

## TWITTER KNOWS: UNDERSTANDING THE EMERGENCE OF TOPICS IN SOCIAL NETWORKS

Lachlan Birdsey  
Claudia Szabo

School of Computer Science  
The University of Adelaide  
Adelaide, Australia

Yong Meng Teo

Department of Computer Science  
National University of Singapore  
Singapore

### Abstract

Social networks such as Twitter and Facebook are important and widely used communication environments that exhibit scale, complexity, node interaction, and emergent behavior. In this paper, we analyze emergent behavior in Twitter and propose a definition of emergent behavior focused on the pervasiveness of a topic within a community. We extend an existing stochastic model for user behavior, focusing on advocate-follower relationships. The new user posting model includes retweets, replies, and mentions as user responses. To capture emergence, we propose a *RPBS (Rising, Plateau, Burst and Stabilization) topic pervasiveness model* with a new metric that captures how frequent and in what form the community is talking about a particular topic. Our initial validation compares our model with four Twitter datasets. Our extensive experimental analysis allows us to explore several “what-if” scenarios with respect to topic and knowledge sharing, showing how a pervasive topic evolves given various popularity scenarios.

### 1 INTRODUCTION

Systems with a large number of components, and complex and dynamic interconnections are ubiquitous (Reynolds 1987, Zhan et al. 2008) and exhibit properties that are irreducible to the behavior of their components (Darley 1994, Deguet et al. 2006, Li et al. 2006, Davis 2005, Johnson 2006, Mogul 2006). These properties, called emergent properties, are increasingly becoming important as software systems grow in complexity, coupling, and geographic distribution (Bedau 1997, Holland 1999, Johnson 2006, Mogul 2006). Examples of emergence include connection patterns in data extracted from social networks (Chi et al. 2009), trends in big data analytics (Fayyad and Uthurusamy 2002), router synchronization problems (Floyd and Jacobson 1993), and load-balancer failures in a multi-tiered distributed system (Mogul 2006). Emergent properties may have undesired and unpredictable effects and consequences, and unpredictable systems are less credible and difficult to manage; therefore, techniques for the identification and validation of emergent properties are becoming crucial. A plethora of examples of emergent properties have been identified and classified but most examples focus on simple, small systems, such as flocks of birds or the game of life, or in-house distributed systems that are not available for study (Chen et al. 2007, Holland 1999, Kubik 2003, Mogul 2006, Szabo and Teo 2012a, Szabo and Teo 2012b, Szabo and Teo 2013).

In recent years, social networks such as Twitter and Facebook have become an important communication tool, used by a large percentage of the population with access to the Internet. They are regularly used to disseminate information with respect to current or future events, to help causes and individuals, and to act as the base platform for various political or community movements (Bakshy et al. 2012). These movements are often identifiable on Twitter via hashtags such as #gamergate, #climatechange, or #BlackLivesMatter. Hashtags are also often used in short-term campaigns such as #CancelColbert or #SaveConstantine. Within these communities, the emergence of established knowledge or under-

standing can be observed, either in the form of agreement on a particular topic, or as the sharing of URLs. In this paper, we focus on Twitter as it represents a significant and widely used information dissemination network, with 284 million active users generating 500 million tweets every day.

In this paper, we define emergence in Twitter as the adherence of groups of users to the opinions of various topic advocates. We propose a first step to identify emergence by analyzing how pervasive a topic is within a group of Twitter users. Understanding how a Twitter topic became pervasive, trending, and the subsequent underlying knowledge within a Twitter community can have significant benefits for social analysis, and for identifying conflicts and opinions of the general public with respect to various social matters. It can also be used to identify credible users or, conversely, users that spread rumors (Castillo et al. 2011). Determining whether a topic has become convention in a social network has been studied using information diffusion modeling techniques. Proposed methods may require a complete analysis of the network (Shamma et al. 2011) or only information about nodes involved in discussion at a particular time (Yang and Leskovec 2010). However, they all require a number of additional resources to understand the cause of diffusion such as news feeds. This is troublesome on a large-scale, volatile network such as Twitter, especially considering that data collection tools have their collection rates limited by the Twitter API. Several methods have been proposed to understand the influence of a user (Yang and Leskovec 2010) and to predict how influential a topic or a user will be based on several factors such as geographical proximity and the interest of news media outlets (Toole et al. 2012, Myers et al. 2012, Java et al. 2007).

We propose the modeling and simulation of Twitter user behavior using a advocate-follower model adopted from Hogg et al. (2013). We extend this model to consider various response types, not just collaboration, and propose a topic pervasiveness model called RPBS (Rising, Plateau, Burst, and Stabilization) that establishes how important the topic is to the community of users under study. The main contributions of our work include:

- A model for the prediction of topic pervasiveness within a set of Twitter users based on the posting behavior of two types of users, *advocates* and *followers*.
- A comprehensive “what-if” analysis identifying the conditions under which certain topics, such as #gamergate or #climatechange become pervasive within the selected advocate-follower user group.

