# NERCCS 2018:
# First Northeast Regional Conference on Complex Systems

April 11 (Pre-Conference Event), 12–13 (Main Conference), 2018     Binghamton, NY

**Abstracts / Papers**

# Abstracts and Papers for the First Northeast Regional Conference on Complex Systems

Center for Collective Dynamics of Complex Systems
Binghamton University, State University of New York
`http://coco.binghamton.edu/nerccs/`

April 11–13, 2018
Binghamton, NY

# Contents

## IV    Contributed Talks 4: Network Dynamics     41

## V    Contributed Talks 5: Social Networks     60

# IX   Posters            102

**Part I**

# Thursday April 12th 11:00am-12:15pm Contributed Talks 1: Network Structure (Symposium Hall)

# Sampling Community Structure in Dynamic Social Networks

Humphrey Mensah     Sucheta Soundarajan

Syracuse University, Syracuse, NY

{hamensah,susounda}@syr.edu

When studying dynamic networks, it is often of interest to understand how the community structure of the network changes, giving us insight into questions such as: when did a group split or how long did it take for a group to split? However, before studying the community structure of dynamic social networks, one must first collect appropriate network data. In this paper we present a crawling technique to crawl the community structure of dynamic networks when there is a limitation on the number of nodes that could be contacted.

## 1   Introduction

Researchers are interested in a wide variety of problems related to communities in dynamic social networks, including understanding their growth, dissolution, and merging behaviors. However, before studying such questions, a researcher must first obtain an appropriate representation of real network data. Because typical social networks may contain millions or billions of nodes, it can be a challenge to collect adequate data within a reasonable amount of time, due both to the computational efforts required to collect such data as well as API rate limits imposed by the companies owning the data. Given such a scenario, a data collector must make the most of a limited query budget. We consider two different problem settings: (1) The query budget limits the total number of queries that can be made over the entire timeline (e.g., queries cost money, and we have a fixed amount of money for the entire sampling process). (2) There is a query limit for each timestep (e.g., queries take time, and each day has a limited amount of time).

## 2   Proposed Approach

This work proposes a novel algorithm (DYNSAMP) for sampling a dynamic network such that the community similarity between the true and sampled networks is maximized. The intuition behind DYNSAMP is that the current snapshot of a graph may be similar to an earlier snapshot; or if not, portions may be similar.

DYNSAMP begins by obtaining a sample for the first day of the sampling process with an allocated number of queries. For subsequent days, a fraction of the budget allocated for that day is used to obtain a graph called the *startup graph*. The startup graph is then compared to previously discovered graphs to determine if they are similar. If similar, a portion of the budget allocated for that day is saved for future use. If not similar, the entire allocated budget for the day is used. If there is saved budget, it is used to perform extra queries to grow the graph.

Figure 1 shows the performance of DYNSAMP on one dataset (Autonomous System) with a setting where



Figure 1: Similarity comparison between sampled graph and true graph for different timesteps on the Autonomous System.

there is a limit over the entire timeline. Similar results were observed in other datasets using the same setting. Currently, the dissimilarity threshold is assumed to be fixed across all time steps. A possible future work could be a dynamic definition of the dissimilarity threshold. Another direction could be the inference of dynamic behavior of an unseen node.

# Evaluation of Community Similarity based on Hierarchical Distance

Ricky Laishram    Sucheta Soundarajan

Syracuse University, Syracuse, NY, {rlaishra,susounda}@syr.edu

## Abstract

In network analysis, there are a number of techniques for calculating the similarity between two sets of communities, such as Jaccard Similarity, Mutual Information etc. are used. However these measures do not account for the "closeness" of the different communities, and as result, they can be misleading. In this paper, we examine this problem and propose a method of computing the community quality based on the distances in hierarchical community.

## 1    Introduction

A common task in network analysis is community comparison. For example, if two community detection algorithms identify two different sets of communities, how similar are those results? Common community similarity metrics include Jaccard Similarity, Mutual Information, etc. However, when considering hierarchical community structure, these measures do not account for the "closeness" of the different communities, and so can be misleading. In this paper, we examine this problem and propose a method of computing community quality based on the hierarchical distance.



Figure 1: Similarity comparison for different real-world networks.

For example, consider the results shown in Figure 1, which depict the similarities between communities found by multiple runs of the Louvain modularity maximization algorithm on the same graph, across 14 real-world graphs. If we compare these communities using Jaccard similarity, we see that the community similarity may be very low. The Louvain method is non-deterministic, and so some variation in results is expected, but these results are shockingly low. Results for other standard measures are not shown because of space limitations, but are similar. We postulate that this occurs not because the detected communities are actually so dissimilar, but because the comparison metrics fail to take into account the hierarchical structure of the communities.

## 2    Community Hierarchical Distance

To address the problem described in Section 1, we introduce the *Community Hierarchical Distance (CHDist)*. Suppose that $\mathbb{C}$ and $\mathbb{C}'$ are two sets of communities in a graph $G$. The idea behind CHDist is that if a node $u$ is in community $C \in \mathbb{C}$, but in a different community $C' \in \mathbb{C}'$, the penalty for this should be based on the change in modularity if $C$ and $C'$ are merged. (Other measures of community quality can also be used in place of modularity.)

