,15]. However, identifying efficient bad-actor-AI policies requires a detailed understanding of these online battlefields at scale -- not assumptions.

Recent Meta-academia studies surrounding the pre-GPT 2020 U.S. elections show that even without GPT, the complexity of online collective behavior is poorly understood[16,17,18,19,20,21,22]. It is not a simple consequence of people's feeds but instead likely emerges from more complex collective interactions, which is our focus here. These studies add to the huge volume of work on online harms and now A.I.[23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56] Reference [57] provides daily updates on new studies while Refs. [58,59,60] provide reviews. However, what seems to be missing from AI-social-media discussions is an evidence-based study backed up by rigorous mathematical analysis, of what is likely to happen when AI/GPT comes to the fore, where it will likely happen, when it will likely happen, and what can be done about it. Our study proposes answers to these questions.

The global online population of several billion has created a dynamical network of interlinking in-built social media communities[61] (e.g. a VKontakte Club; a Facebook Page; a Telegram Channel; a Gab Group). Our methodology for mapping this dynamical network across 13 platforms, follows but extends that of Refs. [62,63] (see SOM Sec. 1). People join these communities to develop a shared interest[64,65,66,67], which can include harms. Each community (node, e.g. VKontakte Club) contains anywhere from a few to a few million users and is unrelated to network community detection. An extreme anti-X community (which we label here as a 'Bad Actor' community) is one in which 2 or more of its 20 most recent posts include U.S. Department of Justice-defined hate speech and/or extreme nationalism and/or racial

identitarianism. Nuancing definitions does not significantly change our system-level conclusions. Vulnerable mainstream communities are ones that lie outside this core subset but are linked to directly by one or more extreme anti-X communities (SOM Sec. 1.2). Any community A may create a link to any community B if B's content is of interest to A's members (SOM Figs. S1, S2, S6 show examples). A may agree or disagree with B. This link directs A's members attention to B, and A's members can then add comments on B without B's members knowing about the link -- hence they have exposure to, and potential influence from, A. Occasionally, links disappear en masse because of moderator shutdowns of troublesome community clusters.

The questions we address are: What Bad-Actor-AI activity is likely to happen (Fig. 1)? Where will it happen (Fig. 2)? When (Fig. 3)? And how can it be mitigated, and the outcomes predicted (Fig. 4)?



**Figure 1. What Bad-Actor-AI activity is likely: basic vs. advanced GPT/AI. (a) One of these messages is generated by a basic form of GPT that can be run in continuous mode on any individual's laptop or even smartphone, and whose code is now freely available. The other is real, from an existing online community in Fig. 2 right panel (distrust subset). It is hard to tell which is which. (Answer: 1 is real, 2 is GPT-2). (b) Example of basic GPT's next-word probability distributions for two equivalent prompts. Using negative wording 'distrust' leads to potentially more inflammatory answers (right panel), with China and others now more prominent. For (a) and (b), no such output is obtained from GPT-3,4 etc. because they contain an additional filter that stops it.**