## 2 RELATED WORK

Emergence in social networks has been defined by several authors. Kooti et al. (2012) and Suagwara (2014) define emergence in social networks as being the widespread adoption of a norm or convention. These conventions could be a particular notation, for example, using *RT* instead of *ReTweet* when posting on Twitter, or the use of a particular taxonomy to categorize topics by using hashtags or certain keywords. Taxonomies created by users, called *folksonomies*, are emergent processes whereby users collaborate to categorize topics or things using freely chosen keywords (McAfee 2006, Mika 2007). These *folksonomies* exist on several social networks. For example, on Digg, simple keyword tags are used to drive *folksonomies* (McAfee 2006). On Twitter, the *folksonomies* are created through the use of particular hashtags. The information diffusion community defines emergence in social networks differently, instead focusing on the appearance and popularity. Myers et al. (2012) considers that topics may appear and become popular when a number of users begin to use similar information such as a particular URL. Cataldi et al. (2010) define the appearance and usage of a topic if that topic becomes popular at a particular time interval, but is not popular in any previous time interval. Their analysis suggests that terms related to a topic fluctuate before becoming widely adopted. Yang and Leskovec (2010) define that a topic has become popular when a community adopts a particular way of expression, for example by using a specific hashtag or keyword.

Information diffusion studies whether a topic, taxonomy, or notation has become convention in a social network. Shamma et al. (2011) developed a model of *peaky topics* to determine when the discussion of a topic reaches peak levels or the topic becomes persistent and part of regular conversation. *Peak*

*topics* requires the observation of the entire social network over a period of time to construct a complete corpus, as well as the comparison with external sources such as news feeds to validate the cause and time accuracy of the topic peak. Yang and Leskovec (2010) developed a *Linear Influence Model (LIM)* which aims to determine the global influence of a user over time in an implicit network. *LIM* achieves this by calculating individual influence functions for each user to determine the volume of users that adopt a piece of information at a certain time. The authors test *LIM* on a Twitter dataset to determine the influence of a hashtag from an external source across several sets of users and observe that user adoption of hashtags is largely driven by external sources. Lu and Yang (2012) re-defined an algorithm originally used to track the momentum of stocks, to track the momentum of trends and topics on Twitter to determine whether a topic was about to disappear or become very popular. However, Lu & Yang focus on external influence of news topics and prediction of when a topic is likely to become a trend. Java et al. (2007) use a log-likelihood ratio to determine how important a term is on a particular day of the week. Their model is able to determine the users within a community who contribute heavily to a particular set of terms and are able to influence other communities depending on their interaction. Several authors posit that external factors, such as media and news sources, contribute to the appearance of topics in social networks. Myers et al. (2012) analyze Twitter data collected over a period of one month to determine the appearance of information in seemingly random parts of the networks. They established that less than one third of the information on Twitter is affected by events outside the network. Toole et al. (2012) identified that real-world geographical proximity and news media are important factors that influence the adoption of topics such as innovations within a social network. In contrast, our RPBS model determines what state a topic is in according to our topic pervasiveness metric, while considering topics that may have external or internal influence.

Several models have been proposed to predict the behavior of users on social networks. Hogg et al. (2013) designed a model to predict when a user will retweet the post of a topic advocate depending on the user's interest in both the topic and the advocate. Their model uses a stochastic approach as the activity of each user on Twitter is largely unpredictable. Hogg et al. (2013) also consider that a user may view the post of an advocate depending on the position of the post in the timeline of the user. Gatti et al. (2014) propose an agent-based-model to predict how each user behaves when posting messages related to a topic. By modeling users that follow a common user, they are able to predict how each user reacts to the posts of the common users. Their approach employs text mining techniques such as sentiment analysis to train the parameters of the model. An example of a study on how knowledge becomes convention in a Twitter community is the work of Kooti et al. (2012) on analyzing Twitter data to determine when a certain notation for retweets became convention. Their study shows that the notation is driven initially by early adopters, who are users who are heavily involved with the platform at the time. Over time, the different versions of the notation spread throughout the network until they became widely used. Romero et al. (2011) analyze Twitter data and determine that the influence of a topic is controlled by its persistence and the exposure that each user has to the topic. They determine that topics based around political ideas tend to be the most persistent. This analysis is used to create a simulation model to examine the effects that a larger active topic user-base, with more hashtag combinations, has on the influence of a topic. Their method is based on Latent Dirichlet Allocation to determine the style of a particular Twitter user according to four categories: *style*, *substance*, *status*, and *social* (Ramage et al. 2010). This information can then be used to provide enhanced filtering of Twitter users. To address the current limitations, we extended the model proposed by Hogg et al. (2013) to include all possible Twitter response types and consider when a user may make a post on the particular topic or not, by considering how each user interacts with the topic.

### 3 PROPOSED APPROACH

This section presents our two models and our topic pervasiveness metric. The first model captures the posting behavior of Twitter users. The second model defines the state of a Twitter topic within a community, focusing on how pervasive or intensely discussed the topic is within the community.