Let $\mathcal{H}_{\mathbb{C},G}$ be the hierarchy of communities in $G$ with the elements of $\mathbb{C}$ as the leaves. For $C_0, C_1 \in \mathbb{C}$, let $\eta\left(C_0, C_1, \mathcal{H}_{\mathbb{C},G}\right)$ be the normalized height of the smallest $C_\cup \in \mathcal{H}_{\mathbb{C},G}$ such that $C_0 \cup C_1 \subset C_\cup$. For a node $u$, let $\gamma(u, \mathbb{C})$ be the community $C \in \mathbb{C}$ such that $u \in C$. Let $V_G$ denotes the set of all nodes in $G$.

For $C \in \mathbb{C}$, let us define, $\beta\left(C, \mathbb{C}'\right) = \underset{C' \in \mathbb{C}'}{argmax}|C \cap C'|$. Then,

$$\delta_H(\mathbb{C}, \mathbb{C}') = \frac{1}{|V_G|} \sum_{u \in V_{G_0}} \eta\left(\gamma\left(u, \mathbb{C}'\right), \beta\left(\gamma\left(u, \mathbb{C}\right), \mathbb{C}'\right), \mathbb{C}'\right)$$

Then we define the Community Hierarchical Distance between $\mathbb{C}$ and sample $\mathbb{C}'$ as the harmonic mean of $\delta_H(\mathbb{C}', \mathbb{C})$ and $\delta_H(\mathbb{C}, \mathbb{C}')$. As seen in Figure 1, the community hierarchical distance (denoted in red color) is much closer to the best theoretical value of $1.0$ for all the networks considered.

# A Comparison of Community Detection Techniques Across Thematic Twitter Emoji Networks

Ryan Hartman, S.M. Mahdi Seyednezhad, Diego Pinheiro,

Josemar Faustino and Ronaldo Menezes

BioComplex Lab, Department of Computer Science
Florida Institute of Technology, Melbourne, Florida, USA
{rhartman, mseyednezhad, dsilva, jcruz, rmenezes}@biocomplexlab.org

### Abstract

Emojis are emerging as an alternative way to interact and communicate online, and their large-scale adoption has the potential to reveal hidden patterns of human communication and social interactions. In this work, we investigate the hypothesis that emojis are a form of language. By building networks of emoji co-occurrence, we examine the diversity of the community structure of such networks with regards to predefined categories of emojis. Using four different techniques of community detection, we validate our hypothesis on six Twitter data sets: five from specific topics and one random data set. Our results demonstrate that the community structure of emojis is more diverse when they are used in non-random topics such as politics and sports, and that Stochastic Block Models appears to extract communities with higher diversity.

## 1 Introduction

Online social networks have rapidly evolved and attracted a significant number of users. In these social networks, users mainly convey their feelings and emotions via short messages, which can be challenging and facilitate the occurrence of misunderstandings. Although these users have previously employed *emoticons* (e.g., ";-)", ";p") along with text to express their feelings, the need for pictographs to supplement text led to the emergence of emojis.

Emojis are grouped into seven categories,*Smiley-People, Animals-Nature, Food-drink, Activities, Objects, Symbols, and Flags* and can be found on social media sites such as Instagram, Facebook, and Twitter. Currently, more than half of the posts on Instagram contain emojis, and attract a 17% higher interaction rate when emojis are used [8]. In 2015, emojis were announced as the greatest growing language in the United Kingdom [3], and given the ever-increasing usage of emojis in social media, the assessment of emojis as a form of language deserves further attention.

One way of analyzing emojis is by investigating the sequences formed by their co-occurrence in social media posts. Users may communicate and express their feelings using a more diverse set of emojis such that emojis of different categories will appear together with a higher likelihood. By building networks of co-occurring emojis, we can unveil their community structure and subsequently assess the diversity of their community structure with respect to categories of emojis. In this work, we use the assessment of communities proposed by Hartman et al. to examine the diversity of communities of emojis from six data sets of different thematics using four techniques of community detection. By analyzing the structural properties of such communities as well as their resemblance to available metadata, our work sheds light on the idea that emojis are a form of language.

## 2 Related Work

Recent studies on emojis can be divided into two major approaches. In the first approach, researchers aim to understand the meaning of emojis. Barbieri et al. [2] investigated the meaning of Twitter emojis by examining the likelihood of the pairwise appearance and measuring how often emojis convey the same meaning. Novak et al. [12] drew a sentiment map of the 751 most frequently used emojis and found high frequency of usage associated with positive tweets. Wijeratne et al. [16] created a dictionary to make a machine readable sense inventory for emoji. In order to create octuples representing the meaning of the emoji, they used the Unicode, description, image, and keywords attached to the meaning of the emoji.