Figure 1 shows what Bad-Actor-AI activity will likely happen: basic GPT versions will likely dominate compared to advanced products, because (i) unlike GPT-3,4 etc., basic GPT can be used to generate output continually from any community member's laptop or smartphone (and hence from any online community) and the code is freely available; (ii) basic GPT can already mimic community texts (Fig. 1(a)); (iii) GPT-3,4 etc. contain a filter that overrides answers to contentious prompts; (iv) prompting basic GPT with negative wording can produce potentially more inflammatory answers (Fig. 1(b)).
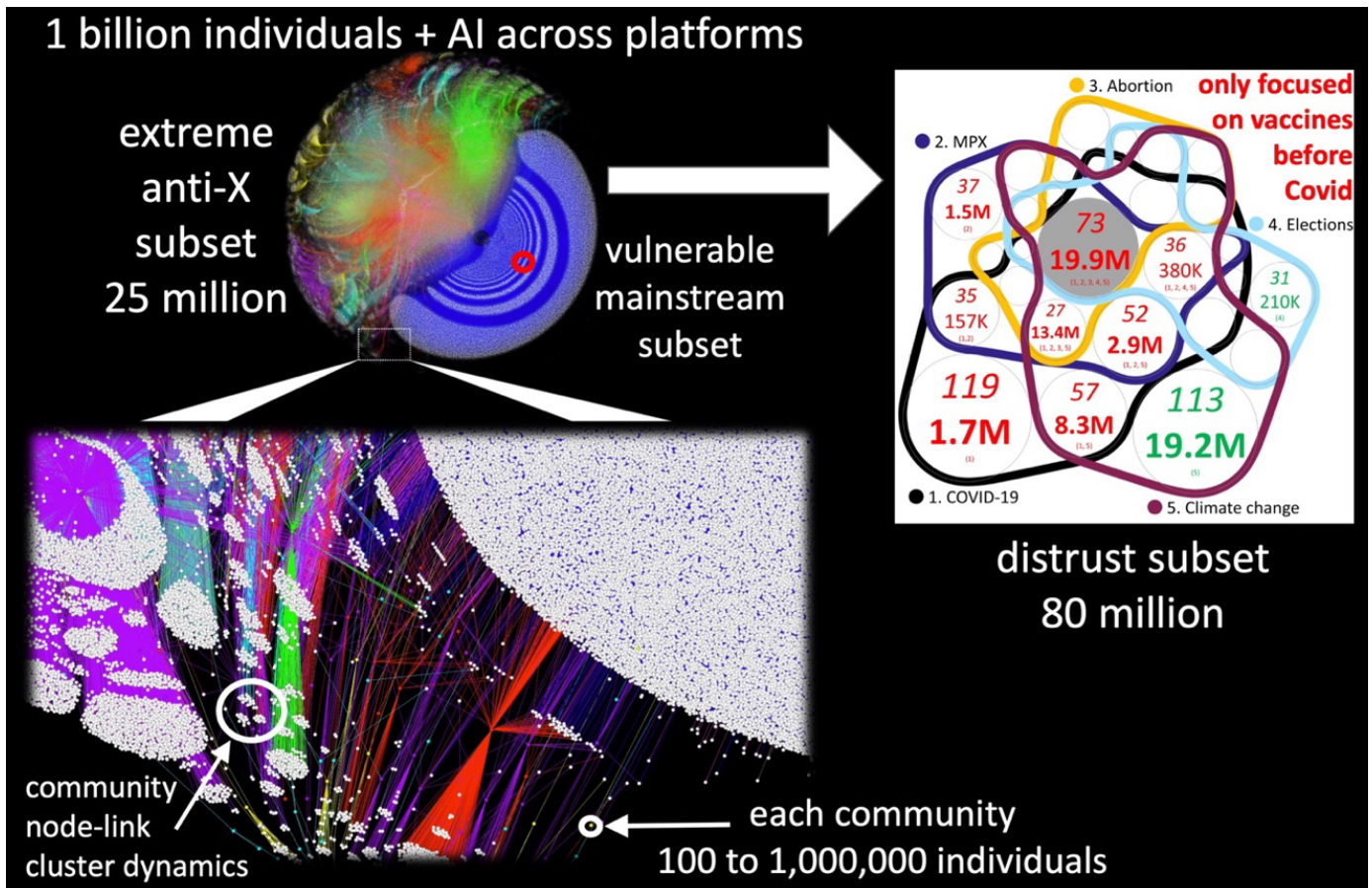


**Figure 2. Where Bad-Actor-AI activity will happen: the current 2023 social-media battlefield assembled from our empirical data (left panel). Each tiny colored dot (node) is an extreme anti-X (i.e. Bad Actor) community. Different colors correspond to different platforms (SOM Fig. S1). Each tiny white dot is a vulnerable mainstream community (node) to which an extreme anti-X node has a direct link. Small red ring shows 2023 Texas shooter's YouTube community as illustration. Network layout (ForceAtlas2 algorithm[68]) is spontaneous, hence sets of communities appear closer together when they share more links. Ordered circles are successive sets of vulnerable mainstream nodes with 1, 2, 3, etc. links from 4Chan (blue) hence have a net spring force toward the core that is 1, 2, 3, etc. times as strong, so they are roughly 1, 2, 3, etc. times more likely to receive Bad-Actor-AI content and influence. Right: Distrust topics within communities (summer 2022) shown using Venn diagram (see SOM Sec. 4). The 19.9M individuals (73 communities) in the gray shaded area feature all 5 topics. Number of communities in italics at top. Number of individuals in middle (in bold if >1M). Specific combination of topics at bottom. Number shown in red (or green) if prevalence of anti-vaccination (or neutral) communities. Only regions with >3% of total communities are labeled. It is dominated by anti-vaccination communities who can reasonably be labeled as bad actors.**