### 3.1 Modeling the Posting Behavior of Twitter Users

Similar to the model proposed by Hogg et al. (2013), our model relies on posts made by *advocate* users, which are users with a large following, and a large amount of posts on a single topic. These users are likely to be the instigators of trends or the crucial links in the passing of information within the network. For example, a user in one of our datasets has 813 posts on the topic of #climatechange, representing 37.2% of their total posts. The user has 101,695 followers, with an average of 33.96 retweets and 13.70 favorites per post. The model proposed by Hogg et al. (2013) defines two behaviors for advocates and normal Twitter users respectively and is a statistical model that infers users' behaviors based on collected data. In contrast to existing research (Lerman 2007, Hogg and Lerman 2009, Iribarren and Moro 2009, Castellano et al. 2009) the model focuses on individual user behaviors rather than collective responses. A normal user  $u$  from a population of advocate followers receives a new post  $p$  from an advocate  $a$ . When  $a$  makes a new post, the post appears in  $u$ 's feed; however  $u$  may not view it or interact with it immediately. By the time of  $u$ 's later visit,  $u$ 's friends, i.e. other accounts  $u$  follows, will have generated a number,  $L$ , of newer posts, moving  $p$  to position  $L + 1$  in  $u$ 's feed.  $u$  may examine enough of this list to view  $p$ . Once viewed,  $u$  may decide to respond to it, by retweeting it to their followers.

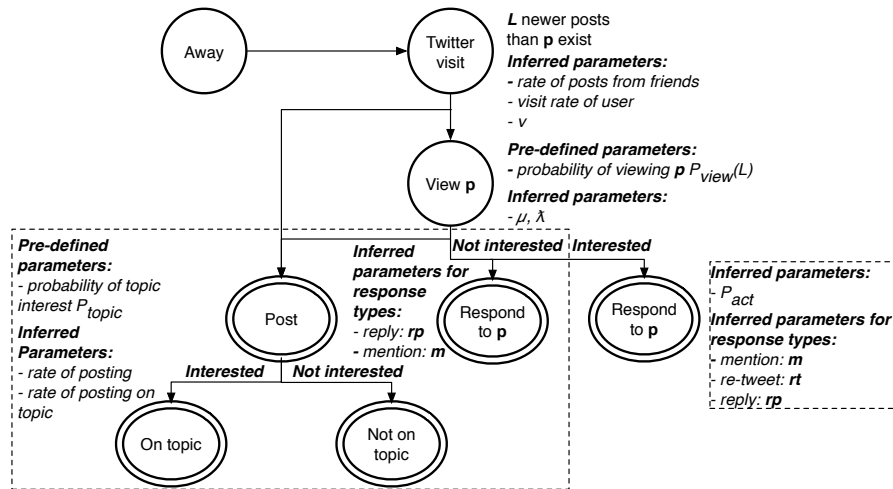


Figure 1: State machine modeling user behavior

We extend this model in two ways. Firstly, we include three types of responses, namely, retweets, replies, and mentions, and we infer different probabilities for each from the data. We distinguish between these categories as our insight suggests that collaboration through mentions and responses aids the establishment of knowledge with a community. Our future work aims to employ sentiment analysis to further distinguish between agreement and disagreement, which are mostly present in mentions and retweets. This would help identify community stances on various social issues. Secondly, if the user is not interested in a topic, the Hogg model assumes that the user will not respond. In contrast, we assign a probability to the user responding to a tweet even if they are not interested. We also consider the case that a user may not respond to the advocate but may still post, either on topic or not. We believe this captures more closely real life relationships between advocates and their followers. An extended state machine is shown in Figure 1, with dotted lines showing our enhancements. Each state in Figure 1 is annotated with the probability of user  $u$  to perform the action defined by the state. This ensures that the likelihood of  $u$  performing an action can be computed and that model parameters can be estimated from collected Twitter data.

The training phase estimates four parameters, namely  $v, \mu, \lambda,$  and  $P_{act}$ , and calculates each user's rate of posting, rate of responding to an advocate, and rate of receiving posts, amongst others. The four estimated parameters are used to define three probability distributions for each user  $u$ . The three distributions are defined similarly to those presented in Hogg et al. (2013). These parameters are obtained by performing maximum likelihood estimation similar to the procedure performed in Hogg et al. (2013),

with the exception that we consider all users that have posted on topic using any type of response, as opposed to just considering users that retweet at least one advocate post. We maximize the logarithm of an equation that considers all users, to maximize the probability of observed responses,  $P_{respond}(u)$ . We consider that the estimated parameters form an advocate profile as they apply to each follower of a particular advocate. At each view of Twitter, each user will see  $L$  posts made after the advocate's last post. For each  $u$ ,  $L$  is generated from a geometric distribution,  $P_{posts}$ .  $\mu$  and  $\lambda$  define the shape of the Inverse Gaussian Distribution for  $P_{view}(L)$ , which determines the probability of  $u$  to view post  $p$  at position  $L + 1$ .  $P_{topic}$ , the probability  $u$  is interested in a topic is calculated from a Beta Distribution with  $\alpha$  set to the number of posts  $u$  has made, and  $\beta$  set to the number of posts on topic  $u$  has made.  $P_{interesting}$  is the product of  $P_{act}$  and  $P_{topic}$ .  $P_{visible}$  is the sum of the product  $P_{posts}$  and  $P_{view}(L)$  for each new post that has appeared in the user's timeline.  $P_{respond}$  is the product of  $P_{visible}$  and  $P_{interesting}$ . Each user's  $rp$ ,  $rt$ , and  $m$  values are calculated as the fraction of replies, retweets, and mentions respectively, that are on topic, of the number of the user's total tweets, and represent the likelihood that the user would perform each type of action.