The second approach attempts to analyze the collective behavior of users based on emoji usage. Novak et al. [12] found that the inter-annotator agreement of tweets containing emojis were higher than the ones without emojis. More interestingly, they acknowledged that users normally use emojis at the end of tweets, and the rank of emojis did not change between different languages. Seyednezhad et al. [15] extracted a network of emojis based on their co-occurrence in tweets from two different datasets. They stated the emoji with the maximum edge betweenness could give us a hint about the underlying subject in which the tweets were collected. This work was generalized by Fede et al. [5] by experimenting with more data sets which contained directed networks. They concluded that important emojis are topic dependent. Lu et al. [11] created a network of emojis by point-wise mutual information (PMI). Their findings pointed to a strong correlation between social indicators and patterns of emoji usage.

Using networks of emojis, we can extract the structure of related emojis using community detection techniques. Community detection techniques aim to identify the building blocks of networks and their structural properties. It has been applied to networks of protein interaction, food web, genetic disorders, gene expression, and social networks [6]. However, the most efficient techniques for exploring communities may yield different results [10, 9]. Hence, recent works have used community detection algorithms in a holistic approach [9], which includes a comparative analysis of multiple techniques as well as the the resemblance of extracted communities to available metadata. This is the direction we pursue in this work.

## 3 Data and Methods

### 3.1 Data Collection and Curation

For this study, we collected tweets from Twitter based on different topics at different time periods. The goal is to cover a diverse set of topics, allowing us to examine the effect of such diversity on communities, extracted by state-of-the-art community detection techniques. Moreover, we add to the analysis a topic-free dataset. This data contains tweets randomly sampled from the Twitter feed, without the use of tracking keywords. The random data allows us to observe any possible bias due to using topic-based data. Table 1 shows further information about the datasets used in this work.

In order to show the network statistics are correlated with the emojis' frequency of usage, we calculated the Spearman's rank correlation between the frequency and weighted degree of the nodes. The highly correlated ranks suggest network characteristics such as weighted degree can explain the frequency of emoji usage.

Table 1: Six data sets collected from Twitter. The topics of the data sets covers several areas of interest. The Spearman's rank correlation is calculated for each data set between the frequency and weight-degree of the nodes.

| Label | Dataset | Characteristics | # Tweets (Millions) | % Containing emojis | Collection period | Spearman's rank correlation |
|---|---|---|---|---|---|---|
| $D_1$ | G-20 | Surnames of G-20 countries' leaders | 10.6 | 7% | Aug. 24 - Sep. 24, 2014 | 0.94 |
| $D_2$ | Organ | Organ transplantation terms | 2.5 | 9% | Oct. 2015 - Apr. 2017 | 0.85 |
| $D_3$ | rioSports | Sports in the 2016 Rio Olympics | 1.8 | 1% | Aug. 05  Aug. 21, 2016 | 0.95 |
| $D_4$ | rioTerms | "Olympics" in different | 5.8 | 1% | Aug. 05  Aug. 21, 2016 | 0.92 |
| $D_5$ | WWC | Women's World Cup 2015 | 10.7 | 1% | Jun. 06 - Jul. 05, 2015 | 0.91 |
| **$D_6$** | **randSample** | **2 months samples from Twitter** | **168.5** | **< 1%** | **Dec. 13, 2016 - Jan. 31, 2017** | **0.97** |

In summary, we have data related to politics ($D_1$), health ($D_2$), sports ($D_3$, $D_4$, and $D_5$), as well as a random collection of tweets ($D_6$). The random sample $D_6$ has the lowest percentage of tweets containing emojis, while the organ transplantation $D_2$ collection has the greatest amount.

## 3.2  Network Construction

The main focus of this paper is on comparing prominent community detection algorithms and the characteristics of the communities they uncover for a variety of datasets. Differing from previous works [15], here we consider that the order of emojis appearing in a tweet is fundamental and hence better represented using directed links.



Figure 1: Directed network of emojis. We create a connection from emoji to emoji in the order they appear in a tweet. This process is repeated for every tweet in the dataset. Then we accumulate all the sub-networks extracted from tweets into a main directed network of emojis.

A directed network of emojis gives us an opportunity to study the collective usage of emojis on social media. Additionally, different sequences of emojis may reflect different feelings expressed by users. For example, someone tweeting "I loved ❤️ this place until that horrible 😔 incident happened", the meaning is different from another tweet such as "This place is sometimes horrible 😔, but I love ❤️ it anyway!". Note that the order of the emojis is related to the sentiment being expressed. In order to build the directed network, we connect each emoji to every subsequent one appearing in the same tweet. Figure 1 shows the process of making directed weighted links between emojis.

### 3.3 Community Detection Techniques and Evaluation Criteria

Since the emoji co-occurrence networks are weighted and directed, we used state-of-the-art algorithms that support networks with these features [6]. Table 2 describes the selected algorithms and their respective approach to identifying communities. For each emoji network that was constructed, all algorithms are applied and the characteristics of the communities found are then analyzed.

