**ORIGINAL PAPER**

# Hidden messages: mapping nations' media campaigns

**Keeley Erhardt[1]** [ORCID] **· Alex Pentland[1]**

## Abstract

Powerful actors have engaged in information control for centuries, restricting, promoting, or influencing the information environment as it suits their evolving agendas. In the Digital Age, information control has moved online, and information operations now target the online platforms that play a critical role in news engagement and civic debate. In this paper, we use a discrete-time stochastic model to analyze coordinated activity in an online social network, representing the behaviors of accounts as interacting Markov chains. From a dataset of 31,521 tweets posted by 206 accounts, half of which were identified by Twitter as participating in a state-linked information operation, we evaluate the coordination, measured by the apparent influence, between pairs of state-linked accounts compared to unaffiliated accounts. Our analysis reveals that state-linked actors demonstrate significantly higher levels of coordination among themselves compared to their coordination with unaffiliated accounts. Furthermore, the degree of coordination observed between state-linked accounts is more than seven times greater than the coordination observed between unaffiliated accounts. Moreover, we find that the account that represented the most coordinated activity in the network *had no followers*, demonstrating the power of our modeling approach to unearth hidden connections even in the absence of explicit network structure.

✉ Keeley Erhardt
keeley@mit.edu

Alex Pentland
pentland@mit.edu

1    Massachusetts Institute of Technology, 75 Amherst St, Cambridge, MA 02139, USA

᠘ Springer

# 1 Introduction

The rate of online media consumption has dramatically increased and individuals' online social networks (OSNs) are an ever more popular source for news content. State and non-state actors desiring to manipulate the information environment have adapted to this trend, launching information operations targeted at a range of online platforms. Since October 2018, Twitter has publicly identified more than 40 state-linked information operations attributed to over 20 countries targeted at its platform Twittter (2022). From 2017 through mid-2021, Facebook similarly took down and reported over 150 information operations originating from more than 50 countries Facebook (2021). An information operation can be characterized as coordinated activity aimed at a strategic objective that is fundamentally deceptive in nature (Erhardt and Pentland 2021). This deception may not necessarily imply explicitly false information (e.g., out-of-context images, agenda-setting, or flooding the information environment with superfluous messaging to confuse and distract (Starbird et al. 2019; King et al. 2017)).

Much of the literature in the mis/disinformation space focuses on detecting information operations through content-based features (Alizadeh et al. 2020; Rheault and Musulan 2021), or network-based approaches (Vargas et al. 2020). Other studies examine the temporal patterns of post activity (Luceri et al. 2020; Magelinski and Carley 2020). In this paper, we choose to instead revisit the *influence model*, first proposed in Asavathiratham (2001). This model is most similar to the temporal approach in Magelinski and Carley (2020) but has the advantage of being able to distinguish the directionality of apparent influence rather than producing an undirected account to account coordination graph. The *influence model* describes the dynamics of networked, interacting Markov chains. A Markov chain is a method for generating a sequence of random variables in which the current value is always probabilistically dependent on only the most recent previous value.

In this context, we choose to model individual social accounts as Markov chains with random variables representing post activity for a given user. With the influence model, we can measure the coordination between pairs of accounts based on post activity alone. From these coordination measures, it is possible to quantify hidden connections between accounts and, potentially, inauthentic activity. We focus on the coordination aspect of information operations for a few reasons. First, it alleviates some privacy and bias concerns associated with moderation. Second, an influence modeling approach is more language and media agnostic than content-based alternatives. Third, unlike network-based methods, this approach does not require access to the underlying network structure.

Our contributions are as follows. First, we present a novel application of the *influence model* for detecting accounts engaged in an information operation. Second, we demonstrate how state-linked accounts can be distinguished from other accounts in a network based on their coordinated post activity alone. And third, we have published an open-source Python library that efficiently implements the influence model and supports the learning of its parameters from sequences of observations. The package is available at https://pypi.org/project/influence-model/.