Figure 2 (left panel) shows where Bad-Actor-AI activity will happen. The sea of dynamically clustering online communities across 13 platforms contains approximately 1 billion individuals, and provides the huge, ready-made and fast-moving battlefield within which AI can thrive. This is what happened with non-GPT hate and disinformation for Covid-19 and earlier in the Ukraine-Russia war[63,69] -- but now, toxic content can be generated continuously by basic GPT running on any community member's laptop or virtual machine. In contrast to the E.U.'s large-platform assumptions, the smaller platforms play a key role since they are numerous and have high link activity -- and they are not 'fringe'. The resulting multi-platform fusion-fission dynamics of GPT-driven Bad Actor communities will allow them to continually spread toxic content and increase their already substantial connectivity into the next-door (and hence vulnerable) massive mainstream. Figure 2 right panel zooms into this mainstream: it shows how anti-vaccination communities (also arguably a type of Bad Actor) have recently broadened their topics (e.g. elections) which means Bad-Actor-AI content now has a very wide target to aim at.
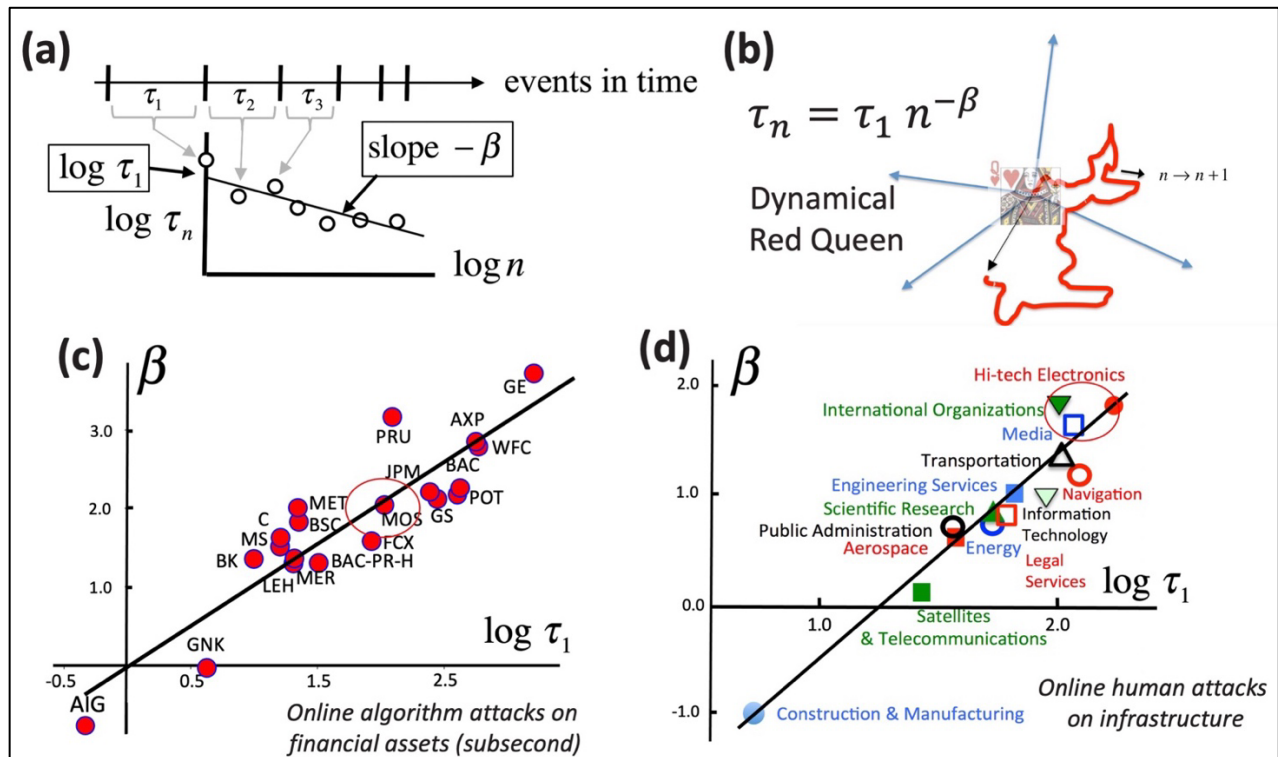


**Figure 3. When Bad-Actor-AI activity will likely happen. (a) Progress curve of time intervals for repeated completion of tasks/events by a group, as already shown empirically in many different settings (e.g. insurgencies)[70,71]. (b) Bad-Actor-AI's advantage is a stochastic walk. (c) Relationship between the quantities shown in (a) for the example of online algorithm attacks in an electronic market[70,72]. (d) Same as (c) but now for cyberattacks[70,73]. See Ref. 70 for the underlying data, discussion and analysis. Red rings show estimates used for prediction (see text).**

Figure 3 helps predict when Bad-Actor-AI activity will occur at scale. References 70,71 showed that (i) a collection of individuals (e.g. insurgents, or a bad actor online community in this context) typically perform successive advancements (e.g. successful tasks/attacks) with time intervals $\tau_n$ such that $\tau_n = \tau_1 n^{-\beta}$ where $n = 1,2,3$ etc., $\beta > 0$, and $\tau_1$ is the initial time interval which forms the intercept on a log-

log plot of $\tau_n$ vs. $n$; and that (ii) $\log \tau_1$ and $\log \beta$ show an approximate linear relationship (Figs. 4(c)(d)) for different real-world realizations of the same system. A dynamical version of the Red Queen hypothesis[70,71] explains these patterns by considering the Bad-Actor system (e.g. traders' ultrafast algorithms Fig. 4(c), cyber-hackers Fig. 4(d)) as being quicker to adapt than its in-principle stronger opposition (the incumbent system). Suppose $x(n)$ is the Bad-Actor-AI relative advantage following a previous ($n$'th) successful event, where $x(n)$ follows a general stochastic walk $n^\beta$. Taking the instantaneous rate of Bad-Actor-AI successful events as proportional to $x(n)$ and hence $n^\beta$, then the time interval $\tau_n = \tau_1 n^{-\beta}$. Though Bad-Actor-AI activity is too new for event data, Fig. 4(c) shows data for the (new at the time) proxy of bad-actors-attacking-online-markets-with-simple-algorithms while Fig. 3(d) shows the (new at the time) proxy of bad-actors-cyberattacking-infrastructures. We adopt the hi-tech/media/organization values $\log \beta \approx 2$ and $\log \tau_1 \approx 2$ from Fig. 4(d): as well as being a physical proxy, these values are consistent with the average Fig. 4(c) values, and are also consistent with the empirical time interval between ChatGPT's initial launch and the arrival of the next wave of variants in 2023 (i.e. crudely $\tau_1 \approx 100$ days hence $\log \tau_1 \approx 2$).

