

Combining interventions to reduce the spread of viral misinformation

Joseph B. Bak-Coleman^{a,b,c}, Ian Kennedy^{a,d}, Morgan Wack^{a,e}, Andrew Beers^{a,f}, Joey Shafer^a, Emma S. Spiro^{a,c,d}, Kate Starbird^{a,f}, and Jevin West^{a,c}

^aCenter for an Informed Public, University of Washington, Seattle, WA 98195

^beScience Institute, University of Washington, Seattle, WA 98195

^cThe Information School, University of Washington, Seattle, WA 98195

^dDepartment of Sociology, University of Washington, Seattle, WA 98195

^eDepartment of Political Science, University of Washington, Seattle, WA 98195

^fHuman Centered Design and Engineering, University of Washington, Seattle, WA 98195

ABSTRACT

Mis- and disinformation online pose a range of threats, from subverting democratic processes to undermining public health measures^{1–3}. Proposed solutions range from encouraging more selective sharing by individuals, to platform removal of false content and accounts that create or promote it^{4,5}. How, whether, and which strategies to implement depends on their relative and combined ability to reduce viral misinformation spread at plausible levels of enforcement. Here we provide a framework to evaluate interventions aimed at reducing viral misinformation online both in isolation and when used in combination. We begin by deriving a generative model of viral misinformation spread, inspired by research on infectious disease. Applying this model to a large corpus of misinformation events that occurred during the 2020 US election, we reveal that commonly proposed interventions—including removal of content, virality circuit breakers, nudges, and account banning—are unlikely to be effective in isolation without extreme censorship. However, our framework demonstrates that a combined approach can achieve a substantial ($\approx 50\%$) reduction in the prevalence of misinformation. Our results challenge claims that combating misinformation will require new ideas or high costs to user expression. Instead, we highlight a practical path forward as misinformation online continues to threaten vaccination efforts, equity, and democratic processes around the globe.

1 Introduction

2 Misinformation—i.e. false information—has become a pervasive
3 feature of online discourse, resulting in increased belief in
4 conspiracy theories, rejection of recommended public health
5 interventions, and even genocide^{1,3,6,7}. Academics and those
6 working in industry have proposed a host of potential solutions,
7 ranging from techniques for detecting and removing
8 misinformation to empowering users to be more discerning in
9 their sharing habits^{4,5,8}. Despite an abundance of proposed
10 interventions, misinformation remains a pervasive feature of
11 digital life^{1,2,9}.

12 During the election, the spread of false information was
13 often characterized by brief (i.e. hours, days) periods of
14 rapid growth in discussion and sharing³. During these events,
15 engagement (i.e. all discussion and sharing) exhibits viral,
16 disease-like dynamics—self-replicating, endogenous growth
17 stemming from a limited number of initial sources^{10,11}. Some
18 of these incidents quickly died out, while others had multiple
19 waves, spread to other platforms, and often became consolidated
20 into broader narratives. Reducing early virality provides
21 a source of promise for successful interventions, as disrupting
22 viral spread may have cascading effects on narrative consolidation
23 and future engagement. Unfortunately, the rapid growth inherent to viral misinformation makes it challenging
24 to respond to in a timely manner.

25 Of urgent need is quantitative comparison between proposed interventions' ability to reduce the spread of misinfor-

mation. Crude approaches like outright removal and banning of either content or accounts will certainly work if applied in excess, yet come with costs to freedom of expression and force private entities to be arbiters of truth. For judicious use, questions arise about how soon and how much removal is necessary for a meaningful effect. Interventions that rely on empowering individuals to consume and share more discerningly have shown promise in experimental contexts, but it remains unclear what impact they will have at scale^{4,4}.

Beyond comparison, we lack an understanding of when—and indeed whether—multiple interventions can act synergistically to reduce the spread of misinformation. Toward this goal, we derive and parameterize a generative model of misinformation engagement (i.e. total discussion and sharing of posts related to false information) using a large corpus of tweets collected during the 2020 election in the US³. We rely on this model to examine the efficacy of misinformation interventions both in isolation and when deployed in combination. Finally, we examine how the spread of misinformation during viral periods impacts subsequent engagement.

Results

Overview and Model

We begin by deriving a generative model of misinformation spread that relies on a few simple assumptions. First, that a users' audience can be approximated by their follower count^{12,13}. Next, that the spread of misinformation resem-

