# Mapping Information Flows on Twitter

## Gilad Lotan

VP Research and Development, SocialFlow
416 W 13th St, #203, New York NY 10014-1178
gilad@socialflow.com

## Abstract

Social network services like Twitter and Facebook have created an expectation that you interact with your customers, followers and friends. There's an expectation to connect rather than broadcast, listen and engage in conversations. But how can we expect to interact with our invisible audience when we can't really see whose there?

For the first time in history, there is a plethora of information produced by people's actions. We can now observe a friend take another's recommendation to purchase an item, or a powerful stream of clicks to content that we choose to curate. Social media professionals are jumping on the bandwagon and attempting to quantify social interactions by using terms like influence, reach, trust and klout. But even though data is more visible than ever, it is still representative of people's complex reasoning mechanism, changing relationships, timing and logic.

This paper looks at two different ways to analyze and display characteristics of online audiences on Twitter through information flows. By visualizing flows, it is possible to "put a face" to an audience, seeing interactions between interconnected users. By replaying a representation of the series of events, it is possible to note key moments in the act of information dissemination.

## General Terms

Design, Experimentation

## Keywords

Social Media, Data Analysis, Information Flows, Audience, Social Graph, Social Analytics, Twitter, Visualization

## Introduction

Social media, and in particular Twitter, has offered a new venue for studying information and communication flows. While dominant genres of social media encourage reciprocal sharing both blogging and micro-blogging have been shown to enable rapid information flow. Marlow[1] argued that information flows on blogs were a combination of broadcast diffusion and media "contagion" – emphasizing the person-to-person dissemination of information. Likewise, Kwak et al. [2] concluded that the non-reciprocal nature of information sharing on Twitter meant that it operated more like news media, an information-sharing network, and less like a social network.

Information-sharing behavior has been studied for decades with the key theory of two-step flow of communication first developed by Katz and Lazarsfeld [3]. They determined that mass media had very little effect on how citizens voted, and that the disproportionately greater influence came from people with whom they regularly associate. These individuals were termed "opinion leaders". Wu et al. [4] tested the two-step flow theory on Twitter for normal traffic and found strong similarities in their information diffusion models.

One way of conceptualizing information flow on Twitter is through the frame of "information cascades," or situations where "it is optimal for an individual, having observed the actions of others ahead of him, to follow the behavior of the preceding individual without regard to his own information" (Bikhchandani, Hirshleifer & Welch [5]). On Twitter, information cascades are easily amplified through the common practice of "retweeting" content, or reposting the content while referencing either the source of the content or the last person who shared it (Boyd, Golder, Lotan [6]). The popularity of hashtags (Romero, Meeder & Kleinberg [7]) makes it easier for participants to follow content on a particular topic. Finally, Twitter's "trending topic" feature highlights content that is collectively on a topic that is statistically outstanding within the data. Thus, if people suddenly start talking about Egypt, Egypt is a trending topic on the front page of Twitter. The features that Twitter provides – and the way that participants use them – make it easy for information cascades to occur.

A number of online services have attempted to provide solutions for data visualizations. Many Eyes [8] lets users easily put together a visualization based on any given or pre-loaded datasets. Another service, Visualizing.org [9], acts as a social hub for people to create and share visualizations around given datasets. While these projects are important in supporting the general public's ability to

explore data, they do not at all deal with real-time, constantly updating, ephemeral content – such as is common in social spaces.

## Information Flows

In the social web, information spreads through people, networks of friends, fans and followers. Social network sites create compelling spaces where users feel comfortable to hang out, interact, consume, poke and publish. Social interactions lubricate the flow of information within these spaces, creating a plethora of dynamics. These spaces are filled with endless streams of content, encouraging users to participate, add to, consume from and redirect content. As information flows by, users grab content when it is most relevant, valuable, entertaining or insightful, and at times, choose to pass it on.

Because information flows through networks of people, attention has become a scarce commodity. Media companies no longer control people's attention, but rather fighting for a smaller portion of the pie. True power lies in understanding how information flows and its effect of where people choose to focus their attention.

In order for messages to propagate through social networks, people along the way must be attentive to the pieces of information, see them at the right time, and pass them onwards. Whether you're interested in socializing or in selling a product, understanding people's habits around information consumption and production is imperative to attaining people's attention and building an audience. By leveraging the publicly available data around people's practices, we can create services that shed a light on people's habits and preferences. Additionally, by mining this data over time, we can infer their value in affecting information flows.