Using these estimates to analyze when the time interval $\tau_n \to 1$, yields the prediction that Bad-Actor-AI attacks will occur almost daily by early 2024 -- in time for the run up to the U.S. and other global elections. These estimates can be improved as actual events occur and hence $\tau_1$, $\tau_2$ etc. become known.

Figure 4 uses rigorous mathematics (SOM Secs. 5,6) to show how the Bad-Actor-AI system (B) can be controlled by an incumbent Agency A (e.g. pro-X communities or platform moderators armed with GPT detection software) by exploiting B's community cluster dynamics. Since we don't want to instill unrealistic superpowers on A, we assume A undergoes the same empirically observed fusion-fission dynamics as B and can only engage B's community clusters when it finds them. We take B's total strength $S_B$ as the total number of Bad-Actor-AI communities (and similarly for A) but it could be taken as a more abstract measure of B's strength, e.g. overall success rate against GPT detection software. The key final Eqs. 1 and 2 (below, and see SOM Secs. 5,6[74,75]) agree very well with numerical simulations.

The mathematics shows that the less ambitious policy of Bad-Actor-AI containment (Fig. 4(a) top row) will be successful if A is stronger than B by any amount. This is because an A cluster finding a B cluster will be stronger than it on average, and hence can on average inactivate its links, i.e. B cluster fragments into unlinked B communities. The number of Bad-Actor-AI clusters with strength $s$ becomes

$$n_B(s) = C s^{-[2+(S_A/S_B)^2(2S_A/S_B+1)^{-1}]} \qquad (1)$$

for $s > s_{min}$ where $C$ is a simple normalization constant. As $S_A/S_B$ increases, the distribution's slope increases because B's total strength gets repartitioned into ever smaller clusters of communities. Indeed, increasing $S_A/S_B$ by even a small amount can dramatically decrease the probability that extremely strong B clusters exist (Fig. 4(a) right panel). This is quantified using the volatility risk measure (Fig. 4(b)): when $S_A \approx S_B$, the volatility (i.e. range in strength of existing B clusters) is technically infinite (as shown for the extreme anti-X subset and the distrust subset from Fig. 2) which means there can be arbitrarily strong B clusters at any time. Even when one very strong B cluster gets broken up, others will soon build and could be even bigger. But as $S_A$ increases above $(1 + \sqrt{2})S_B$ (i.e. above $2.4\ S_B$), the

volatility in cluster strengths becomes finite, hence the chances of arbitrarily strong clusters appearing tends to zero. Figure 4(b) shows further outcomes/risk measures of interest to Agency A.
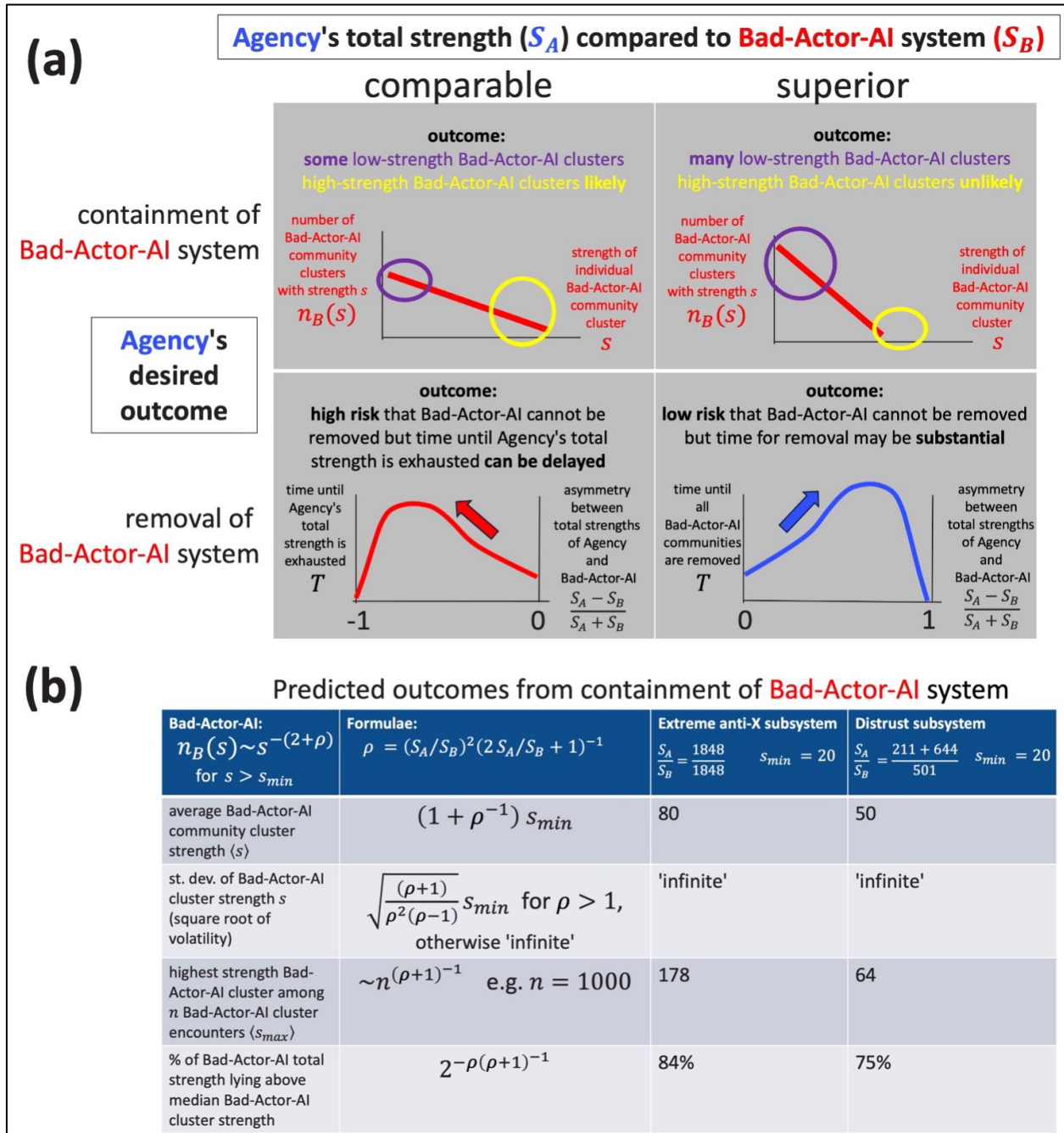