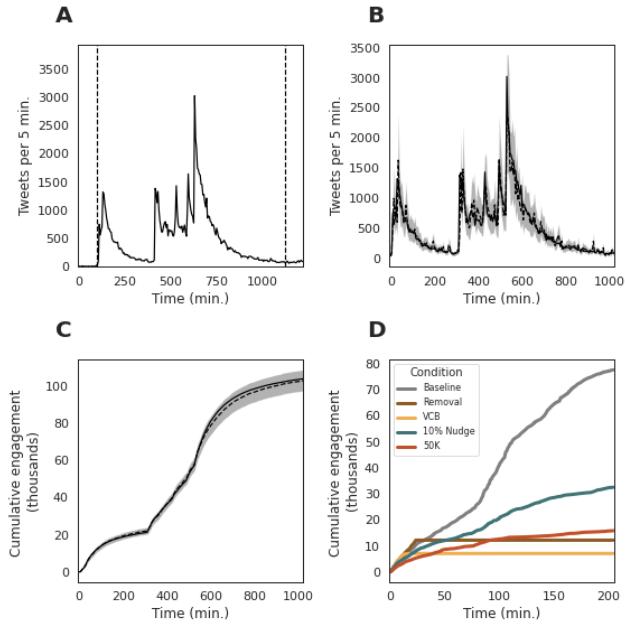


Figure 1. Overview of our model and analysis **A** An event segmented from a larger incident (dashed lines). **B** Time-series for a single event, here a later-recanted story about a poll worker in Pennsylvania admitting to ballot tampering. Dashed-line: expected value, shaded region: 89% credible interval (CI). **C** Cumulative engagement as a measure of total misinformation, lines and shading are as in B. **D** Model-simulated platform interventions for a single event. Lines indicate median cumulative engagement over 100 simulations. Grey: baseline, purple: 10% "nudge", orange: banning, yellow: virality circuit breaker, green: outright removal of content

bles a simple contagion whereby only a single interaction is required for transmission^{14,15}. We further assume that spread during a viral event occurs endogenously within a single platform rather than via other processes (e.g cable news, cross-platform). Finally, we assume that discussions on a topic decay over time as new topics replace them and reach saturation. These phenomena can be captured by a minimally parameterized branching process model, such that:

$$\begin{aligned} \mathbb{E}[y_t] &= \exp(\alpha + \beta v_{t-1}) \\ v_t &= v_{t-1} \delta e^{-\lambda t} + x_t \\ x_t &= \log\left(\sum_{j=1}^{y_t} F_{j,t}\right) \end{aligned} \quad (1)$$

Where y_t are the posts (i.e. retweets, tweets, replies, quote tweets) at time t , α is the baseline rate of discussion, and β is the effect of virality, v . Virality is a proxy for the total number of users at a given point in time that are exposed to and may propagate misinformation. Virality decays as an exponential function via δ and λ . Here, δ captures the baseline rate of

decay per time step, and λ controls the way in which that decay changes over the lifetime of an event. This could either be due to algorithmic processes favoring new content or, for very large events, user saturation. Every time step, for each of y_t accounts, the log sum (x_t) of their followers, F_j is added to virality for the subsequent time step. Our model bears similarity to those used to evaluate interventions and superspreading in infectious disease¹⁶.

Rather than solving this model analytically, we instead rely on a computational approach with parameters estimated from 216 events ($\approx 6M$ Posts) of rapid misinformation spread observed online during the 2020 US election (See Methods). This allows us to draw from an empirical distribution of follower counts specific to a given event and study the effect on engagement of banning users in a manner that is conditioned on their real-world patterns of behavior. We estimated the parameters of the model for each event using Bayesian inference, generating posterior predictive time-series for evaluating fit for each event (Fig 1B-C, Methods). Using the data-derived parameters, we simulated the impact of platform interventions on cumulative engagement across all events (See Methods, Fig. 1D). Posterior predictive plots for all 216 events are presented in the SI.

Fact-Checking and time-lagged approaches

We begin by considering the impact on user engagement of approaches in which a platform applies policies that target a specific instance of misleading or false information—in this case individual posts. Among the more commonly employed strategies during the 2020 US election was to identify specific misinformation and take action, ranging from applying a label to outright removal³. These approaches share a common feature of requiring time before action is taken. Time is necessary to not only identify the misinformation, but also to decide on an appropriate response.

In an extreme case, a platform could remove or hide all content matching search terms related to an emerging misinformation incident. To simulate this, we ran our model until time t , at which point growth stopped entirely (Fig 2A). Our results indicate that outright removal can indeed be effective, producing a dramatic 93.5% median reduction in total posts (i.e. tweets, replies, quote-tweets, and retweets) on the topic if implemented within 30 minutes (89% C.I. [92.4, 94.2]). Even with a 4 hour delay, our model indicates reductions of 50.5% (89% C.I. [47.6, 52.9], Table S1). These effects generously assume that platforms are able to monitor, detect, sufficiently fact-check (for ethical considerations) and implement a full removal response within the specified time-frame. As such, the efficacy is dramatically reduced if only a fraction of events lead to action (Fig 2B, Table S2).

A more plausible approach could involve "virality circuit breakers" which seek to reduce the spread of a trending misinformation topic without explicitly removing content, for example by suspending algorithmic amplification¹⁷. As this approach is less challenging from a platform public relations