## 2 The influence model

The influence model describes the relationships between networked Markov chains in terms of the "influence" chains have on one another. The model is made up of a network of interacting Markov chains each associated with a node in a network. At the network level, nodes are referred to as sites and their connections are described by the stochastic network matrix $D$. At the local level, each site has an internal Markov chain $\Gamma(A)$ and assumes one of the statuses of $\Gamma(A)$ at any given discrete-time instant. These statuses are represented by a length-$m$ status vector $\vec{s}$, an indicator vector containing a single 1 in the position corresponding to the present status and 0 everywhere else:

$$\vec{s}_i'[k] = [0...010...1]. \tag{1}$$

Each chain evolves according to its own status and the statuses of its neighbors. Updating the status of the $i$th site in the influence model takes place in three stages:

1.  The $i$th site, site$_i$, randomly selects one of its neighbors to be its determining site; site$_j$ is selected with probability $d_{ij}$.
2.  The status of site$_j$ at time $k$, $\vec{s}_j[k]$, fixes the probability vector $\vec{p}_i[k+1]$ that is used in (3) to randomly select the next status of site$_i$.
3.  The next status $\vec{s}_i[k+1]$ is realized according to $\vec{p}_i[k+1]$.

A state-transition matrix $A_{ij}$ describes how the state-transition probabilities of site$_j$ depend on the previous status of site$_i$. $A_{ij}$ is an $m_i \times m_j$ non-negative matrix with rows summing to 1. $A$ is a matrix with $A_{ij}$ in its $(i, j)$th block. From the stochastic network matrix $D$ and the state-transition matrix $A$, one can compute the influence matrix $H$ that describes the "influence" exerted by and on each site in the network. $H$ represents the joint status of groups of sites in the influence model and is given by the generalized Kronecker product of $D'$ and $\{A_{ij}\}$:

$$H = D' \otimes \{A_{ij}\}. \tag{2}$$

The influence model has been applied to a number of problems, ranging from modeling failures in a power grid to recognizing functional roles in meetings (Asavathiratham et al. 2001; Dong et al. 2007) For more detail on the model, its properties, and applications, we refer readers to Asavathiratham et al. (2001) and (Pan et al. 2012).

### 2.1 The influence-model library

In conjunction with this paper, we have published an open-source Python library that provides a complete implementation of the influence model. The library supports defining new influence models and generating observations through applying the model's evolution equations. It is also provides an implementation to reconstruct

an influence model from a sequence of observations, learning the parameters *D*, *A*, and *H*. This implementation is demonstrated on the toy example presented in Basu et al. (2001):

```python
import numpy as np

leader = Site("leader", np.array([[1], [0]]))
follower = Site("follower", np.array([[0], [1]]))
D = np.array([
    [1, 0],
    [1, 0],
])
A = np.array([
    [.5, .5, 1., 0.],
    [.5, .5, 0., 1.],
    [.5, .5, .5, .5],
    [.5, .5, .5, .5],
])
model = InfluenceModel([leader, follower], D, A)
initial_state = model.get_state_vector()
next(model)
next_state = model.get_state_vector()
```

## 3 Data

In this paper, we analyze an information operation targeted at Twitter and attributed to the People's Republic of China (PRC). The operation focused on promoting Chinese Communist Party (CCP) narratives related to the treatment of the Uyghur population in Xinjiang. In December 2021, Twitter published a representative sample of accounts and tweets associated with this state-linked information operation, including 31,269 tweets from 2016 unique accounts Twittter (2022). The tweets begin April 20, 2019 and end April 5, 2021. We augment this dataset with "unaffiliated" accounts and tweets, defined as accounts and tweets still available on the Twitter platform as of March 2022 (not removed as part of a suspected information operation or for another platform violation). Tweets from unaffiliated accounts were collected using the Twitter Search API v2, selecting for tweets posted between April 20, 2019 and April 5, 2021 with at least one of the keywords or hashtags (case-insensitive): "xinjiang", "uighur", "uighurs", "uyghur", "uyghurs", "uygur", "uygurs", "uigur", or "uigurs". This search query returned a total of 14,728,582 tweets from 2,665,001 unique accounts.