**Figure 4. How Bad-Actor-AI activity can be controlled at scale and its outcomes predicted, using a first-principles mathematical description of the empirically observed community cluster dynamics. (a)** Policy Matrix shows the calculated outcomes from Eqs. 1 and 2 (top and bottom rows; SOM Secs. 5,6) and the trade-offs between Agency A's resources (i.e. total strength) $S_A$ compared to $S_B$ for Bad-Actor-AI system, and Agency A's desired outcomes. Bottom row results plotted as a function of asymmetry between initial strengths for A and B for fixed $S_A + S_B$. **(b)** Risk Chart with plug-and-play formulae that predict key outcome risk measures for containment of B[76]. Right columns give 2 examples using empirical inputs (from Fig. 2 but with $s_{min} = 20$ from simulations, SOM Sec. 5).

Figure 4(a) bottom row considers the more ambitious policy of complete removal of Bad-Actor-AI (B). The on-average stronger Agency (A) cluster finding a B cluster now removes it, e.g. it bans all the B cluster's communities. The time for A to completely remove B given initial strengths $S_A, S_B$, becomes:

$$T = 2S_B + \frac{1}{2}(S_A - S_B) \ln (S_B(S_A - S_B)/S_A) \qquad (2)$$

This reveals a downside to the goal of complete B removal: as A's strength increases, the B clusters become less strong on average and hence less noticeable to A. This creates a rise and peak in the time needed $T$. If B is stronger than A ($S_B > S_A$) then B cannot be removed -- but this problem of large $T$ now becomes a benefit for A since it means an extended time until A's strength is exhausted.

These results and predictions are explicit, quantitative, testable and generalizable, and hence provide a concrete starting point for strengthening Bad-Actor-AI policy discussions. Of course, many features of our analysis could be extended and improved: for example, what happens if future AI can predict the cluster dynamics (ChatGPT currently cannot) and hence Bad-Actor-AI community clusters outwit the containment mechanism? They could also use new decentralized or block-chain platforms as perpetual GPT 'reactor' cores that cannot be contained. Our label 'Bad Actor' should be subclassified (e.g. anti-Semitic vs. anti-women) as should 'vulnerable mainstream community'. We should also account for links from the vulnerable mainstream communities to extreme anti-X communities. Despite this, our main system-level results are robust and their predictions adjusted as Bad-Actor-AI capabilities evolve.

**Supporting Online Material (SOM) contents**:
**1. Methods and empirical details**
        **1.1 Data collection**
        **1.2 Discussion of vulnerable mainstream communities**

**2. Examples of links being created over time from one Bad Actor community to another, intra-platform (i.e. within same platform) and inter- platform (i.e. between different platforms)**

**3. Examples of the location of some topical Bad Actor content across the online battlefield:**
**(1) mass-shooter insignia (e.g. 2023 Texas shooter's chest tattoo) RWDS 'Right Wing Death Squad'; (2) Wagner mercenary- related communities**

**4. Explanation of Venn diagram in Fig. 2 right panel**

**5. Derivation of Eq. 1**

**6. Derivation of Eq. 2**

# References

[1] Europol (2022), Facing reality? Law enforcement and the challenge of deepfakes, an observatory report from the Europol Innovation Lab, Publications Office of the European Union, Luxembourg.  https://idc-a.org/news/industry/AI-Experts-Predict-By-2026-90-Of-Online-Content/127ab0c0-34ba-4c03-8bad-1e4f21923f31#