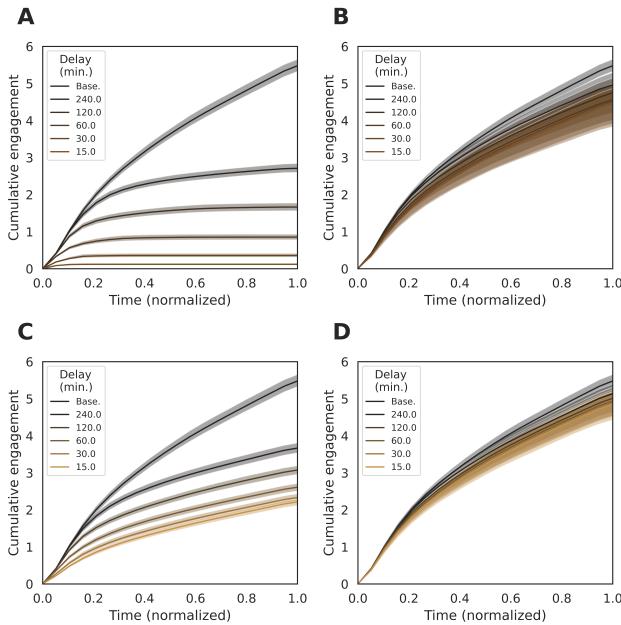


Figure 2. **A** The impact of outright removal on of all misinformation-related posts following a delay specified in minutes. Time (x-axis) is normalized to the duration of the event. **B** As in A if only 20% of events are removed. **C** The impact of applying a virality circuit breaker that reduces virality by 10% to all misinformation events after a specified period of time. **D** As in C, if the VCB is applied to only 20% of events.

and ethics standpoint than outright removal, a lower threshold for fact-checking would enable quicker response times. We simulate the impact of virality circuit breakers by reducing virality after a period of time such that $\hat{v}_t = v_t * (1 - p)$, where p is the proportional reduction in virality (See [Methods](#)).

Through simulations, we reveal how virality circuit breakers can have similar efficacy to outright removal even if the amount by which virality is reduced is small (Fig. 2B, Table S3). For instance, a 10% reduction in virality, implemented four hours after the start of an event, can reduce the spread of misinformation by nearly 33.0% (89% C.I. [29.7 36.3]). As with outright removal, however, the efficacy is primarily limited by the proportion of events for which the platforms take action (Fig 2D, Table S2)

136 Nudges and reduced reach

137 A drawback of fact-checking based approaches is that they
138 are most applicable to transparently false or readily falsifiable
139 claims³. Many instances of misinformation involve claims
140 that are either partly true, or require non-trivial time to debunk.
141 Depending on implementation, time-lagged responses may
142 further require that users are receptive to the intervention, or
143 do not find ways around removal or platform action.

144 These challenges motivate approaches that leverage individual discretion to reduce the spread of misinformation⁵. For

instance, encouraging users to consider accuracy has been shown to reduce rate at which individuals share misinformation by 10-20%⁴. AI based approaches could likewise scan drafts of posts and warn users if they appear to be amplifying misinformation.

A central question is whether a modest reduction in individual sharing behavior can lead to a more dramatic change in overall rates of misinformation. Agent-based models support this notion across a range of network topologies⁴. From the perspective of our model, nudge-based approaches can be simulated by maintaining the parameters from the initial model fit while proportionally reducing the following of every user that discusses an incident. Recall [Eqn. 1](#): from the perspective of the model a nudge, η , can be implemented such that $\hat{F}_{j,t} = (1 - \eta)F_{j,t}$ where $0 < \eta < 1$.

Using our model to simulate nudges, we find that they can indeed reduce the prevalence of misinformation (Fig 3A, Table S5). Nudges that reduce sharing by 5, 10, 20, and 40% result in a 6.6, 12.4, 22.6, and 38.9% reduction in cumulative engagement, respectively (Table S5). The median effect tends to be larger than the nudge suggesting a degree of feedback whereby the individual effect of a nudge is compounded in the misinformation dynamics.

169 Account Banning

In our dataset, several accounts shared or amplified misinformation across multiple incidents³. Moreover, some of these repeat offenders had out-sized audiences when compared to the average Twitter user—ranging hundreds of thousands to millions of followers. While removal of repeat offenders during the election was rare, several were removed following the violent insurrection at the US capital on January 6, 2021. A question remains whether the removal of these accounts, or account-focused policies in general, will have a meaningful impact on misinformation. While large-followings often confer engagement, it remains possible that there is sufficient sharing from smaller accounts to ensure the spread of misinformation even in the absence of the larger audience removed accounts¹⁸.

One challenge in modeling account removal is that there likely exists non-trivial relationships between account-size, propensity to share misinformation, and the timing at which certain accounts amplify narratives. A large account that regularly shares misinformation in the first five minutes will have an out-sized effect compared to a smaller account that occasionally shares misinformation hours later. To account for this, our model samples from the empirical follower-count distribution in a given time step. Further, as identities of individuals are known, we can remove specific accounts and simulate total engagement (See [Methods](#)). In other words, our simulations are conditioned on unseen patterns of, and variation in, individual behavior without explicitly quantifying the differences in individual behavior. Through this, our model and simulations exhibit robustness to considerable unmeasured real-world complexity.

We begin by considering the consequences of account removals ($N = 1504$) that occurred in early 2021. We seek to answer whether previously implemented account removal is sufficient to curb misinformation going forward. Our simulations reveal that the removal of these accounts from our dataset reduces total engagement with misinformation by 12.0%, (89% C.I. [8.4, 15.8], Fig 3C). This is comparable in efficacy to a 10% reduction in sharing of misinformation (i.e. a nudge) impacting all accounts in the absence of removal.