Table 2: Community Detection Algorithms

| Acronym | Name | Approach | Description |
|---|---|---|---|
| $\mathcal{IM}$ | InfoMap | Bottom-up | Builds a map of information flow in the network using a random walk. Finding a community is equivalent to minimizing the flow representation by applying a compression technique [14]. |
| $\mathcal{BM}$ | Stochastic Block Models | Top-down | Applies maximum likelihood estimation to infer the latent block division in the empirical network. Such inference is equivalent to the entropy minimization of the network ensemble [13]. |
| $\mathcal{LP}$ | Label Propagation | Hybrid | Based on belief propagation, where each node spreads its label to its neighbors. Convergence of labels uncovers the community structure [7]. |
| $\mathcal{LM}$ | Louvain modularity | Bottom-up | Works by optimizing network modularity, which is the tightness of node connectivity into modules/communities in the empirical network relative to a null model [4]. |

To begin, we examined the size characteristics of the communities found by these four algorithms. Then, we apply an unsupervised evaluation by computing the communities' conductance. The conductance $C$ of a community $k$, measures the ratio between the intragroup and intergroup connectivity of the communities [6] and is computed as shown in Equation 1.

$$C(k) = \frac{\sum_{i \in k, j \notin k} w_{ij}}{\sum_{i \in k, j} w_{ij}} \ , \tag{1}$$

where $w_{ij}$ is the weight of the link connecting nodes $i$ and $j$. In this sense, well-structured communities exhibit a higher volume of edges between nodes within a community compared to edges going to the outside of the community.

We also conduct a supervised evaluation using the idea of rank stability for community detection [9] which was proposed as a way to measure the homogeneity of a particular community using the attribute values of nodes within that community as follows:

$$E(n) = -\sum_{t=1}^{L} p_{kt}(n) \log_2[p_{kt}(n)] \ , \tag{2}$$

where $p_{kt}(n)$ is the proportion of nodes in community $k$ that are associated with attribute $t$ in their rank $n$. In this work, we only have one attribute ($n = 1$) which can be one of the six categories of emojis.

## 4 Results

The statistical and structural properties of extracted communities can vary depending on the community detection technique. For instance, the number of communities extracted on large-scale networks significantly vary depending on the community detection technique [9]. We organized

our results in three parts. First, we characterize two structural properties of major importance, namely, community size and conductance. Then, we characterize the communities of emoji networks with regards to emoji categories known a priori to shed light in the context of emojis as language. Lastly, we present how these macroscopic characteristics are related and how a careful exploration of such relationship has the potential to help us gain insights on the structure and function of emojis.

Overall, emoji networks exhibit a well-defined community structure with regards to their size which is slightly shifted depending on the data set (Figure 2, left). Exceptionally, $\mathcal{LP}$ appears to find communities with typically greater size, it identifies the least number of communities, whereas other techniques find ten times more communities; this result is consistent with previous work [9]. Although the distribution of nodes within communities is an important aspect when identifying groups of interrelated emojis, we also need to quantify the extent in which extracted communities exhibit desirable structural properties.

Despite the lack of a general definition of a community, the number of links running between nodes within the community (i.e., internal edges) should be larger than the number of links running from nodes within the community to nodes outside the community (i.e., external edges). Conductance extends such definition for weighted networks (Equation 1). The conductance of communities vary depending on the technique and data set (Figure 2, middle). $\mathcal{IM}$ and $\mathcal{LM}$ show a more similar conductance distribution when compared to $\mathcal{BM}$. Precisely, $\mathcal{BM}$ has a higher likelihood of identifying communities with greater typical conductance. Similarly, $\mathcal{IM}$ and $\mathcal{LM}$ are likely to identify communities with moderate values of typical conductance. Lastly, $\mathcal{LP}$ is likely to identify communities with lower conductance.

Besides analyzing the structural properties of communities, we are also interested on evaluating communities to gain insight on the usage of emojis as a language. Here, we apply an entropy based metric to explore the levels of meaning in emoji usage. Our assumptions are that the higher the diversity of emoji categories within a community, the higher the level of meaning conveyed by these emojis.

We carry out this analysis by assessing the extent in which emojis within a community resemble the official emoji categories, using the rank entropy of communities [9]. The rank entropy (Equation 2) varies depending on the employed technique and underlying data set (Figure 2, right). Overall, the highest rank entropies are exhibited by communities extracted using $\mathcal{BM}$ as well as communities extracted from data set $D_1$. Conversely, the lowest entropies are exhibited by communities extracted using $\mathcal{LP}$ as well as communities extracted from the random data set $D_6$. $\mathcal{LP}$ presents some exceptions to aforementioned statements.

We can characterize communities of emojis by their structural properties, such as community size and conductance, as well as by their resemblance with available metadata using the rank entropy. Besides the independent characterization of community size, conductance, and rank entropy, we can also examine the relationship between these characteristics and unveil additional properties of the extracted communities. Indeed, these characteristics are related to each other (Table 3). Overall, communities with greater size tend to be moderately associated with lower conductance and higher rank entropy; however, there is a lack of a significant association between conductance and rank entropy. After controlling for community size, conductance appears to be associated with rank entropy mainly depending on the underlying data set. For instance, it is moderate in data sets $D_1$ and $D_2$, and it is absent in the random data set $D_6$.