To ensure a reasonable number of observations (tweets) for each account, we only consider tweets from accounts in the top one percent of accounts by total number of tweets. This means that an account must tweet at least 60 times over the two-year period to be included in the analysis. After downselecting tweets to only those posted by the most prolific accounts, we are left with 10,889 tweets from 103 state-linked accounts and 6,231,955 tweets from 27,003 unaffiliated accounts. From

these unaffiliated accounts, we randomly select 103 accounts (corresponding to the number of state-linked accounts) and their associated tweets to analyze. Our final dataset then includes 31,521 tweets from 206 accounts (50% state-linked and 50% unaffiliated).

## 4 Methodology

Each account in our dataset is represented as a site in a network graph. The two classes of accounts (state-linked and unaffiliated), as well as the true network structure (the follower-following relationships), are not known a priori. Our goal is to quantify the "influence" that determines the status of each site in the network using observed behaviors.

### 4.1 Constructing observations

Sites interact by posting messages (tweets), the observed behavior. If a site posts a message at discrete-time instant $k$, we consider the site "active" at time $k$. At any given time, a site can be in one of two states, *Active* or *Inactive*. We choose to discretize tweets into 1-h time blocks to ensure enough granurity to differentiate explicitly coordinated behavior from topics that begin to trend, while still ensuring a reasonable number of accounts are likely to be *Active* at any given time. The sequence of observations for each account represents the account's status over time.

---

**Algorithm 1:** Construct a sequence of observations for each site

1   <u>function GetObservations</u> ($posts, accounts, start, end$)
     **Input**   : All posts, accounts, and the time range of interest
     **Output:** Mapping from accounts to observations
2   delta ← 1 hour
3   **foreach** $account \in accounts$ **do**
4     **while** $start < end$ **do**
5       k ← time range from start to delta
6       **if** $account\ posted\ at\ time\ k$ **then**
7         status ← 1
8       **else**
9         status ← 0
10      AddToObservations (account, status)
11      start ← start + delta

---

Given we expect coordinated actors to collectively promote similar narratives, we are less interested in overall post activity and more interested in post activity by topic. We choose a simple definition for "topic": any entity is a topic. Each message contains zero or more entities, defined as hashtags, URLs, or user mentions. We first extract all entities from posts and then construct observation sequences for each entity individually, across all sites. For example, for the entity #hashtag, we only

consider an account *Active* if the account posts a message that includes #hashtag. We exclude any entities that were used as search terms in collecting accounts from the Twitter API. And, we normalize URLs by stripping the protocol, subdomain(s), and any query parameters.

## 4.2 Learning the state-transition matrices

In the influence model, the status of each site varies over time based on the "influence" of the other sites in the network. This influence is represented in part by the state-transition matrices covered previously. Given sequences of observations for each site, we can reconstruct the state-transition matrices using a maximum-likelihood estimate, similar to the approach in Basu et al. (2001). Each state-transition matrix is $2 \times 2$ representing the two possible statuses, *Active* and *Inactive*. If $site_j$ perfectly follows the behavior of $site_i$ (positive coordination), then $A_{ij}$ is the identity matrix. To obtain a scalar coordination measure for each state-transition matrix, we compute the Frobenius inner product of $A_{ij}$ and the identity matrix. The coordination measure can range [0, 2]. Zero represents maximum positive coordination, $site_i[k-1] = site_j[k] \, \forall \, k$, and two represents maximum negative coordination, $site_i[k-1] \neq site_j[k] \, \forall \, k$. By averaging these coordination measures across all entities, we can determine the master state-transition matrix for each pair of sites.