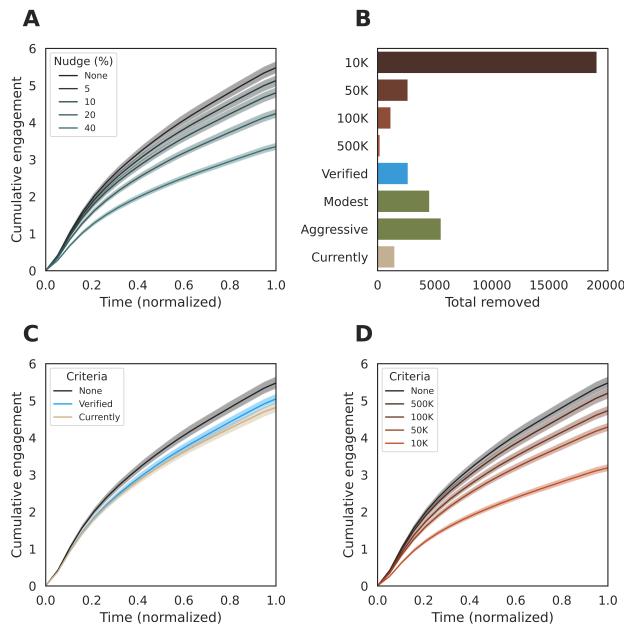


Figure 3. **A** The effect of nudges that inoculate a percentage of the population against spreading misinformation. Shown is the cumulative total engagement across all events, with time normalized to the duration of the event. **B** Number of accounts that are either currently removed, or would have been removed under a three-strikes policy **C** The effect of accounts removal for either those that are currently banned (orange) or those banned following a three strikes rule applied solely to verified accounts (blue). **D** As in A and C, yet showing the impact of enacting three-strikes policies with varying thresholds.

We next consider a "3 Strikes" rule in which accounts are removed from the platform after they are detected in three distinct incidents of misinformation (i.e. topics, regardless of number of posts for a given topic). For these simulations, any interaction or amplification of misinformation (i.e. tweets, retweets, quote tweets) would be counted as a strike. A policy focused solely on original content could be gamed by using large accounts to amplify disposable smaller accounts. This type of policy would avoid banning accounts that were swept up by a given piece of misinformation and tweeted repeatedly, while focusing on those that spread misinformation more broadly. Applied solely to verified accounts, we observe a 7.8% drop in cumulative engagement (89% C.I. [4.1, 11.7])

which likewise is similar in efficacy to a nudge rolled out across the board (Table S6, Fig. 3C). If, instead of verification, the policy is applied based on the number of followers an account has, pronounced effects are only observed when the threshold is quite low ($\approx 10,000$ followers) requiring large numbers of accounts to be removed (Fig. 3B and D, Table S7)

Combined Approaches

All of the approaches above exhibit some efficacy reducing engagement with viral misinformation. Unfortunately, each strategy tends to become maximally effective in impractical regions of parameter space. Outright removal of misinformation is particularly effective, yet it is difficult to imagine that more than a small fraction of misinformation can be easily removed. Virality circuit breakers face similar challenges, albeit to a lesser extent. For nudges that minimally impact user experience yet improve individual discretion, effects far beyond $\approx 20\%$ are unlikely without a major breakthrough in information literacy or social psychology⁴. In the case of banning specific accounts, low follower thresholds increase the number of accounts removed, and thus costs and challenges, super-linearly.

We therefore consider a combined approach relying on only modest implementations of each of the strategies studied above. Specifically, viral circuit breakers are employed for 5% of content, reducing virality (v_t) by 10%, and enacted after 120 minutes (i.e. for ≈ 11 events). Among the content that is subjected to a viral circuit breaker, 20% is subsequently removed outright after four hours (≈ 2 events). We further assume a 10% reduction in individual sharing of misinformation resulting from a nudge. Finally, accounts that have been removed remain banned, and a 3-strokes policy is applied to verified accounts and those with more than 100K followers. Our model reveals that even a modest combined approach can result in a 36.3% (89% C.I. [31.4, 41.8]) reduction in the total volume of misinformation (Figure 4A, Table S8)

We additionally consider a more aggressive version of a combined policy, applying viral circuit breakers to 10% of content and reducing virality by 20% while cutting response times in half. We further assume a 20% nudge, and reducing the threshold for the 3-strokes policy to 50K followers. This more aggressive approach dramatically reduced misinformation by 48.8% (89% C.I. [43.4, 54.9], Figure 4D, Table S9). Similar efficacy from standalone approaches would either be impossible ($> 40\%$ Nudge) or require draconian removal of content and accounts.

One limitation of our model is that it relies on assumptions specific to periods of viral misinformation spread. In our dataset, only 40% of posts occur during the largest event for a given incident. Yet 48% of engagement occurs after the largest event. While our model cannot provide direct insight into how interventions will impact engagement during these periods, we can gain indirect insight by considering the relationship between the size of an event and subsequent discussion.