In order to identify information flows on Twitter, I extract an ordered set of near-duplicate tweets (in many cases a retweet chain, but not necessarily). Then, I identify flows by finding very similar tweets using the shingling method for string comparison (Manning, Raghavan, & Schütze, 2008), which converts a string of text (such as a tweet) into a fingerprint summary of the words it comprises. This fingerprint can then be efficiently compared against other strings (other tweets) to find near-duplicates. This methodology parallels the one used in *Retweet Revolution* – visual analysis of tweets surrounding the 2009 Iranian election protests – example below.

## Example Cases

The following examples display two different methods to analyze and visualize information flows on Twitter.
Retweet Revolution
Retweet Revolution is a visualization application, that displays the most popular conversation threads that were passed amongst Twitter users at the time of the events following the Iranian elections in June of 2009. Twitter

served as an incredibly engaging mechanism to disseminate information on the riots and protests that were taking place around the world. Its realtime qualities enables information to rapidly spread between users, while its personal style drives a sense of emotional involvement to the events. This piece aims to help viewers grasp which of the messages were chosen to be passed on by millions of twitter users, and how they were manipulated along the way

By far, the messages with imminently important information, received an overwhelming amount of retweets: from posting proxy IP addresses to passing on a plea to wear green. However, as the events played out, users learned to post messages without linking to the origin as a means of protecting the Iranian sources. Retweeting is revolutionizing the way people connect to news and newsworthy events.

The applet displays 372 of the most popular threads extracted from a pool of over 230,000 messages posted on Twitter between June 14th and June 24th. It is important to remember that this is only a sampling of the Twitter data polled from the public timeline at regular intervals. Nevertheless, the content is substantial enough to identify trends and get a sense for people's practices.

The visualization begins on the 14th of June, displaying twitter messages as they come in throughout the night after election day in Iran. Each tweet is analyzed and placed within its corresponding thread, which then grows taller. As the time moves forwards, it is possible to see the different threads appear from the right. By clicking on one of the threads, the applet enters a focus screen that displays the chosen thread's network structure. The yellow node in the middle represents the earliest published tweet, corresponding to this thread that was found by the script. Many times it serves as the central node - the starting point of the conversation. When focusing on a thread, it is possible to see how a message was retweeted from one user to another, along with how the message content was changed as it was passed.
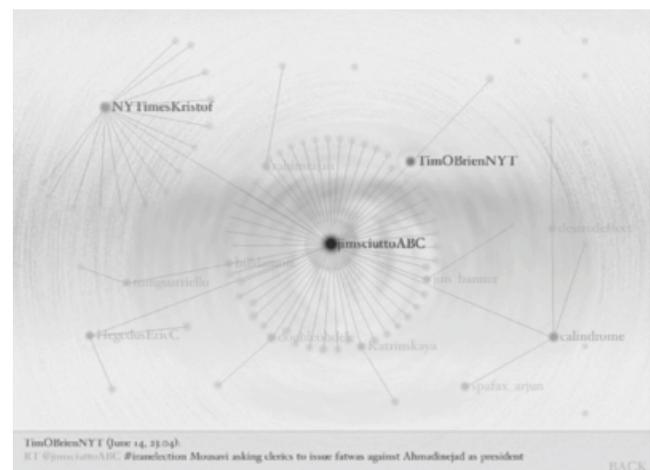


*Image 1: @jimsciuttoABS writes: #iranelection Mousavi asking clerics to issue fatwas against Ahmadinejad as president.*

A number of the users who retweeted the message in the image above (@NYTimesKristof, @HegedusEricC, @calindrome) were retweeted themselves by their followers. In this structure, it is evident that there is one major hub who is also the source of the communication. Nodes that appear with no connection to any other nodes on the screen are similar enough to be part of the same thread, yet either did not include an attribution to another user or did not use an RT format that the script could detect.
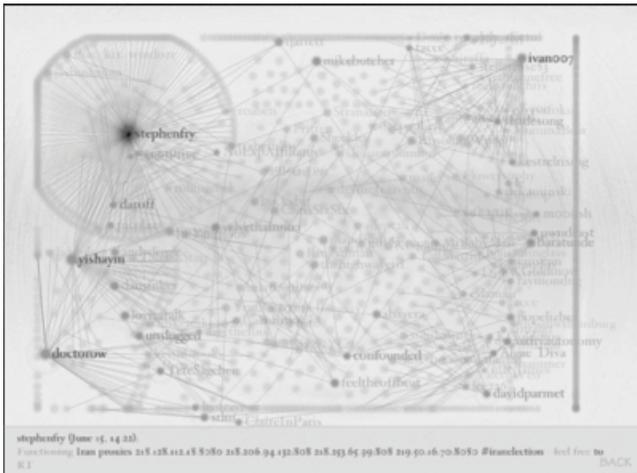