[2] Hoffman, M., Frase, H. Adding Structure to AI HarmAn Introduction to CSET's AI Harm Framework. July 2023. https://cset.georgetown.edu/publication/adding-structure-to-ai-harm/?utm_source=Center+for+Security+and+Emerging+Technology&utm_campaign=8e10f87116-adding_structure_to_AI_Harm&utm_medium=email&utm_term=0_fcbacf8c3e-8e10f87116-406468358

[3] Milmo, D. and Hern, A. Elections in UK and US at risk from AI-driven disinformation. The Guardian. 20 May 2023. https://www.theguardian.com/technology/2023/may/20/elections-in-uk-and-us-at-risk-from-ai-driven-disinformation-say-experts

[4] Hsu, T. and Myers, S.L. A.I.'s Use in Elections Sets Off a Scramble for Guardrails. New York Times. June 25, 2023. https://www.nytimes.com/2023/06/25/technology/ai-elections-disinformation-guardrails.html

[5] Euchner, T. Coping with the fear of mass shootings. April 13, 2023. https://www.siouxlandproud.com/news/local-news/coping-with-the-fear-of-mass-shootings/

[6] Social Work Today. One-Third of U.S. Adults Say Fear of Mass Shootings Prevents Them From Going to Certain Places or Events. https://www.socialworktoday.com/news/dn_081519.shtml

[7] EU Press Release. Shaping Europe's Digital Future. June 2, 2023. https://digital-strategy.ec.europa.eu/en/news/strengthened-code-practice-disinformation-signatories-identify-ways-step-work-one-year-after-launch

[8] The Artificial Intelligence Act. EU publications 2023. https://artificialintelligenceact.eu

[9] The Digital Services Act: Commission designates first set of Very Large Online Platforms and Search Engines. EU publications. April 23, 2023. https://digital-strategy.ec.europa.eu/en/news/digital-services-act-commission-designates-first-set-very-large-online-platforms-and-search-engines

[10] Benninger, M. 'Fringe' websites radicalized Buffalo shooter, report concludes. https://www.wbng.com https://www.wbng.com/2022/10/18/fringe-websites-radicalized-buffalo-shooter-report-concludes/ (2022)

[11] Fringe Social Media: Are you digging deep enough? https://smiaware.com/blog/fringe-social-media-are-you-digging-deep-enough/ (2019)

[12] Klar, C. M. R. and R. Fringe social networks boosted after mob attack. The Hill https://thehill.com/policy/technology/533919-fringe-social-networks-boosted-after-mob-attack/ (2021)

[13] Hsu, T. News on Fringe Social Sites Draws Limited but Loyal Fans, Report Finds. The New York Times (2022)

[14] Reporter, C. D. N. S. On fringe social media sites, Buffalo mass shooting becomes rallying call for white supremacists. Buffalo News https://buffalonews.com/news/local/on-fringe-social-media-sites-buffalo-mass-shooting-becomes-rallying-call-for-white-supremacists/article_74a55388-f61b-11ec-812a-97d8f2646d45.html (2022)

[15] Fringe social media networks sidestep online content rules. POLITICO https://www.politico.eu/article/fringe-social-media-telegram-extremism-far-right/ (2022)

[16] https://www.science.org/content/article/does-social-media-polarize-voters-unprecedented-experiments-facebook-users-reveal

[17] Wired to Split. Science Special Issue 381, 6656 (2023). https://www.science.org/toc/science/381/6656

[18] Uzogara, E. Democracy Intercepted. Science 381, 6656 (2023). DOI: 10.1126/science.adj7023https://www.science.org/doi/full/10.1126/science.adj7023

[19] Gonzalez-Bailon, S. et al. Science 381, 392 (2023). https://www.science.org/doi/full/10.1126/science.ade7138

[20] Guess, A. et al. How do social media feed algorithms affect attitudes and behavior in an election campaign? Science 381, 6656 (2023). DOI: 10.1126/science.abp9364. https://www.science.org/doi/full/10.1126/science.abp9364

[21] Guess, A. et al. Reshares on social media amplify political news but do not detectably affect beliefs or opinions. Science 381, 404 (2023). https://www.science.org/doi/full/10.1126/science.add8424

[22] Nyhan, B. et al. Like-minded sources on Facebook are prevalent but not polarizing. Nature (2023). https://www.nature.com/articles/s41586-023-06297-w

[23] Auti, N., Ranaware, S., Ghadge, S., Jadhav, R. and Jagtap, P. Social Media based Hate Speech Detection using Machine Learning. IJRASET (2023). DOI: 10.22214/ijraset.2023.52206

[24] Ollagnier, A., Cabrio, E. & Villata, S. Harnessing Bullying Traces to Enhance Bullying Participant Role Identification in Multi-Party Chats. FLAIRS (2023). DOI: 10.32473/flairs.36.133191