## 5 Results

We find that the accounts engaged in the most coordinated activity are overwhelmingly the accounts controlled by state-linked actors. Additionally, we discover that the accounts at the center of networks of coordination would not have been identifiable through analysis of the more traditional follower-following relationship network (even if it were available), as these accounts predominantly had few to no followers.

### 5.1 Account clusters

To assess clusters of accounts with high-levels of coordinated activity, we construct a coordination network from the pairwise coordination measures. A directed edge $(i, j)$ in the coordination network represents that $site_i$ exhibits apparent influence on $site_j$ with an edge weight equal to one minus the coordination measure. We are primarily interested in positive coordination—when an account mimics the behavior of another account—so only create an edge if the coordination measure is less than one (recall that zero corresponds to maximum positive coordination). This filtering means that not all accounts are represented in the coordination network. If an account does not positively "influence" another account and is not itself positively "influenced", it will be absent. We find that the clusters of accounts with high degrees of coordination are primarily controlled by state-linked actors, and that each cluster is typically made up of all state-linked or all unaffiliated accounts. This

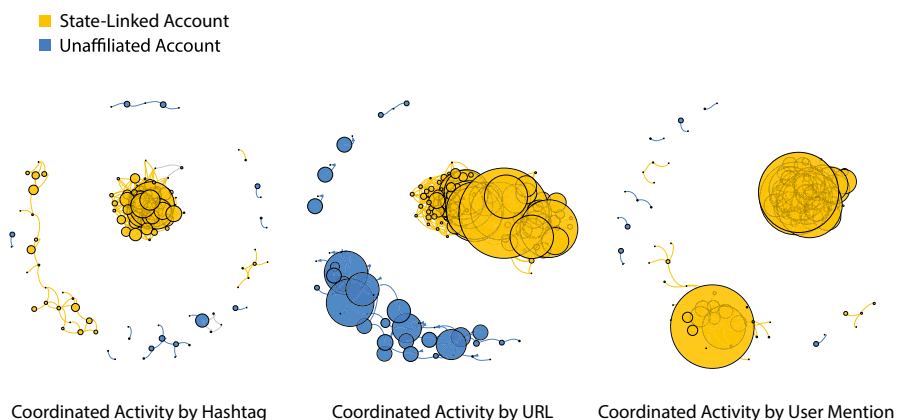**Table 1** Summary of coordinating accounts by class

| | % Coordinating accounts | | | Avg. Coordination score | | |
|---|---|---|---|---|---|---|
| | Hashtag | URL | User mention | Hashtag | URL | User mention |
| State-linked | 91.26 | 98.06 | 93.20 | 0.44 | 0.88 | 0.96 |
| Unaffiliated | 27.18 | 42.72 | 14.56 | 0.14 | 0.70 | 0.09 |

corresponds with our intuition that accounts will exhibit discernible differences in the accounts that they coordinate with based on class membership.

As show in Table 1, we observe differences in how coordination is expressed when we examine the three entity types individually. In all cases, state-linked accounts make up the majority of the accounts engaged in coordinated activity and almost exclusively coordinate with other state-linked accounts. The unaffiliated accounts are most represented in the network through URL shares, potentially due to the rapid rate at which emerging news stories can diffuse through an OSN.

## 5.2 Coordinated link sharing

For the state-linked accounts, an English-language article from Xinhua News Agency, the official state press agency of the PRC, revealed the most coordinated activity. The story condemned sanctions imposed by the United States for alleged human rights violations in Xinjiang. For the unaffiliated accounts, a Chinese-language Facebook post from the Photographic Society Of Hong Kong Media Limited (PSHK Media) describing the "sinicization" of the Uyghur population in Xinjiang



■ State-Linked Account
■ Unaffiliated Account

Coordinated Activity by Hashtag    Coordinated Activity by URL    Coordinated Activity by User Mention