### 3.2 Modeling Twitter Topic Intensity

To advance the understanding of emergence, we propose the RPBS (Rising, Plateau, Burst, Stabilization) model to define the states of a topic within a Twitter group, as shown in Figure 2.

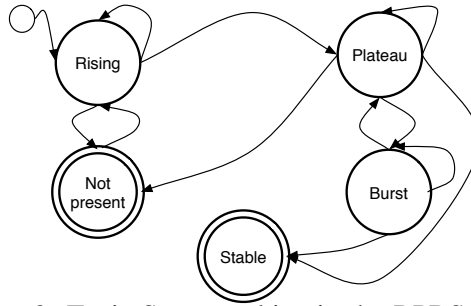


Figure 2: Topic State machine in the RPBS Model

When a topic is in the *Rising* state, the intensity of the topic is increasing. The topic intensity ( $I$ ) for the interval  $[x, y]$  is defined as:

$$I_{[x,y]} = \frac{t_{[x,y]} * n}{T_{[x,y]} * u_{[x,y]}}$$

Where  $T$  is the total number of tweets in the interval,  $t$  is the number of tweets on topic,  $n$  is the total number of users, and  $u$  is the number of unique users posting on the topic. In a rising state, the topic intensity for the interval  $[y, z]$  is greater than that for the previous interval  $[x, y]$ :

$$Rising : I_{[y,z]} - I_{[x,y]} > \epsilon$$

Once the intensity of the topic stops increasing, we define it as reaching a *plateau* state:

$$Plateau : I_{[y,z]} - I_{[x,y]} < \epsilon$$

The plateau state is sometimes followed by another increase in intensity, showing a *burst* state. This is a generalized version of the excitement phase in Riemer et al. (2012) but we allow for a burst state to occur after a gradual increase of interest as opposed to a sudden spike. It is similar to trend momentum turning positive from negative in Lu and Yang (2012) as we consider the burst state occurring after the topic intensity becomes more positive but we allow for when the prior momentum is stable for a prolonged period. The burst state is also similar to persistent terms in Shamma et al. (2011) as we consider the burst state occurring after a sustained interest in the topic but we allow for the burst state to occur even if the topic interest has not reached its peak previously. A burst state, in definition, is similar to the rising state:

$$Burst : I_{[y,z]} - I_{[x,y]} > \epsilon$$

Lastly, the topic becomes widespread and accepted within the community as it reaches a *stabilization* state:

$$\text{Stabilization} : I_{[y,z]} - I_{[x,y]} < \epsilon$$

#### 4 EXPERIMENTAL ANALYSIS

This section presents our experimental analysis. We employ four datasets for training, focused on two discussion topics, namely, #gamergate and #climatechange, with advocates posting for and against the topic respectively. Table 1 presents a summary of each dataset. To collect the data, we employed Twitter’s API to identify an advocate for/against each topic. An advocate, in our definition, is a user who has at least 10% of their collected posts dedicated to that topic. We determined whether the advocate was for or against the topic by reading their tweets. Once an advocate was identified, for each of their followers, a maximum of 3,200 tweets were collected, as this was the maximum number permitted by the API.

Table 1: Twitter Datasets

Dataset	Topic	# Tweets	# Users	# Tweets on Topic	# Users on Topic
D1	#gamergate	80,987	58	1,548	33
D2	#gamergate	45,105	35	52	8
D3	#climatechange	54,213	74	3,279	55
D4	#climatechange	587,760	978	26,372	637

We first perform an initial validation of our training module. For each dataset, the training module extracts the relevant parameters for the user population. The simulation is then run using these parameters, and a score of topic pervasiveness is calculated. The Topic Pervasiveness Range (TPR) represents a per-step score for the topic according to the state in the RPBS model. This simulated TPR is compared with that calculated from the real data. In the second stage of our experimental analysis, we run “what-if” scenarios to better understand the conditions, with respect to Twitter topology and other simulation parameters, under which a topic becomes pervasive within the user groups. For three of our datasets (D1 to D3), the advocate profiles generated from the training phase are relatively similar, with their  $\mu$ ,  $\lambda$ ,  $\nu$ , and  $P_{act}$  parameters having standard deviations of 1.27, 1.27, 4.73, and 0.02 respectively. Moreover, with the exception of D4, the profiles are similar to the profile shown in Hogg et al. (2013).

##### 4.1 Initial Validation

In our validation, we compare the evolution of the TPR calculated for the existing datasets and that of the TPR within our simulation. We perform our validation using  $\epsilon$  values of 0, 0.001, and 0.01, and interval sizes 1, 2, 5, and 10, for 25 runs each. We define an interval size,  $k$ , for an interval  $[x,y]$ , as  $k = y - x$ . We define the similarity between the TPR calculated from a simulation and from the real data as being the fraction of identical states occurring at the same point in time. The results are summarized in Table 2.

Table 2: TPR similarity with  $k = 10$  and  $\epsilon = 0$

Dataset	Mean	Median	$\sigma$
D1	0.50	0.50	0
D2	0.16	0.15	0.03
D3	0.56	0.60	0.13
D4	0.80	0.80	0

We determined that the best overall results were achieved using an interval size of 10 and  $\epsilon$  set to 0. Figure 3 contains a comparison of topic states from our best performing simulation, D4, with the topic states from the respective real dataset.