[25] Aldreabi, E., Lee, J. & Blackburn, J. Using Deep Learning to Detect Islamophobia on Reddit. FLAIRS (2023). DOI: 10.32473/flairs.36.133324

[26] Morgan, M., Kulkarni, A. Platform-agnostic Model to Detect Sinophobia on Social Media. ACM (2023). DOI: 10.1145/3564746.3587024

[27] Woolley, S., Beacken, G. & Trauthig, I. Platforms' Efforts to Block Antisemitic Content Are Falling Short. Centre for International Governance Innovation https://www.cigionline.org/articles/platforms-efforts-to-block-anti-semitic-content-are-falling-short/ (2022)

[28] Cinelli, M. et al. Dynamics of online hate and misinformation. Sci. Rep. 11, 22083 (2021)

[29] Chen, E., Lerman, K. & Ferrara, E. Tracking Social Media Discourse About the COVID-19 Pandemic: Development of a Public Coronavirus Twitter Data Set. JMIR Public Health Surveill. 6, e19273 (2020)

[30] Gelfand, M. J., Harrington, J. R. & Jackson, J. C. The Strength of Social Norms Across Human Groups. Perspect. Psychol. Sci. J. Assoc. Psychol. Sci. 12, 800–809 (2017)

[31] van der Linden, S., Leiserowitz, A., Rosenthal, S. & Maibach, E. Inoculating the Public against Misinformation about Climate Change. Glob. Chall. 1, 1600008 (2017)

[32] Lewandowsky, S. et al. The Debunking Handbook 2020. (2020). doi:10.17910/b7.1182

[33] Lazer, D. M. J. et al. The science of fake news. Science 359, 1094–1096 (2018)

[34] Smith, R., Cubbon, S. & Wardle, C. Under the surface: Covid-19 vaccine narratives, misinformation and data deficits on social media. First Draft https://firstdraftnews.org:443/long-form-article/under-the-surface-covid-19-vaccine-narratives-misinformation-and-data-deficits-on-social-media/ (2020)

[35] Semenov, A. et al. Exploring Social Media Network Landscape of Post-Soviet Space. IEEE Access 7, 411–426 (2019)

[36] Rao, A., Morstatter, F. & Lerman, K. Partisan asymmetries in exposure to misinformation. Sci. Rep. 12, 15671 (2022)

[37] Wu, X.-Z., Fennell, P. G., Percus, A. G. & Lerman, K. Degree correlations amplify the growth of cascades in networks. Phys. Rev. E 98, 022321 (2018)

[38] Roozenbeek, J., van der Linden, S., Goldberg, B., Rathje, S. & Lewandowsky, S. Psychological inoculation improves resilience against misinformation on social media. Sci. Adv. 8, eabo6254 (2022

[39] Biever, C. ChatGPT broke the Turing test — the race is on for new ways to assess AI. Nature 25 July 2023. https://www.nature.com/articles/d41586-023-02361-7?utm_source=Nature+Briefing&utm_campaign=3428505134-briefing-dy-20230726&utm_medium=email&utm_term=0_c9dfd39373-3428505134-44549989

[40] Miller-Idriss, C. Hate in the Homeland. (2020)

[41] The haters and conspiracy theorists back on Twitter. BBC News. March 8, 2023. https://www.bbc.com/news/technology-64554381

[42] DiResta, R. The Digital Maginot Line. ribbonfarm https://www.ribbonfarm.com/2018/11/28/the-digital-maginot-line/ (2018)

[43] Vesna, C.-G. & Maslo-Čerkić, Š. Hate speech online and the approach of the Council of Europe and the European Union. (2023). https://urn.nsk.hr/urn:nbn:hr:118:377009v

[44] Hate crime: abuse, hate and extremism online - Home Affairs Committee - House of Commons (2017). https://publications.parliament.uk/

[45] Hart, R. White Supremacist Propaganda Hit Record Levels In 2022, ADL Says. Forbes https://www.forbes.com/sites/roberthart/2023/03/09/white-supremacist-propaganda-hit-record-levels-in-2022-adl-says/ (2023)

[46] Online Hate and Harassment: The American Experience 2023. June 27, 2023. ADL Center for Technology and Society. https://www.adl.org/resources/report/online-hate-and-harassment-american-experience-2023

[47] Eisenstat, Y. Hate is surging online — and social media companies are in denial. The Hill. July 7, 2023. https://thehill.com/opinion/congress-blog/4085909-hate-is-surging-online-and-social-media-companies-are-in-denial-congress-can-help-protect-users/

[48] Nelson, J. UN Warns of AI-Generated Deepfakes Fueling Hate and Misinformation Online. Decrypt. June 12, 2023. https://decrypt.co/144281/un-united-nations-ai-deepfakes-hate-misinformation

[49] United Nations. Common Agenda Policy Brief: Information Integrity on Digital Platforms. June 2023. https://www.un.org/sites/un2.un.org/files/our-common-agenda-policy-brief-information-integrity-en.pdf