**Fig. 1** The coordination between accounts. An edge is colored yellow if it represents coordinated activity between a state-linked account and another state-linked account, blue if the coordination is from an unaffiliated account to an unaffiliated account, and gray if the edge connects accounts from different classes. The size of a node is scaled by the total "influence" the account exerts. The network is visualized using the Fruchterman-Reingold layout, which places adjacent nodes physically close to one another. (Color figure online)

**Fig. 2** Coordinated activity between state-linked and unaffiliated accounts before and after Mike Pompeo ▶
publicly accuses the PRC of genocide. Nodes are scaled according to the amount of "influence" they
exert on the network, and the network is visualized using the Fruchterman-Reingold layout. (Color figure
online)

by CCP officials revealed the most coordination. The post accused CCP officials of
coercing the ethnic, Muslim minority into celebrating a traditional Chinese holiday
and consuming pork. Interestingly, Facebook blocks redirects to PSHK Media's site
from its platform and, as of the writing of this paper, the site appears to have been
suspended by its hosting provider.

### 5.3 Top influencer

Averaging across the hashtag, URL, and user mention coordination networks pro-
duces a new network consisting of 81 accounts, 75 state-linked accounts and six
unaffiliated. In this network, we find that one account exhibits a much higher degree
of coordination than any other account. This "top influencer" is state-linked, and
exclusively coordinates with other state-linked accounts. Interestingly, this account
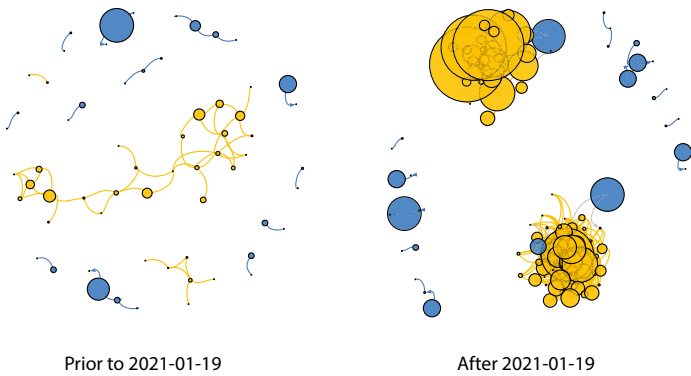did not follow any other users and *had no followers*.

The account posted 87 times during the two-year long PRC information opera-
tion. 59 tweets included a hashtag, the most popular being "xinjiang", "xinjian-
gonline", and "stopxinjiangrumors". 28 included URLs, referencing stories from
eight news or informational sites owned by the Chinese government in addition
to the People's Daily, a newspaper of record for the CCP. 71 of the user's tweets
contained user mentions. The tweets range from argumentative, countering alle-
gations of state-mandated sterilizations and forced labor in Xinjiang, to upbeat,
describing the happy, peaceful, and productive lives of people in the region. This
account was often one of the first to begin sharing new content at the beginning
of a campaign, e.g., a previously unshared URL that was then picked up by other
state-linked accounts, contributing to its high influence score.
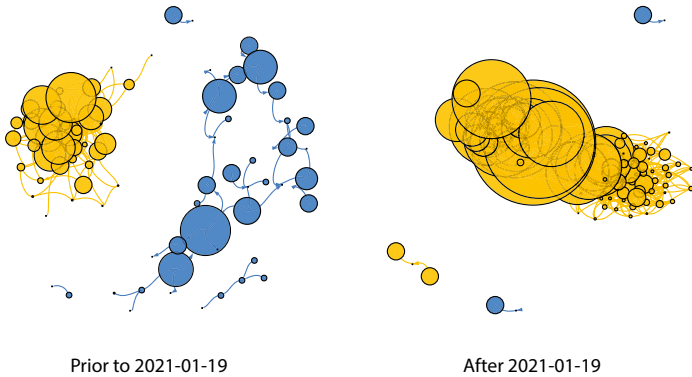
### 5.4 State-linked tweet activity