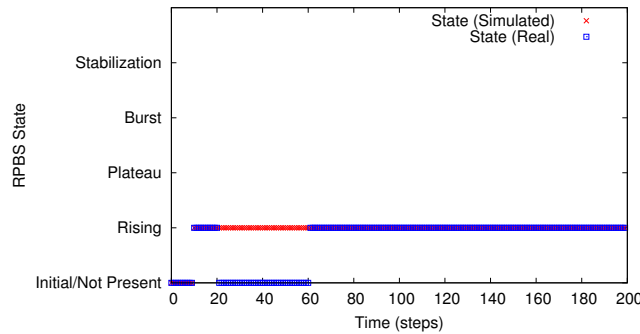


Figure 3: RPBS states for D4: Simulated & Real Data

## 4.2 Topic Emergence

We experiment with three “what-if” scenarios. In the first scenario, we consider the case in which an advocate becomes substantially more popular. We recreate this effect by increasing the number of followers that an advocate has by 10, 100, and 1,000 times. In the second scenario, an advocate becomes substantially less popular by reducing the number of its followers to a half, a fifth, and a tenth of the original follower count. The third scenario mimics the case when a topic becomes more popular by adding users from an additional advocate on the same topic. We also consider how the second advocate’s followers would behave using the original advocate’s profile, and the effects of a random advocate profile for both these cases.

We run each scenario with four  $L$  values namely, user-generated ( $L_u$ ), 1, 10, and 100, for each parameter variation unique to the scenario, for 25 times. We chose the  $L$  values 1, 10, 100 because they represent a user visiting Twitter instantly after an advocate posts, shortly after, or long after, respectively. For each experiment, we examine the final topic state that is determined by the RPBS model. We set the topic intensity interval size to 10 and  $\epsilon$  is set to 0 across all experiments.

### 4.2.1 Scenario 1: Popular Advocate

For our first scenario, we increase the number of followers of each advocate by 10, 100, and 1,000 times, and employ the  $L$  values generated from the dataset, as shown in Table 3.

Table 3: Final RPBS states for Scenario 1;  $L$  is generated from the dataset

Dataset	Follower Multiple	Final Topic State
D1	10	Rising
	100	Rising
	1,000	Rising
D2	10	Rising
	100	Rising
	1,000	Rising
D3	10	Rising
	100	Rising
	1,000	Rising
D4	10	Rising
	100	Rising
	1,000	Rising

We observe that by substantially increasing the number of followers for each advocate, each simulation terminates while the RPBS model is in the *Rising* state. We find this same behavior across all of our different  $L$  settings. We observe the behavior that as an advocate accrues followers, discussion of their respective topic by the followers does not stagnate.

### 4.2.2 Scenario 2: Less Popular Advocate

For our second scenario, we employ for the  $L$  values, the values as extracted from the dataset, with advocate follower numbers being reduced to a half, a fifth, and a tenth, as shown in Table 4. Figure 4 illustrates the case whereby our RPBS model for D1 reaches the final *Stabilization* state. This particular RPBS model behavior highlights a case where the particular topic has become emergent.

Table 4: Final RPBS states for Scenario 2;  $L$  is generated from the dataset

Dataset	Follower Reduction	Final Topic State
D1	1/2	Rising
	1/5	Stabilization
	1/10	Not Present
D2	1/2	Not Present
	1/5	Not Present
	1/10	Initial
D3	1/2	Rising
	1/5	Burst
	1/10	Burst
D4	1/2	Rising
	1/5	Rising
	1/10	Rising

We observe that reducing the number of followers for each advocate produces interesting results. Intuitively, by reducing the number of users, the persistence of a topic is more volatile. This is apparent for the experiments conducted using datasets D1 and D3. For the experiments conducted using D2, reducing the number of followers causes the discussion of the topic to halt, or in some cases to never begin. For D3, we observe that the size of the population discussing the topic directly influences the persistence of the topic. This is shown by the RPBS model progressing to the *Burst* state when the follower numbers are a fifth and a tenth of the original number of followers, as opposed to when the number of followers is half and the RPBS model remains at the *Rising* state. For D4, the advocate has such a large number of followers who post on topic, that when we reduce the number to a tenth of the original size, more than 60 followers still remain. This suggests that the followers who post on topic, do so frequently.

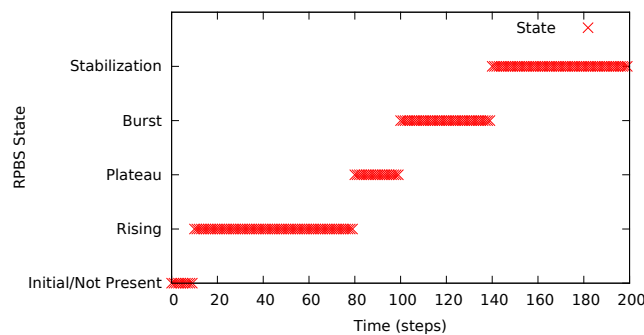


Figure 4: RPBS states over time for D1 Scenario 2

### 4.2.3 Scenario 3: Topic Becomes Popular

For our third scenario, we present the results whereby all users'  $L$  values are extracted from the dataset, with followers from both the original advocate and an additional similar advocate, or followers from only the additional advocate, as shown in Table 5. We chose each additional advocate to have the same stance on the topic as the original advocate. This is to ensure that we use followers that have similar dispositions to the followers from the original advocate. We also test using a different advocate profile which contains



the  $\mu$ ,  $\lambda$ ,  $\nu$ , and  $P_{act}$  values presented in Hogg et al. (2013), which are 14, 14, 38, and 0.12 respectively. Figure 5 illustrates the case where the RPBS model for D2 terminates at the *Burst* state.