*Image 2: Massive dissemination of information about functioning IP addresses.*

In the example thread above, users were passing around information about the functioning proxy IP addresses that could be used from within Iran. @mandelbrot5 was the very first person caught sending out these proxy numbers (June 14th, 19:03). In the early morning of June 15th, a variety of users (@ivan007, @yishym, @doctorow) retweeted a very similar message describing the proxy addresses, finally reaching @stephenfry who posted the proxy addresses the following afternoon. The line of nodes along the outskirts represent either tweets that have no RT attribution (but clearly belong to this thread) or use syntax that my algorithm cannot understand.
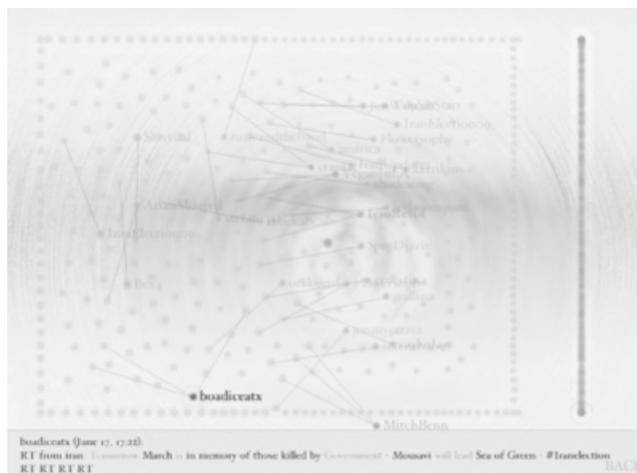


*Image 3: Eventually users intentionally don't attribute information sources for fear of Iranian gov't crackdown.*

By June 19th, Twitter users were already trained not to mention the source of their message, and instead, labeled their source as "RT from iran:". The thread above displays post around the organization of a memory march. In this example, it is impossible to make out any central source or authority from which this message originated.

## Seeing a Twitter Hashtag Spread

*#CheeringForTheYankeesIsLike* is a hashtag created by *@mattsly* the morning of October 26th. He submitted the following snarky message:
'*Go Phillies. #CheeringForTheYankeesIsLike hoping investment bankers get really huge bonuses of at least 8 figures*'

Hoping to entertain his friends, and possibly get others to participate. Matt had 182 followers at the time, not sizeable by any means on Twitter. Little did he expect that some 9 hours later, 271 different users, most of whom have no connection to him whatsoever, would participate, posting around 500 messages in total. How did this happen and what prompted this message to spread?

About an hour after Matt sent out his first message, one of his followers, @lizzieohreally, wrote the following message:

> @jaketapper? @abcdude? ...Hoping someone w/ more Twitter than I can help popularize #CheeringForTheYankeesIsLike (via @mattsly).

Lizzie clearly understood that in order to get many others to play, she would have to get someone with a large set of followers to participate. Lizzie had only around 500 followers at the time, so posted this message in an attempt to seek @jaketapper or @abcdude's attention.

Sure thing, some twenty minutes later, @abcdude see's the message and adds his own variation to the meme: *#cheeringfortheyankeesislike pulling for Regina George in "Mean Girls."*

He enjoys it so much that he promptly posts another message and attaches the hashtag. *@abcdude* is a new york based correspondent for ABC news. He dubs himself a Red Sox fan and a cosmic power broker. Not as cosmic as Lizzie had hoped, but still, he has some 7,000 followers, which could certainly help give the meme some traction. We see a small spike after *@abcdude*'s participation, and by now, some 3 hours after Matt sent the original message, there have been 34 different messages posted with this unique hashtag.

But it wasn't until @jaketapper joined in that the conversation really took off. The hashtag came to Jake's attention after @DetourJazz, whom he follows, participated. Jake reacted by posting: *'RT @DetourJazz: #cheeringfortheyankeesislike rooting for "Craterface" in Grease to beat Danny (via @Laura_Martin)'*. He then added a new message that he posted to his followers. Jake is a senior White House correspondent for ABC news with

25

over 30,000 followers. Before he took part in this meme, new posts appeared at a frequency of one every 5 minutes. Immediately after he joined, we see a sharp rise in participation, with multiple messages from a variety of users every minute.

The image below is a network graph where every user taking part in the meme is represented by a red circle and Matt, whom first started the meme, is shown in yellow. Edges represent the person who most likely influenced the other to first participate. This is an initial way to look at the information spread. One noticeable aspect is the dense cluster of users around @*jaketapper*.