On January 19, 2021, his last day in office as U.S. secretary of state, Mike
Pompeo released a press statement accusing China of "ongoing" genocide per-
petuated against the Uyghur population in Xinjiang Pompeo (2021). The state-
ment appears to have triggered a dramatic uptick in tweet activity from state-
linked accounts, with many of the tweets directly mentioning Pompeo. To assess
whether the level of tweet activity has an effect on the coordination measures
that we study, we computed the same networks as in Fig. 1, this time subdividing
tweets into two groups: tweets posted before January 19, 2021 (the low-activity
period) and tweets posted after (the high-activity period). The findings remain
consistent with those observed for the entire two-year period. However, after
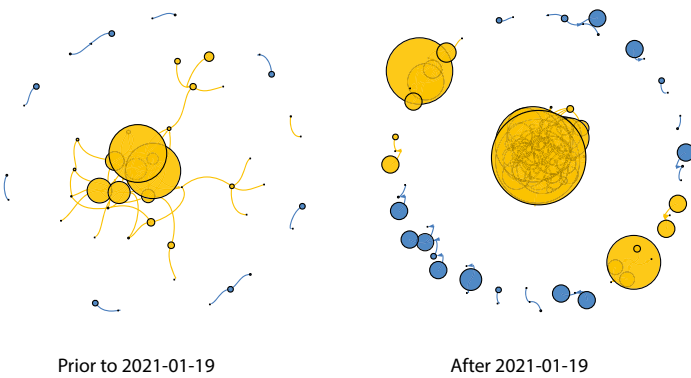Pompeo's public statement, when state-linked tweet activity reaches its peak, the

■ State-Linked Account
■ Unaffiliated Account



Prior to 2021-01-19                    After 2021-01-19

Coordinated Activity by Hashtag



Prior to 2021-01-19                    After 2021-01-19

Coordinated Activity by URL



Prior to 2021-01-19                    After 2021-01-19

Coordinated Activity by User Mention

coordination measures between state-linked accounts are even more pronounced. On average, the coordination scores show a 4.36 times increase in the hashtag network, a 1.47 times increase in the URL network, and a 2.09 times increase in the user mention network (Fig. 2).

## 6 Discussion

Analyzing network dynamics appears to show promise as a detection mechanism for uncovering explicitly orchestrated activity on OSNs. In this paper, we use the influence model to describe the relationships between the accounts engaged in discussing the Uyghur ethnic group and/or Xinjiang, an autonomous territory in northwest China where a large number of Uyghurs reside. We would like to run the same analysis on larger datasets from a broader range of information operations. Twitter has released dozens of datasets containing accounts and tweets from over 40 state-linked information operations. We are interested in how our model performs on this wide-range of campaigns.

Further, we plan to explore how recent statistical methods might be used in new ways to uncover the network dynamics that underlie coordinated activity. For example, there have been recent advancements in research into nonlinear "causal" analytics and subgraph detection (Wang et al. 2018). Similar to the approach presented in this paper, in which post activity is treated as an observed behavior, dynamical features such as posting behavior may prove useful for detecting hidden subgraphs, such as networks of state-linked actors.

## 7 Conclusion

We believe that this work represents a unique approach to detecting coordinated information operations, rooted in a well-studied model with broad utility. We show how the proposed approach successfully discriminates state-linked accounts involved in a real-world information operation from unaffiliated Twitter users, without access to user information, post content, or the underlying network structure. By considering network dynamics alone, it is possible to mitigate potential privacy and bias concerns that can arise through content moderation on OSNs. Further, network dynamics can reveal the hidden influence that underlies coordinated inauthentic activity, much of which is not discernible from standard measures of social influence online.