Table 5: Final RPBS states for Scenario 3;  $L$  is generated from the dataset

Dataset	Followers	Profile	Final Topic State
D1	New Advocate	Original	Not Present
	Combined	Original	Rising
	New Advocate	Hogg et al.	Not Present
	Combined	Hogg et al.	Rising
D2	New Advocate	Original	Rising
	Combined	Original	Rising
	New Advocate	Hogg et al.	Rising
	Combined	Hogg et al.	Burst
D3	New Advocate	Original	Rising
	Combined	Original	Rising
	New Advocate	Hogg et al.	Rising
	Combined	Hogg et al.	Rising
D4	New Advocate	Original	Rising
	Combined	Original	Rising
	New Advocate	Hogg et al.	Rising
	Combined	Hogg et al.	Rising

We observe that replacing the followers of an advocate with those from a similar advocate, as well as changing the advocate profile, has little effect on the topic persistence. While these appear to be counterintuitive, especially in the case of using a different advocate profile, we speculate two possibilities. Firstly, advocates of the same topic are followed by users that tend to behave in similar ways. This is highlighted in all experiments with the exception of those conducted on D1. Secondly, the values present in the advocate profiles do not influence the stochastic processes in our user posting model strongly enough.

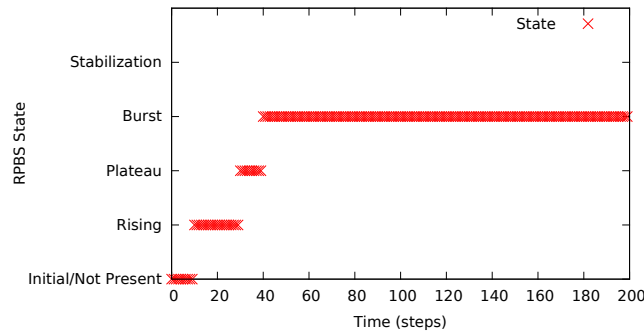


Figure 5: RPBS states over time for D2 Scenario 3

### 4.3 Discussion

Across our experiments, we found that setting each user’s  $L$  values did not alter the topic persistence greatly when compared to the topic persistence whereby each user’s  $L$  value was generated from the dataset. Two main conclusions can be drawn from this. Firstly, as we consider users that follow a particular advocate, each follower may be invested enough in the topic that the number of posts between the latest advocate post and when the follower logs in does not influence user posting. Conversely, a user may follow an advocate but exhibit little or no interest in the topic by not contributing, such as a user that follows an advocate so they receive updates about a topic without posting on it. Secondly, the stochastic influence across our simulation may not be greatly affected by selecting an  $L$  value as opposed to letting it be generated.

Our experiments were conducted with relatively small datasets, with the smallest containing 45K tweets and 35 users, and the largest containing 587K tweets and 978 users. In the scope of the topics we focused on, these datasets are very small for their respective communities. Our initial validation determined that

the accuracy of our simulated data was greatest when more users were present, and the least when fewer users were present. This highlights two main issues with our models. Firstly, our user behavior model performs better when there are more users for our training phase. This is to be expected, but if we were to focus on a topic that has a dedicated but small community, our user behavior model may not provide enough information to train an accurate model. Secondly, it suggests that our RPBS model may not be fine-tuned enough to track topic persistence when used with a subset of, or a very small, community. Alternatively, it may suggest that the interval size and/or the  $\epsilon$  value needs to be determined dynamically. Despite the obvious correlation between number of users in a dataset and accuracy, our most interesting results appeared in Scenario 2, where each advocate had heavily reduced follower counts.

The topics we focused on, #climatechange and #gamergate, are moderately different, with respect to their existence and usage. The #climatechange topic has been around for a considerable amount of time and is likely to be greatly influenced from outside sources such as news sources and other websites. This suggests that many Twitter users have been exposed to the #climatechange topic. Conversely, the #gamergate topic originated on Twitter in mid-late 2014 and has only permeated a small, albeit very active, subset of the total Twitter population and has also had limited external presence.

While our user posting model can capture the key behaviors of Twitter users, it has three main drawbacks. Firstly, our model dictates that a user will make a post regardless of the topic by following a uniform random distribution depending on how many times a day the user inspects Twitter. It is highly likely that other distributions could accurately determine when a user makes a post not on topic. This would require deeper analysis of Twitter user behavior and would require estimating more parameters. Secondly, a few pertinent features of users are not utilized, such as rapid succession of posts, how involved a user is with their own followers, and the influence of an advocate's followers on the advocate. Thirdly, our model is unable to utilize information not easily attainable from Twitter data such as how strongly a user is affected by external sources, multiple person user accounts, and bots or automated users.