Our data demonstrate that the size of an event is strongly

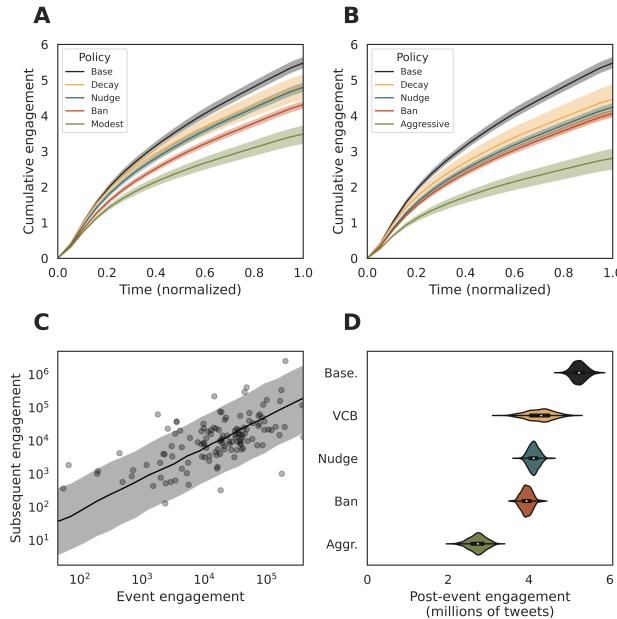


Figure 4. **A** The impact of a modest combined approach to intervention (described in text, green) and each intervention applied individually (as per legend) **B** The impact of a more aggressive combined approach (described in text, green) and each intervention applied individually (as per legend) **C** Relationship between engagement within the largest viral event for a given incident and subsequent engagement **D** Expected post-event engagement given action taken during an event.

predictive of subsequent engagement (Fig. 4C, Bayesian Log-Normal Regression $\beta = .95$, 89% C.I. [93, 97], Table S10). Using this relationship, we can estimate subsequent discussion based on simulated, intervention-adjusted engagement during the largest event (See Methods). Doing so, we reveal that impacts on subsequent discussion are likely to be similar in magnitude to the impacts of the intervention during an event (Fig. 4D, S7).

Conclusions

Our derived model, grounded in data, provides quantitative insight into the relative efficacy of proposed interventions. Through simulation, we reveal that proposed interventions are unlikely to be effective if implemented individually at plausible levels. Effective removal of content or virality circuit breakers would require large teams, rapid turn-around times, and place content decisions squarely in the hands of private organizations. Nudges are promising but unlikely to be a panacea at known levels of efficacy⁴. Banning is the most plausible solution, but would require draconian removal of tens of thousands of users.

Fortunately, our results show that combining interventions at plausible levels of enforcement can be effective at reduc-

ing misinformation. While it is unsurprising that multiple interventions outperform individual approaches, our paper provides necessary insight into the magnitude of that difference. The efficacy of a combined approach is dependent not only on the nature of individual interventions, but how they interact with one another, the dynamics of misinformation spread, event duration, user sharing behavior, user follower counts, and how these factors change throughout the course of a disinformation campaign. In fitting our model to a large corpus of events during an active period of mis- and disinformation, our results are conditioned on much of this complexity. Further, by drawing from the empirical distribution of users' follower counts, our model indirectly and implicitly accounts for unseen behavioral patterns of users and changes to their follower counts over time.

What remains unclear is how changes in the magnitude of events will impact longer-term dynamics of misinformation and translate to a reduction in harm. If implemented in tandem, it may prove a sufficient shock to collapse the misinformation ecosystem altogether, as shock-induced collapse is a central feature of complex systems¹⁹. For instance, subsequent events likely depend on the size of previous events, and breaking that feedback could lead to greater than expected gains. However, this same body of literature suggests that an insufficient shock may yield only short-term changes as the system re-organizes and adapts. For this reason, rolling out policies individually and insufficiently may make the problem harder to solve in the long term.

We note that the results presented here rely on a simplified model of events on a single platform in what is a highly complex, multi-platform system. These types of simplifications are an inherent limitation of any approach, short of risky, large-scale experimentation. However, abstract models of complex processes have proved essential to predicting the benefits of interventions on complex systems, from the mitigating the spread of disease to stabilizing ecosystems^{16,20}. Models provide particular utility when experiments are unethical and impractical, and costs of inaction are high. Given the substantial risks posed by misinformation in the near term, we urgently need a path forward that goes beyond trial and error or inaction. Our framework highlights one such path that can be adopted in the near-term without requiring large-scale censorship or major advances in cognitive psychology and machine learning.

1 Methods

1.1 Data Collection and Processing

Our dataset was collected in real-time during the 2020 US election. We relied on an evolving set of keywords to collect data from Twitter's API. Keywords were updated in response to new narratives, for instance adding "sharpiegate" and related terms after false narratives emerged about the use of Sharpie markers invalidating ballots. Working with the Electoral Integrity Partnership, we catalogued instances of misinformation that were either detected by the team or

298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341

352 reported by external partners³. This led to a large corpus of
353 tickets associated with validated reports of misleading, viral
354 information about election integrity.