Figure 2: Community size ($S$, left), conductance ($C$, middle), and rank entropy ($E$, right) of emoji networks characterized across six data sets of different thematics $D_{1...6}$ using four community detection techniques, namely, infomap $\mathcal{IM}$, block models $\mathcal{BM}$, label propagation $\mathcal{LP}$, and louvain modularity $\mathcal{LM}$. For details, see [1].

## 5   Conclusion

Emojis are an emerging form of communication, and emoji networks can exhibit structural properties with significant detail on the function of emojis and how they relate to language. In this work, we investigate the hypothesis that emojis are a form of language by building networks of emojis co-occurring on social media posts and subsequently analyzing the diversity of their community

Table 3: Relationship between community size $S$, conductance $C$, and rank entropy $E$ as measured by Pearson correlation. Community size $S$ is negatively correlated with conductance $C$ and positively correlated with rank entropy $E$. Even after controlling for community size, there is a lack of significant correlation between conductance and rank entropy in all techniques, except for label propagation.($\hat{\rho}_{SE \cdot S} = 0.12$, $p > .1$). However, further looking such correlation in each data set across multiple techniques, conductance is positively correlated with rank entropy, which is strongest in $D_1$ and weakest in the random data set $D_6$.

| | $\mathcal{IM}$ | $\mathcal{BM}$ | $\mathcal{LP}$ | $\mathcal{LM}$ | $D1$ | $D2$ | $D3$ | $D4$ | $D5$ | $\mathbf{D6}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| $\hat{\rho}_{SC}$ | $-0.14$ | $-0.12$ | $-0.24$ | $-0.05$ | $-0.19$ | $-0.14$ | $-0.14$ | $-0.17$ | $-0.13$ | $\mathbf{-0.19}$ |
| $\hat{\rho}_{SE}$ | $0.41$ | $0.37$ | $0.56$ | $0.46$ | $0.36$ | $0.31$ | $0.36$ | $0.35$ | $0.34$ | $\mathbf{0.31}$ |
| $\hat{\rho}_{CE}$ | $-0.03$ | $0.01$ | $-0.03$ | $-0.03$ | $0.11$ | $0.12$ | $0.10$ | $0.05$ | $0.01$ | $\mathbf{-0.05}$ |
| $\hat{\rho}_{CE \cdot S}$ | $0.01$ | $0.05$ | $0.12$ | $-0.02$ | $0.20$ | $0.18$ | $0.16$ | $0.12$ | $0.06$ | $\mathbf{0.01}$ |

structure. To gain insights on how emojis are used, we compare the diversity of communities from specific topics such as politics and sports with that of random. We find that users tend to communicate on social media using emojis of different categories. In this sense, the Stochastic Block Models would be more suitable since it is capable of finding more diverse (i.e., higher rank entropy) and well-formed communities (i.e., higher conductance). Yet, other possibilities to build emoji networks remain to be explored such as those based on risk ratio, pointwise mutual information and $\Phi$-correlation.

## Acknowledgements

## References

[1] For supplemental figures and details visit: `https://osf.io/5j489/?view_only=3214c22087ca40caba43f0b6101f7e91`.

[2] Francesco Barbieri, Francesco Ronzano, and Horacio Saggion. What does this emoji mean? a vector space skip-gram model for twitter emojis. In *Language Resources and Evaluation conference, LREC, Portoroz, Slovenia*, 2016.

[3] Anna Doble. Uk's fastest growing language is... emoji, 2015.

[4] Nicolas Dugué and Anthony Perez. Directed Louvain : maximizing modularity in directed networks. Technical report, Université d'Orléans, November 2015.

[5] Halley Fede, Isaiah Herrera, SM Mahdi Seyednezhad, and Ronaldo Menezes. Representing emoji usage using directed networks: A twitter case study. In *International Workshop on Complex Networks and their Applications*, pages 829–842. Springer, 2017.

[6] Santo Fortunato. Community detection in graphs. *Physics Reports*, 486(3-5):75–174, 2010.

[7] Chris Gaiteri, Mingming Chen, Boleslaw Szymanski, Konstantin Kuzmin, Jierui Xie, Changkyu Lee, Timothy Blanche, Elias Chaibub Neto, Su-Chun Huang, Thomas Grabowski, Tara Madhyastha, and Vitalina Komashko. Identifying robust communities and multi-community nodes by combining top-down and bottom-up approaches to clustering. *Scientific Reports*, 5(1):16361, dec 2015.

[8] Julian Gottke. Instagram emoji study: Emojis lead to higher interactions, 2017.

[9] Ryan Hartman, Josemar Faustino, Diego Pinheiro, and Ronaldo Menezes. Assessing the suitability of network community detection to available meta-data using rank stability. In *Proceedings of the International Conference on Web Intelligence - WI '17*, pages 162–169, New York, New York, USA, 2017. ACM Press.

[10] Andrea Lancichinetti and Santo Fortunato. Community detection algorithms: A comparative analysis. *Physical Review E*, 80(5):056117, nov 2009.

[11] Xuan Lu, Wei Ai, Xuanzhe Liu, Qian Li, Ning Wang, Gang Huang, and Qiaozhu Mei. Learning from the ubiquitous language: an empirical analysis of emoji usage of smartphone users. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, pages 770–780. ACM, 2016.