## 5 CONCLUSION

Emergence in social networks is a recent phenomenon that provides significant challenges. As new topics can appear at anytime due to either internal or external causes, being able to capture and determine when they reach certain stages in a topic's evolution may provide useful insight into how a topic is affecting one or more communities. Moreover, modeling the behavior of users on a social network allows us to predict when topics reach a certain stage of evolution. In this paper we present two models. The first, a user posting model, aims to capture how followers of an advocate post either on or off the advocate's topic, and whether the post is in the form of a response. The user posting model reproduces the behavior of a user in a particular topic community. The second, our RPBS model, aims to capture the importance and adoption of a particular topic in a community. As we define emergence in Twitter as the adherence of groups of users to the opinions of various topic advocates, our RPBS model shows promising results. Our future work is two-fold. Firstly, we seek to enhance RPBS with dynamic interval sizing and  $\epsilon$  values. This could, for example, allow us to capture topic pervasiveness when a popular topic becomes dormant and then becomes popular again. Moreover, this may allow us to deploy RPBS into a live Twitter stream. Secondly, to improve our user model with sentiment analysis to determine whether a user agrees or disagrees with an advocate's post and to generate user posts that mimic both the user's posting habits and the post content.

## ACKNOWLEDGMENTS

This work is supported by the National University of Singapore under grant number R-252-000-522-112.

## REFERENCES

Bakshy, E., I. Rosenn, C. Marlow, and L. Adamic. 2012. "The Role of Social Networks in Information Diffusion". In *Proceedings of the 21st International Conference on World Wide Web*, 519–528: ACM.

- Bedau, M. 1997. "Weak emergence". *Philosophical Perspectives* 11:375–399.
- Castellano, C., S. Fortunato, and V. Loreto. 2009. "Statistical Physics of Social Dynamics". *Reviews of Modern Physics* 81 (2): 591–646.
- Castillo, C., M. Mendoza, and B. Poblete. 2011. "Information Credibility on Twitter". In *Proceedings of the 20th International Conference on World Wide Web*, 675–684.
- Cataldi, M., L. Di Caro, and C. Schifanella. 2010. "Emerging Topic Detection on Twitter Based on Temporal and Social Terms Evaluation". In *International Workshop on Multimedia Data Mining*, 4:1–4:10.
- Chen, C., S. B. Nagl, and C. D. Clack. 2007. "Specifying, Detecting and Analysing Emergent Behaviours in Multi-Level Agent-Based Simulations". In *Proceedings of the Summer Computer Simulation Conference*.
- Chi, L., W. K. Chan, G. Seow, and K. Tam. 2009. "Transplanting Social Capital to the Online World: Insights From Two Experimental Studies". *Journal of Organizational Computing and Electronic Commerce* 19:214–236.
- Darley, V. 1994. "Emergent Phenomena and Complexity". *Artificial Life IV* 4:411–416.
- Davis, P. 2005. "New Paradigms and Challenges". In *Proceedings of the Winter Simulation Conference*, edited by M. E. Kuhl, N. M. Steiger, F. B. Armstrong, and J. A. Joines, 1067–1076. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers, Inc.
- Deguet, J., L. Magnin, and Y. Demazeau. 2006. "Elements About the Emergence Issue: A Survey of Emergence Definitions". *ComplexUs* 3:24–31.
- Fayyad, U., and R. Uthurusamy. 2002. "Evolving Data into Mining Solutions For Insights". *Communications of the ACM* 45:28–31.
- Floyd, S., and V. Jacobson. 1993. "The Synchronization of Periodic Routing Messages". In *Proceedings on Communications Architectures, Protocols And Applications*, 33–44. San Francisco, USA.
- Gatti, M., P. Cavalin, S. B. Neto, C. Pinhanez, C. dos Santos, D. Gribel, and A. P. Appel. 2014. "Large-scale Multi-agent-based Modeling and Simulation of Microblogging-based Online Social Network". *Multi-Agent-Based Simulation XIV* 8235:17–33.
- Hogg, T., and K. Lerman. 2009. "Stochastic Models of User-Contributory Web Sites". In *International AAAI Conference on Weblogs and Social Media*, 50–57.
- Hogg, T., K. Lerman, and L. Smith. 2013. "Stochastic Models Predict User Behavior in Social Media". In *HUMAN*, Volume 2, 25–39: ASE.
- Holland, J. 1999. *Emergence, From Chaos to Order*. Basic Books.
- Iribarren, J. L., and E. Moro. 2009. "Impact of Human Activity Patterns on the Dynamics of Information Diffusion". *Physical Review Letters* 103 (3): 038702.
- Java, A., X. Song, T. Finin, and B. Tseng. 2007. "Why We Twitter: Understanding Microblogging Usage and Communities". In *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 Workshop on Web Mining and Social Network Analysis*, 56–65.
- Johnson, C. W. 2006. "What are Emergent Properties and How Do They Affect the Engineering of Complex Systems?". *Reliability Engineering and System Safety* 12:1475–1481.
- Kooti, F., H. Yang, M. Cha, P. K. Gummadi, and W. A. Mason. 2012. "The Emergence of Conventions in Online Social Networks". In *International Conference on Weblogs and Social Media*, 194–201.
- Kubik, A. 2003. "Towards a Formalization of Emergence". *Journal of Artificial Life* 9:41–65.
- Lerman, K. 2007. "Social Information Processing in News Aggregation". *Internet Computing* 11 (6): 16–28.
- Li, Z., C. H. Sim, and M. Y. H. Low. 2006. "A Survey of Emergent Behavior and Its Impacts in Agent-based Systems". In *Proceedings of IEEE International Conference on Industrial Informatics*, 1295–1300.
- Lu, R., and Q. Yang. 2012. "Trend Analysis of News Topics on Twitter". *International Journal of Machine Learning and Computing* 2 (3): 327–332.
- McAfee, A. P. 2006. "Enterprise 2.0: The Dawn of Emergent Collaboration". *MIT Sloan Management Review* 47 (3): 21–28.
- Mika, P. 2007. "Ontologies Are Us: A Unified Model of Social Networks and Semantics". *Web Semantics: Science, Services and Agents on the World Wide Web* 5 (1): 5–15.