355 Tickets that shared a common theme were consolidated
356 into incidents. For each incident, we developed search terms
357 and a relevant date range in order to query posts from our
358 tweet database. Incidents ($N_i = 153$) were generally charac-
359 terized by one or more periods of intense activity followed
360 by returning to a baseline state (Fig 1A). Search terms and
361 descriptions of incidents are provided along with the data.
362 The time-series of all related posts from each incident was
363 grouped into five minute intervals and segmented into distinct
364 periods of increased activity, events ($N_E = 216$, See [Methods](#),
365 Fig 1).

366 **Event segmentation and inclusion criteria**

367 Each of the 154 incidents misinformation was characterized
368 by one or more periods of viral spread (i.e. events). Long
369 periods of low activity between events would violate the as-
370 sumptions of our statistical model, requiring segmentation
371 of incidents into discrete events. We began segmenting by
372 grouping collected posts into 5 minute intervals, and finding
373 the interval within the aggregated time-series that had the
374 highest volume of collected posts. Other peaks in activity
375 were considered part of separate events if they were at least
376 30% of the magnitude of the largest peak (to filter out noise).
377 Event boundaries were determined as the points before and
378 after the peak where the number of posts in 5 minutes was less
379 than 5% of the maximum volume. If this did not occur within
380 the range of data collection, the first (or last) time-point col-
381 lected was used to denote the beginning (or end) of an event.
382 Finally, events were required to last at least an hour (i.e. 12
383 data points)

384 Using this initial corpus of 260 events, our model was fit to
385 each event using PyStan^{21,22}. We fit events separately (rather
386 than hierarchically) as they varied widely in their time scales,
387 magnitudes, and context within the broader 2020 election cy-
388 cle. These factors, combined with computational limitations,
389 precluded a full hierarchical model from being feasible or
390 appropriate. Similarly, our model was unlikely to be appro-
391 priate for all events as it makes assumptions post volume is
392 well predicted by the number of previously exposed accounts
393 on twitter. If, for instance, an incident received substantial
394 news coverage (i.e. Dominion software narratives) our model
395 would likely fail.

396 To safeguard against this, we relied on a number of criteria
397 to ensure model fit to a given event. Events were included
398 in the final analysis if a) the posterior 89% C.I. of total posts
399 contained the observed value and b) the chains successfully
400 converged for all parameters ($\hat{R} < 1.1$) c) The fit did not con-
401 tain divergent transitions and d) the event lasted longer than
402 an hour (i.e. > 12 data points to fit). Other than these criteria,
403 events surrounding the dominion narrative were removed as
404 they involved long periods high volume online discussion.
405 This filtering processes resulted in inclusion of 216 events
406 (81% of total events), and $\approx 6M$ posts. We note that, in the

407 main text, we consistently see a slightly smaller ($\approx 5\%$) num-
408 ber of cumulative posts in our baseline condition. We suspect
409 this is due to non-randomness in the relationship between fol-
410 lower counts and the probability of being involved in a tweet
411 (we sample randomly) and/or the absorbing boundary of zero
412 posts in our model.

413 **Statistical and Computational Model**

414 **1.2 Model Justification**

415 The spread of misinformation online occurs on complex net-
416 works involving aspects of both organic growth and coor-
417 dinated disinformation campaigns. Acceptance of a given
418 misinformation narrative likewise involve a complicated cog-
419 nitive process involving partisan leanings, prior knowledge,
420 attention, the message content, and a host of other factors^{4,23}.
421 At face value, it would appear unlikely that a minimally pa-
422 rameterized model could adequately capture the generative
423 process and provide useful insight. Yet, a similar argument
424 could be invoked regarding the spread of disease which in-
425 volves non-trivial behavioral, fluid, and immune dynamics.
426 Nevertheless, compartment models (e.g. SIR, SEIR) have
427 become essential epidemiological tools in the century since
428 their introduction^{16,24,25}.

429 Models of complex process provide useful insight when
430 they capture the leading-order terms drive a system's dynam-
431 ics²⁶. Here we assume that the dynamics are driven primarily
432 by the number of people previously exposed and declining en-
433 gagement through saturation or replacement with new content.
434 The ability of our model to recreate patterns of engagement
435 provides indirect evidence that it captures key phenomena
436 (See [S1](#)). Finally, we note that previous work has lever-
437 aged epidemiological models to understand the spread of
438 viral memes¹¹.

439 **Statistical Model**

440 We model the growth of misinformation as a branching pro-
441 cess in which posts (and thus virality) in subsequent time steps
442 is a function of activity in previous time steps. Posts y_t at time
443 t are assumed to be distributed as a gamma-poisson mixture
444 (i.e. negative-binomial) with expected value μ_t . Specifically:

$$y_t \sim \text{NegativeBinomial2}(\mu_t, \phi) \text{ for } t = 2 \dots T$$

$$\mu_t = \exp(\alpha + \beta v_{t-1}) \text{ for } t = 2 \dots T$$

$$v_t = v_{t-1} \delta e^{-\lambda t} + x_t$$

$$\alpha \sim \text{Normal}(-3, 3)$$

$$\beta \sim \text{Normal}(0, 3)$$

$$\delta \sim \text{Beta}(2, 2)$$

$$\lambda \sim \text{HalfExponential}(1)$$

$$\phi \sim \text{HalfExponential}(1)$$

$$v_1 = x_1$$

$$x_t = \log\left(\sum_{j=1}^{y_t} F_j + 1\right)$$

445 Where α is the baseline rate of detection for related key-
446 words and β is the effect of virality, v . Virality is calculated as
447 a decaying function of v_{t-1} and the log of the sum of account
448 follower counts F_j for posts in the previous time step. A value
449 of one is to avoid an undefined value in time steps with no
450 followers. The log transform accounts for the link function
451 (\exp) transforming the linear model into an expected value
452 for the Negative Binomial distribution. Given the wide range
453 of possible event shapes, weakly generic, weakly informative
454 priors were chosen for all parameters. Models were fit using
455 HMC in Pystan with default sampling parameters^{21,22}.