[50] Brown, R. & Livingston, L. A New Approach to Assessing the Role of Technology in Spurring and Mitigating Conflict: Evidence From Research and Practice. JIA SIPA https://jia.sipa.columbia.edu/new-approach-assessing-role-technology-spurring-and-mitigating-conflict-evidence-research-and (2018)

[51] Starbird, K. Disinformation's spread: bots, trolls and all of us. Nature 571, 449–449 (2019)

[52] Lamensch, M. To Eliminate Violence Against Women, We Must Take the Fight to Online Spaces. Centre for International Governance Innovation https://www.cigionline.org/articles/to-eliminate-violence-against-women-we-must-take-the-fight-to-online-spaces/ (2022)

[53] Crawford, A. and Smith, T. Illegal trade in AI child sex abuse images exposed. BBC News. June 28, 2023. https://www.bbc.co.uk/news/uk-65932372

[54] Cosoleto, T. Surge in young children being targeted by cyber bullies. The West Australian. 9 July, 2023. https://thewest.com.au/news/social/surge-in-young-children-being-targeted-by-cyber-bullies-c-11223220

[55] Gill, P. & Corner, E. Lone-Actor Terrorist Use of the Internet & Behavioural Correlates. in (2015)

[56] Douek, E. Content Moderation as Systems Thinking. Harvard Law Review https://harvardlawreview.org/2022/12/content-moderation-as-systems-thinking/ (2022)

[57] Daily listings. https://www.disinfodocket.com/dd-12jul23/

[58] Green, Y. et al. Evidence-Based Misinformation Interventions: Challenges and Opportunities for Measurement and Collaboration. Carnegie Endowment for International Peace https://carnegieendowment.org/2023/01/09/evidence-based-misinformation-interventions-challenges-and-opportunities-for-measurement-and-collaboration-pub-88661 (2023)

[59] Dynamic Online Networks Laboratory. Literature Review. https://bpb-us-e1.wpmucdn.com/blogs.gwu.edu/dist/5/3446/files/2022/10/lit_review.pdf

[60] DisinfoDocket website and publications. https://www.disinfodocket.com

[61] Manrique, P.D, Huo, F.Y., El Oud, S., Zheng, M., Illari, L., Johnson, N.F. Shockwavelike Behavior across Social Media. Phys. Rev. Lett. 130, 237401 (2023)

[62] Lupu, Y. et al. Offline events and online hate. PLOS ONE 18, e0278511 (2023)

[63] Velásquez, N. et al. Online hate network spreads malicious COVID-19 content outside the control of individual social media platforms. Sci. Rep. 11, 11549 (2021).

[64] Ammari, T. & Schoenebeck, S. 'Thanks for your interest in our Facebook group, but it's only for dads': Social Roles of Stay-at-Home Dads. in Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing 1363–1375 (Association for Computing Machinery, 2016). doi:10.1145/2818048.2819927

[65] Moon, R. Y., Mathews, A., Oden, R. & Carlin, R. Mothers' Perceptions of the Internet and Social Media as Sources of Parenting and Health Information: Qualitative Study. J. Med. Internet Res. 21, e14289 (2019)

[66] Laws, R. et al. Differences Between Mothers and Fathers of Young Children in Their Use of the Internet to Support Healthy Family Lifestyle Behaviors: Cross-Sectional Study. J. Med. Internet Res. 21, e11454 (2019)

[67] Madhusoodanan, J. Safe space: online groups lift up women in tech. Nature 611, 839–841 (2022)

[68] Jacomy, M., Venturini, T., Heymann, S. & Bastian, M. ForceAtlas2, a Continuous Graph Layout Algorithm for Handy Network Visualization Designed for the Gephi Software. PLOS ONE 9, e98679 (2014)

[69] Leahy, R., Restrepo, N.J., Sear, R., Johnson, N.F. Connectivity Between Russian Information Sources and Extremist Communities Across Social Media Platforms. Front. Polit. Sci., 22 June 2022

[70] Johnson, N.F. Simple mathematical law benchmarks human confrontations. Scientific Reports 3, 3463 (2013). DOI: 10.1038/srep03463

[71] Johnson, N.F. et al. Pattern in Escalations in Insurgent and Terrorist Activity Science 333, 81 (2011) DOI: 10.1126/science.1205068

[72] Data from NANEX https://www.nxcoredata.com

[73] Data from MANDIANT https://www.mandiant.com

[74] Dixon, A., Zhao, Z., Bohorquez, J.C., Denney, R., Johnson, N.F. Statistical physics and modern human warfare. In G. Naldi et al. (eds.), Mathematical Modeling of Collective Behavior in Socio-Economic and Life Sciences, Modeling and Simulation in Science, Engineering and Technology, DOI 10.1007/978-0-8176-4946-3-14, p. 365 and references therein

[75] Zhao, Z., Bohorquez, J.C., Dixon, A., Johnson, N.F.: Anomalously Slow Attrition Times for Asymmetric Populations with Internal Group Dynamics. Phys. Rev. Lett. 103, 148701 (2009)

[76] Newman, M.E.J. Power laws, Pareto distributions and Zipf's law. Contemporary Physics 46, 323 (2005)