[12] Petra Kralj Novak, Jasmina Smailović, Borut Sluban, and Igor Mozetič. Sentiment of emojis. *PloS one*, 10(12):e0144296, 2015.

[13] Tiago P. Peixoto. Hierarchical block structures and high-resolution model selection in large networks. *Physical Review X*, 4(1):1–18, 2014.

[14] M. Rosvall and C. T. Bergstrom. Maps of random walks on complex networks reveal community structure. *Proceedings of the National Academy of Sciences*, 105(4):1118–1123, jan 2008.

[15] SM Mahdi Seyednezhad and Ronaldo Menezes. Understanding subject-based emoji usage using network science. In *Workshop on Complex Networks CompleNet*, pages 151–159. Springer, 2017.

[16] Sanjaya Wijeratne, Lakshika Balasuriya, Amit Sheth, and Derek Doran. Emojinet: Building a machine readable sense inventory for emoji. In *International Conference on Social Informatics*, pages 527–541. Springer, 2016.

# The Interlayer Coupling Resolution Limit in Multilayer Networks

Michael Vaiana[1,2] and Sarah F. Muldoon[1,3]

[1] Department of Mathematics and Computational Data-Enabled Science and Engineering
Program, University at Buffalo – SUNY
[2] mvaiana@buffalo.edu
[3] smuldoon@buffalo.edu

## Abstract

The organization of communities in a multilayer network has important conse-
quences for network structure and function. It has been shown that in traditional net-
works, maximizing the modularity function as a form of community detection provides
a balance of speed and accuracy, and this method has been extended to the multilayer
framework. However, in the multilayer setting, the method of maximum modularity
has an additional parameter, $\omega$, which is not well understood. Here, we expose an
upper bound for $\omega$ beyond which maximum modularity can not detect certain com-
munities. We demonstrate how this upper bound extends the previously known single
layer resolution limit, and we show how this has non-trivial consequences in the mul-
tilayer setting. Our results indicate that multilayer modularity can fail to detect even
the strongest of community structures.

The multilayer modularity function, $Q$, is a quality function which determines how well a
partition groups nodes into communities. The function $Q$ is dependent on an interlayer coupling
parameter, $\omega$, that controls the strength of connections between layers. We show that there is an
upper bound, $\Omega$, such that if $\omega > \Omega$ then certain multilayer communities can not be detected. This
upper bound has several consequences for modularity in multilayer networks. When communities
are small relative to the degree of the network they are undetectable and thus $\Omega$ is an extension of
the previously known single layer resolution limit. Even when communities are formed by disjoint
cliques (the strongest possible community structure), in certain cases they remain undetectable
as communities when maximizing multilayer modularity. We give an explicit formula for $\Omega$ and
through this we establish a relationship between the parameter $\omega$ and the modularity resolution pa-
rameter $\gamma$. As far as we are aware, this is the first time a dependence between these two parameters
has been established. Our results quantify the ability of multilayer modularity to detect communi-
ties and show explicitly how this ability depends on the parameter $\omega$ thereby clarifying its role and
behavior in multilayer networks.

# Enforcement networks in environmental protection

Lawrence De Geest[1]

[1] Department of Humanities and Social Sciences
Wentworth Institute of Technology
degeestl@wit.edu

Common-pool resources (CPRs) are often protected from poaching by decentralized groups of users coordinating under limited information. Poachers can be deterred with costly sanctions, but efficiently allocating sanctions is a challenge because it requires cooperation between group members. Benefits to cooperation are large, but since deterrence is a public good, so too are incentives to free-ride.

This paper takes a network approach to study enforcement with incomplete monitoring. While there is a well-developed literature that studies CPR management in the absence of government regulation, this is the first paper to study CPR management and protection from a shared threat, and the first to study deterrence in an environmental setting using methods from network science.

Using data from a economic experiment, we examine how groups of users coordinate to deter poachers they can perfectly or imperfectly monitor. In a repeated game, groups allocated effort towards harvesting a CPR and deterring poachers with monetary sanctions. Groups communicated to coordinate their decisions. Sanctions formed directed, weighted networks that captured the allocation of enforcement. We recover these networks from each round of play. We then use simulation and empirical methods to describe how these networks emerged and evolved, and how effective they were at deterring poachers.

We find that enforcement networks are non-random but sub-optimal compared to theoretical predictions. This is largely because contributions are driven by a subset of users. As a result, sanctions were ineffective at deterring poachers; estimates show a divergent effect of sanction events (edge formation) and sanction size (edge weight). Moreover, estimates of group payoffs show that that the number of active nodes – the number of group members contributing enforcement – is more important than the number of sanctions or the size of sanctions. Finally, long-run dynamics show that improved monitoring of outsiders leads to more stability, but not more efficiency.

A picture of failure emerges from these results. As groups struggled to coordinate, they dilute the ability for sanctions to deter poachers, and in turn, the incentive to protect the CPR from destruction.