- Mogul, J. C. 2006. “Emergent (mis)behavior vs. Complex Software Systems”. In *Proceedings of the 1st ACM SIGOPS/EuroSys European Conference on Computer Systems*, 293–304. New York, USA.
- Myers, S. A., C. Zhu, and J. Leskovec. 2012. “Information Diffusion and External Influence in Networks”. In *Proceedings of the International Conference on Knowledge discovery and data mining*, 33–41.
- Ramage, D., S. T. Dumais, and D. J. Liebling. 2010. “Characterizing Microblogs with Topic Models”. *International Conference on Weblogs and Social Media* 5 (4): 130–137.
- Reynolds, C. W. 1987. “Flocks, Herds and Schools: A Distributed Behavioral Model”. In *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques*, 25–34.
- Riemer, K., P. Overfeld, P. Scifleet, and A. Richter. 2012. “Oh, SNEP! The Dynamics of Social Network Emergence - the case of Cag Gemini Yammer”.
- Romero, D. M., B. Meeder, and J. Kleinberg. 2011. “Differences in the Mechanics of Information Diffusion Across Topics: Idioms, Political Hashtags, and Complex Contagion on Twitter”. In *Proceedings of the 20th International Conference on World Wide Web*, 695–704.
- Shamma, D. A., L. Kennedy, and E. F. Churchill. 2011. “Peaks and Persistence: Modeling the Shape of Microblog Conversations”. In *Proceedings of the ACM 2011 Conference on Computer Supported Cooperative Work*, 355–358.
- Suagwara, T. 2014. “Emergence of Conventions in Conflict Situations in Complex Agent Network Environments”. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems*, 1459–1460.
- Szabo, C., and Y. M. Teo. 2012a. “An Integrated Approach for the Validation of Emergence in Component-based Simulation Models”. In *Proceedings of the Winter Simulation Conference*, edited by C. Laroque, J. Himmelspach, R. Pasupathy, O. Rose, and A. M. Uhrmacher, 2412–2423. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers, Inc.
- Szabo, C., and Y. M. Teo. 2012b. “Semantic Validation of Emergent Properties in Component-based Simulation Models”. *Ontology, Epistemology, and Teleology for Modeling and Simulation*:319–333.
- Szabo, C., and Y. M. Teo. 2013. “Post-mortem Analysis of Emergent Behavior in Complex Simulation Models”. In *Proceedings of the 2013 ACM SIGSIM Conference on Principles of Advanced Discrete Simulation*, 241–252: ACM.
- Toole, J. L., M. Cha, and M. C. González. 2012. “Modeling the Adoption of Innovations in the Presence of Geographic and Media Influences”. *PLoS ONE* 7 (1): e29528.
- Yang, J., and J. Leskovec. 2010. “Modeling Information Diffusion in Implicit Networks”. In *2010 IEEE 10th International Conference on Data Mining*, 599–608.
- Zhan, B., D. N. Monekoso, P. Remagnino, S. A. Velastin, and L. Q. Xu. 2008. “Crowd Analysis: A Survey”. *Machine Vision and Application* 19 (5-6): 345–357.

## AUTHOR BIOGRAPHIES

**LACHLAN BIRDSEY** is a PhD candidate at the School of Computer Science at The University of Adelaide in Australia. His email address is [lachlan.birdsey@adelaide.edu.au](mailto:lachlan.birdsey@adelaide.edu.au).

**CLAUDIA SZABO** is a Lecturer at the School of Computer Science and an Associate Dean (DI) for the Faculty of Engineering, Computer Science and Mathematics at The University of Adelaide in Australia. Her email address is [claudia.szabo@adelaide.edu.au](mailto:claudia.szabo@adelaide.edu.au).

**YONG MENG TEO** is an Associate Professor with the Department of Computer Science at the National University of Singapore (NUS), and an Affiliate Professor at the NUS Business Analytics Center. He heads the Computer Systems Research Group. His email address is [teoym@comp.nus.edu.sg](mailto:teoym@comp.nus.edu.sg).