*Image 4: Graphing the network.*

The image below is a snapshot of a dynamic applet highlighting the increase in acceleration of information flow. As the timeline moves forward, each profile lights up when they post a new message with the hashtag. The moment that @*jaketapper* chose to participate is evident – there's a clear, sudden spike in participation after his profile picture lights up.



*Image 5: Seeing the flow.*

The following image highlights the social hops that this hashtag traveled between users. Each node is a user, while edges represent the social ties – when there's a

follower/friend relationship between two users, a line is placed between the appropriate nodes. The first column includes only Matt who first used the hashtag. The second row consists of only those people he directly influenced to participate (out of his followers). While there are a total of 9 columns, it is crystal clear that the most important phase happened in the second and third column, when a core cluster of users chose to participate, and a mini tipping point was reached.
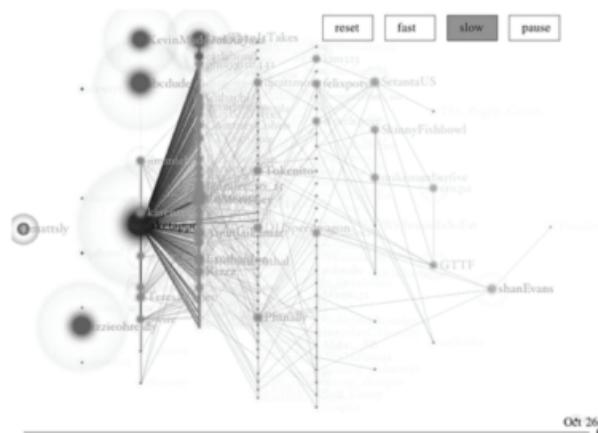


*Image 6: Social distance.*

#*CheeringForTheYankeesIsLike* lasted for a total of 9 hours that day, activated 271 different users and included around 500 messages in total. From looking at this meme, it is clear that on Twitter, there's great advantage to having many followers if one intends to spread a message. It is also clear that having the right followers is key. If it were not for @lizzieohreally who knew to actively pass the message onwards to heavy Twitter users, the meme would never have spread out the way it did. In order to come to these conclusions it was necessary for me to look at social ties in addition to the semantics of the messages posted.

## Conclusion

We all build mental models of our audiences in our head by imagining those that tend to give us attention based on sheer social knowledge and past experience. However, many have reached a point of data saturation. What we need are tools that scale and capture our networks as a whole and not just the individuals. Analysis should dive into how people in our online audiences are interconnected, and ways for us to optimize information flow.

Social network analytics and visualization tools may fundamentally change the way we engage with our online audiences. We need to build better tools that will not only increase efficiency for marketers, but any user attempting to engage with an audience online. Additionally, more work is needed to better understand how information flows between sources. How does information leap across contextual barriers? What is the relationship between

regional actors and global actors? To what degree are users consuming tweets and incorporating that knowledge into articles without retweeting the messages? Which tweets are actually read by followers? How do users choose who they trust?

This piece highlights two visual methods to track information dissemination on Twitter, where by visually representing the observed data, patterns and understanding emerged. These projects are an attempt to move away from the network graph or generic plots, and create visually enticing interfaces for interaction with and extraction of data. There is plenty of work to be done on data analysis, but we must not forget that it is also necessary to focus on the visual representation of this data.

# References

[1] Marlow, C. (2005). *The Structural Determinants of Media Contagion* (Unpublished PhD thesis). MIT Media Lab, United States. Retrieved February 1, 2011, from http://cameronmarlow.com/papers/phd-thesis.

[2] Kwak, H., Lee, C., Park, H., & Moon, S. (2010). What is Twitter, a Social Network or a News Media? *Proceedings of WWW'10.*

[3] Katz, E., & Lazarsfeld, P. F. (1955). *Personal Influence*. Glencoe, IL: Free Press.

[4] Wu, S., Hofman, J. M., Mason, W. A., & Watts, D. J. (2011). Who Says What to Whom on Twitter. *Proceedings of WWW'11.*

[5] Bikhchandani, S., Hirshleifer, D., & Welch, I. (1992). A theory of fads, fashion, custom, and cultural change in informational cascades. *Journal of Political Economy, 100*(5).

[6] boyd, d, Golder, S, and Lotan, G. (2010). Tweet Tweet Retweet: Conversational Aspects of Retweeting on Twitter. *Proceedings of HICSS-43.*

[7] Romero, D., Meeder, B., & Kleinberg, J. (2011). Differences in the mechanics of information diffusion across topics: Idioms, political hashtags, and complex contagion on Twitter. *Proceedings of WWW'11.*

[8] http://www-958.ibm.com/software/data/cognos/manyeyes/

[9] http://www.visualizing.org/