Many environmental problems like CPR management and protection rely on government intervention. At the same time, users are often relied to protect CPRs because enforcement is scarce. What then is the optimal intervention? In the case of environmental protection, a major challenge in enforcement is improving coordination when monitoring of poachers is costly. These results suggest that feedback about the number of users engaging in enforcement may be a cost-effective way to increase protection.

## Acknowledgments

**Part II**

# Thursday April 12th 11:00am-12:15pm Contributed Talks 2: Nonlinear Dynamics (Tree House)

# Student: Stability Analysis of Delayed Complex Systems Under Uncertainty

Yiming Che, Changqing Cheng
Department of Systems Science and Industrial Engineering
State University of New York at Binghamton, Binghamton NY 13902

**Abstract:** Delay systems are ubiquitous in real-world physical and biological systems. Time delay in those systems typically gives rise to rich dynamic behaviors, from aperiodic to chaotic. The stability of such dynamic patterns is of particular interest for process control purposes. While stability analysis under deterministic conditions has been extensively studied, not too many works addressed the issue of stability under uncertainty. Nonetheless, estimate of parameters of complex systems inevitably leads to uncertainty, which will propagate from the input to the system output and could significantly affect stability of system dynamics. Indeed, the uncertainty could lead to divergent behaviors compared to the deterministic study. This is especially true when the system is at the bifurcation point. To this end, we investigated polynomial chaos (PC) to quantify the impact of uncertain parameters on the stability of delay systems. The case studies suggest that uncertainty quantification in delay systems provide richer information for system stability compared to deterministic study. In contrast to the robust yet time-consuming Monte Carlo or Hypercube Design method, PC approach achieves the same accuracy but only with a fraction of the computational overhead.

Keywords: Stability, Delay system, Polynomial chaos

## 1. Introduction

One of the grand challenges in the study of complex systems is the time delay effect [1]. In fact, time delay is ubiquitous and inherent in a plethora of physical and biological systems, from manufacturing to transportation, ecology and neural science, among others [2]. One of the common causes for time delay in feedback control systems is the finite communication speed, as sensors and actuators are rarely collocated. For instance, in steel rolling process, the thickness sensor is usually positioned at a certain distance away from the rolling gap, leading to measurement delay of the thickness, which is consequently used in a feedback control of the process [3]. In the study of glucose-insulin regulatory system, a time delay model is often used to describe the interaction of glucose, insulin and glucose-insulin mixture, to account for the delayed effect (~30-45 minutes) of insulin on glucose production [4] for more effective and personalized treatment. Notably, time delay can give rise to complicated dynamical behaviors. For instance, even simple regulatory gene circuit can dramatically alter the course of system evolution and exhibit rich steady-state behaviors, including limit cycle, aperiodic, weak/strong chaotic, and even intermittent patterns, when the time delay for transcription and translation is considered [5]. This delay-driven model is promising to pinpoint gene expressions associated with certain diseases.

Delay differential equations (DDEs) are among the prevailing tools to represent those dynamic systems with time delay. In DDEs, the evolution of the variable of interest depends on its state at present time $t$ as well as $t - \tau$ in the past, as indicated in Eq. (1) for a linear scalar system with a single discrete delay,

$$\dot{y}(t) = ay(t) + by(t - \tau) \qquad (1)$$

Similar to ordinary differential equations (ODEs), the solution to DDEs can be derived via the characteristic equations. Here, the linear scalar DDE in Eq. (1) has transcendental characteristic equation

$$-\lambda + a + be^{-\lambda\tau} = 0 \tag{2}$$

where $\lambda$ is called the characteristic root or eigenvalue. The solution set is often referred to as the spectrum, which bears the biomarkers of the underlying complex systems. Remarkably, the presence of exponential term $e^{-\lambda\tau}$ in Eq. (2) leads to an infinite number of possible values of $\lambda$, consequently an infinite number of solutions. In other words, the underlying dynamics is embedded in an infinite-dimensional phase space. Thus, this delay term renders solutions of DDEs differ from that of ODEs in a striking manner. It crucially affects the behavior of complex systems, leading to complicated trajectory of dynamics and even chaotic motion [6], and poses a tremendous challenge to local stability analysis of equilibria of such systems, one of the key issues in dynamic systems.

While stability analysis of dynamic systems with delay has been extensively studied [7], the performance under uncertainty has not been well explored. Nonetheless, just like any real-world complex systems, those delay systems are not immune to a wide range of uncertainty in initial and boundary conditions as well as model parameter calibration. There is a pressing need to quantify such uncertainty for reliable feedback control or robust process optimization. This paper investigated the stability of delayed systems under uncertainty using polynomial chaos. The focus of this present research is only on the single discrete delay.

## 2. Stability Analysis of DDEs

Stability of DDEs have been extensively studied in literature, using a variety of approaches, including Laplace transform [8], Lyapunov functions [9], perturbation analysis [10], Lambert W functions [11], semi-discretization [12], and Galerkin approximation [13]. Of those, Lambert W function is of particular interest for first order DDEs, whose characteristic equation has the explicit form of $z \mapsto ze^z$ for scalar argument $z \in \mathbb{C}$. As such, the eigenvalues can be solved in a straightforward way using the Lambert W function, which is defined as the multivalued inverse of $z \mapsto ze^z$. It has an infinite number of branches $W_k(z) \in \{w \in \mathbb{C}: z = we^w\}, k \in \mathbb{Z}$ [14], corresponding to the infinite number of solutions for the DDE in Eq. (1).