456 Computational Model

457 Our computational model relied on the posterior distributions
458 of parameters obtained from fitting our statistical model sep-
459 arately to each event. For each simulation, one sample was
460 drawn at random from the posterior for a given event. At $t = 1$,
461 the model was initialized with the volume of posts and total
462 engagement from the first time step in which any posts were
463 observed. At each subsequent time step, our computational
464 model predicted the number of new posts, y_t , by sampling
465 from a negative binomial distribution as per our statistical
466 model. For each of y_t new posts, we drew a follower count
467 from the actual distribution of accounts that retweeted for that
468 event, at that time step. Doing so allowed us to control for
469 the possibility that some accounts tend to appear earlier in a
470 viral event. This processes was repeated for the duration of
471 the actual event.

472 We simulated removal of misinformation by simply setting
473 $y_{t+1} = 0$ after at a specified intervention time, t . Virality
474 circuit breakers were enacted by multiplying virality at each
475 time step by a constant. For example, a 10% reduction in
476 virality was implemented as $\hat{v}_t = v_t(1 - .1)$. As with the
477 removal, this occurred only after a given time step. In the
478 case of the combined approach, virality circuit breakers (and
479 subsequent removal) were employed at a given probability for
480 each run of the simulation. Nudges were implemented through
481 multiplying follower counts by a constant, reducing the pool
482 of susceptible accounts (i.e. for account j , $\hat{F}_j = F_j(1 - \eta)$).
483 Finally, we implemented a 3-strikes rule by identifying the
484 fourth incident and all subsequent incidents in which a given
485 account appeared in our full dataset. Their follower count was
486 removed from all subsequent simulations.

487 Additionally, our model included a maximum value of
488 twice the observed posts per time interval to account for a
489 rare condition in which long-tail parameters would lead to
490 runaway. This was observed to occur rarely enough to be
491 challenging to quantify (< 1% of model runs), but was im-
492 plemented to reduce upward bias in control conditions. This
493 was done to ensure conservative estimates of efficacy as in-
494 terventions could reduce the possibility for runaway without
495 meaningfully impacting engagement. Such a feature would
496 be expected in any model of a growth process with pareto-like
497 distributions of follower counts and spread at a given time
498 step (i.e. a negative binomial).

499 For the figures show in the main text, and the tables pre-

500 sented in the SI, 500 simulations of all 220 events (110000
501 simulations per condition). For each run, we computed the
502 cumulative engagement, and normalized it across time using
503 linear interpolation to 20 time steps. The 100 simulations were
504 summed across runs, from which we computed the median
505 and credible intervals. All simulations were done in Python.
506

507 Acknowledgements

508 This work was made possible through the generous support
509 from the Knight Foundation, the University of Washington
510 eScience Institute, and Craig Newmark Philanthropies. We
511 also thank our collaborators, the Stanford Internet Observatory,
512 Graphika, DFRLab, and the Electoral Integrity Partnership.
513 We further thank Carl Bergstrom, Iain Couzin, Fernando
514 Rossine, Rachel Moran, and Kolina Koltai for their feedback.
515

516 Author contributions statement

517 J.B-C., A.B., and J.W. conceived of the study. J.S., A.B.,
518 M.W., I.K., E.S., and K.S., developed the dataset. J.B-C. and
519 J.W. wrote the model and simulation code. J.B-C. drafted the
initial manuscript and all authors were involved in subsequent
revision.

520 References

1. Velásquez, N. *et al.* Online hate network spreads malicious COVID-19 content outside the control of individual social media platforms. Tech. Rep. (2020). DOI: [10.21203/rs.3.rs-110371/v1](https://doi.org/10.21203/rs.3.rs-110371/v1).
2. Pennycook, G. & Rand, D. G. Research note: Examining false beliefs about voter fraud in the wake of the 2020 Presidential Election. *Harv. Kennedy Sch. Misinformation Rev.* DOI: [10.37016/mr-2020-51](https://doi.org/10.37016/mr-2020-51) (2021).
3. Election Integrity Partnership. The Long Fuse: Misinformation and the 2020 Election | Stanford Digital Repository. Tech. Rep., Center for an Informed Public, Digital Forensic Research Lab, Graphika, & Stanford Internet Observatory, Stanford Digital Repository (2021).
4. Pennycook, G. *et al.* Shifting attention to accuracy can reduce misinformation online. *Nature* 1–6, DOI: [10.1038/s41586-021-03344-2](https://doi.org/10.1038/s41586-021-03344-2) (2021).
5. Lazer, D. M. J. *et al.* The science of fake news. *Sci. (New York, N.Y.)* **359**, 1094–1096, DOI: [10.1126/science.aao2998](https://doi.org/10.1126/science.aao2998) (2018).
6. Whitten-Woodring, J., Kleinberg, M. S., Thawngmung, A. & Thitsar, M. T. Poison If You Don't Know How to Use It: Facebook, Democracy, and Human Rights in Myanmar. *The Int. J. Press.* **25**, 407–425, DOI: [10.1177/1940161220919666](https://doi.org/10.1177/1940161220919666) (2020).
7. Koltai, K. VACCINE INFORMATION SEEKING AND SHARING: HOW PRIVATE FACEBOOK GROUPS