# References

Alizadeh M, Shapiro JN, Buntain C, Tucker JA (2020) Content-based features predict social media influence operations. Sci Adv 6(30):eabb5824

Asavathiratham C (2001) The influence model: A tractable representation for the dynamics of networked markov chains. PhD thesis, Massachusetts Institute of Technology

Asavathiratham C, Roy S, Lesieutre B, Verghese G (2001) The influence model. IEEE Control Syst Mag 21(6):52–64

Basu S, Choudhury T, Clarkson B, Pentland A (2001) Learning human interactions with the influence model. NIPS, San Francisco

Dong W, Lepri B, Cappelletti A, Pentland A, Pianesi F, Zancanaro M (2007) Using the influence model to recognize functional roles in meetings. In: Proceedings of the 9th International Conference on Multimodal Interfaces, pp. 271–278

Erhardt K, Pentland A (2021) Disambiguating disinformation: extending beyond the veracity of online content. In: Workshop Proceedings of the 15th International AAAI Conference on Web and Social Media

Facebook (2021) Threat report the state of influence operations 2017-2020. Facebook. https://web.archive.org/web/20221121175615/https://about.fb.com/wp-content/uploads/2021/05/IO-Threat-Report-May-20-2021.pdf Accessed 22 Nov 2022

King G, Pan J, Roberts ME (2017) How the Chinese government fabricates social media posts for strategic distraction, not engaged argument. Am Political Sci Rev 111(3):484–501

Luceri L, Giordano S, Ferrara E (2020) Detecting troll behavior via inverse reinforcement learning: A case study of russian trolls in the 2016 us election. In: Proceedings of the International AAAI Conference on Web and Social Media, vol. 14, pp. 417–427

Magelinski T, Carley KM (2020) Detecting coordinated behavior in the twitter campaign to reopen America. In: Center for Informed Democracy & Social-Cybersecurity Annual Conference, IDeaS

Pan W, Dong W, Cebrian M, Kim T, Fowler JH, Pentland A (2012) Modeling dynamical influence in human interaction: using data to make better inferences about influence within social systems. IEEE Signal Processing Mag 29(2):77–86

Pompeo M (2021) Determination of the secretary of state on atrocities in Xinjiang. U.S. Department of State. https://web.archive.org/web/20221121175851/https://2017-2021.state.gov/determination-of-the-secretary-of-state-on-atrocities-in-xinjiang/index.html Accessed 22 Nov 2022

Rheault L, Musulan A (2021) Efficient detection of online communities and social bot activity during electoral campaigns. J Inf Technol Politics 18(3):324–337

Starbird K, Arif A, Wilson T (2019) Disinformation as collaborative work: surfacing the participatory nature of strategic information operations. Proc ACM Hum Comput Interact 3:1–26

Twitter (2022) Transparency report: information operations. Twitter. https://web.archive.org/web/20221118021353/https://transparency.twitter.com/en/reports/moderation-research.html Accessed 22 Nov 2022

Vargas L, Emami P, Traynor P (2020) On the detection of disinformation campaign activity with network analysis. In: Proceedings of the 2020 ACM SIGSAC Conference on Cloud Computing Security Workshop, pp. 133–146

Wang S, Shen C, Badea A, Priebe CE, Vogelstein JT (2018) Signal subgraph estimation via vertex screening. arXiv preprint arXiv:1801.07683

**Keeley Erhardt**  is a Ph.D. candidate at the MIT Media Lab. Her research interests include causal models, dynamic networks, and machine learning. Previously, she was a Director of Engineering at Rebellion Defense and a software engineer at Improbable, Google, Facebook, and Apple. She earned her B.S and M.Eng. degrees in Electrical Engineering and Computer Science from MIT.

**Alex Pentland**  directs MIT Connection Science, an MIT-wide initiative, and previously was academic head of the MIT Media Lab. He is one of the most-cited computational scientists in the world, and is a member of the U.S. National Academy of Engineering. His most recent books are Building the New Economy and Trusted Data, both published by MIT Press, Social Physics, published by Penguin Press, and Honest Signals, published by MIT Press.