- 547 CONTRIBUTED TO THE ANTI-VACCINE MOVE- 599
 548 MENT ONLINE. *AoIR Sel. Pap. Internet Res.* DOI: 600
 549 [10.5210/spir.v2020i0.11252](https://doi.org/10.5210/spir.v2020i0.11252) (2020). 601
 550
- 551 **8.** Pennycook, G. & Rand, D. G. Fighting misinformation 602
 552 on social media using crowdsourced judgments of news 603
 553 source quality. *Proc. Natl. Acad. Sci. United States Am.* 604
 554 **116**, 2521–2526, DOI: [10.1073/pnas.1806781116](https://doi.org/10.1073/pnas.1806781116) (2019).
 555
- 556 **9.** Gollwitzer, A. *et al.* Partisan differences in physical 605
 557 distancing are linked to health outcomes during the COVID- 606
 558 19 pandemic. *Nat. Hum. Behav.* **4**, 1186–1197, DOI: 607
 559 [10.1038/s41562-020-00977-7](https://doi.org/10.1038/s41562-020-00977-7) (2020).
 560
- 561 **10.** Jin, F. *et al.* Epidemiological modeling of news and 608
 562 rumors on Twitter. In *Proceedings of the 7th Workshop on 609
 563 Social Network Mining and Analysis, SNA-KDD 2013*, 1– 610
 564 9, DOI: [10.1145/2501025.2501027](https://doi.org/10.1145/2501025.2501027) (Association for 611
 565 Computing Machinery, New York, New York, USA, 2013).
 566
- 567 **11.** Wang, L. & Wood, B. C. An epidemiological approach 612
 568 to model the viral propagation of memes. *Appl. Math. 613
 569 Model.* **35**, 5442–5447, DOI: [10.1016/j.apm.2011.04.035](https://doi.org/10.1016/j.apm.2011.04.035) 614
 570 (2011).
 571
- 572 **12.** Martin, T., Hofman, J. M., Sharma, A., Anderson, A. 615
 573 & Watts, D. J. Exploring limits to prediction in complex 616
 574 social systems. In *25th International World Wide Web 617
 575 Conference, WWW 2016*, 683–694, DOI: [10.1145/2872427.2883001](https://doi.org/10.1145/2872427.2883001) (International World Wide Web 618
 576 Conferences Steering Committee, Republic and Canton of 619
 577 Geneva, Switzerland, 2016).
 578
- 579 **13.** Arif, A. *et al.* How information snowballs: Exploring 620
 580 the role of exposure in online rumor propagation. In 621
 581 *Proceedings of the ACM Conference on Computer Sup- 622
 582 ported Cooperative Work, CSCW*, vol. 27, 466–477, DOI: 623
 583 [10.1145/2818048.2819964](https://doi.org/10.1145/2818048.2819964) (Association for Computing 624
 584 Machinery, 2016).
 585
- 586 **14.** Kimura, M. & Saito, K. Tractable models for information 625
 587 diffusion in social networks. In *Lecture Notes in Com- 626
 588 puter Science (including subseries Lecture Notes in Arti- 627
 589 ficial Intelligence and Lecture Notes in Bioinformatics)*, 628
 590 vol. 4213 LNAI, 259–271, DOI: [10.1007/11871637{_}27](https://doi.org/10.1007/11871637{_}27) 629
 591 (Springer Verlag, 2006).
 592
- 593 **15.** Bakshy, E., Rosenn, I., Marlow, C. & Adamic, L. The role 630
 594 of social networks in information diffusion. In *WWW'12 - 631
 595 Proceedings of the 21st Annual Conference on World 632
 596 Wide Web*, 519–528, DOI: [10.1145/2187836.2187907](https://doi.org/10.1145/2187836.2187907) 633
 597 (ACM Press, New York, New York, USA, 2012).
 598
- 599 **16.** Lloyd-Smith, J. O., Schreiber, S. J., Kopp, P. E. & Getz, 634
 600 W. M. Superspreading and the effect of individual varia- 635
 601 tion on disease emergence. *Nature* **438**, 355–359, DOI: 636
 602 [10.1038/nature04153](https://doi.org/10.1038/nature04153) (2005).
 603
- 604 **17.** Simpson, E. & Connor, A. Fighting Coronavirus Mis- 637
 605 information and Disinformation - Center for American 638
 606 Progress. Tech. Rep., Center for American Progress 639
 607 (2020).
 608