Concretely, the characterization equation in Eq. (2) can be reformulated as

$$(\lambda - a)\tau e^{(\lambda-a)\tau} = b\tau e^{-a\tau} \tag{3}$$

That is, $W_k(z) = b\tau e^{-a\tau}$ and $z = (\lambda - a)\tau$. Therefore, eigenvalues $\lambda_k = \frac{W_k(b\tau e^{-a\tau})}{\tau} + a$. The stability is determined upon the principal branch at $k = 0$, namely, $\lambda_0 = \frac{W_0(b\tau e^{-a\tau})}{\tau} + a$. The dynamics is unstable, if $\lambda_0$ lies to the right of the imaginary axis, i.e., $Re(\lambda_0) > 0$.

For the second order DDEs, matrix Lambert W function has been investigated, which registered comparable stability chart with the result obtained using bifurcation analysis [15]. Although the matrix Lambert W function is conceptually easy to understand, as it resembles the state transition matrix in linear ODEs, it is only restricted to a certain class of DDEs and the mathematical formulation is usually cumbersome. Thus, discretization and semi-discretization approaches have been developed to approximate $\lambda$ for higher-order DDEs. Among them, temporal finite element method has garnered abundant attentions, for example, in the applications of machining [16].

Nonetheless, most existing works do not consider the uncertainty associated with modeling and parametrization. Indeed, due to limitations in experimental study or calibration, or measurement error, the process parameters cannot be exactly specified and are often modeled as random

quantities in a probabilistic framework. The most straightforward way to quantify the impact of such uncertainty on stability is the Monte Carlo (MC) simulation. It is a brute force model in that it relies on large samples from the underlying random distribution. The system behavior is then evaluated according to the mean response of those sampled realization, and it is often infeasible due to the overwhelming computational overhead involved. While Latin Hypercube design [17] tends to optimize the sampling process in MC, it is still a sampling-based approach, and only marginally relieve the computational cost. On the other hand, polynomial chaos (PC) expansions have arisen as an efficient alternative to represent stochastic quantities as spectral expansions of orthogonal polynomials.

The PC is based on the original Wiener's theory of homogeneous chaos [18]. That is, a random variable with finite second-order moment can be expressed as a convergent series of polynomials in a sequence of random variables. Remarkably, it offers fast exponential convergence rate, and is a cheap alternative to MC simulations.

## 3. Stability Analysis of with uncertainty

In this section, we present PC using two case studies, first and second order DDEs: The first-order delay logistic equation, widely used to model the growth of population with limited resources; second-order delay vibration equation, commonly seen in the metal cutting field.

According to the Wiener-Askey scheme [19], PC expansion with Hermite polynomial basis or Hermite chaos expansion has been used effectively to solve uncertainty with Gaussian inputs, and Jacobian for beta distribution, and so on. The polynomial bases are orthogonal in that

$$\langle \phi_m, \phi_n \rangle = \int \phi_m(\xi) \phi_n(\xi) \rho(\xi) d\xi = \Delta_{mn} \delta_{mn}, \quad m, n \in \mathbb{N} \tag{4}$$

where the $\delta_{mn} = \begin{cases} 0, & m \neq n \\ 1, & m = n \end{cases}$ is the Kronercker delta, $\rho(\xi)$ is the weighting function for the corresponding distribution. The integration in Eq. (4) is oftentimes approximated using Gaussian quadrature rules. With the designated quadrature points and weights $\{(w^{(i)}, \xi^{(i)}): i \in \mathbb{N}\}$,

$$\Delta_{mn} = \sum_{i=1}^{N} w^{(i)} \phi_m(\xi^{(i)}) \phi_n(\xi^{(i)}) \tag{5}$$

3.1 Delayed logistic equation

A biological population with sufficient resources (e.g., food, space to grow and no threat from predators) grows at a rate proportional to the population size. Specifically, when a population is far below the carrying capacity of the ecosystem, it tends to grow exponentially. However, most populations are constrained by environmental capacity. Further, the feedback about the limitations usually comes with a delay due to various factors such as generation and maturation periods, differential resource consumption rate, and response time to the changing environments (e.g., competing species) [20]. Hutchinton [21] incorporated the effect of delays into the logistic equation, to model a variety of complicated dynamic behaviors, from limit cycle to chaotic behaviors, as indicated in Eq. (6). Here, $a$ denotes the rate of increase and $K$ the carrying capacity of the ecosystem. Our primary interest here is in the local stability of the population. For simplicity, we assume $K = 1, a = \frac{1}{2}$.

$$\dot{y}(t) = ay(t)\left(1 - \frac{y(t-\tau)}{K}\right) \rightarrow \dot{y}(t) = \frac{1}{2}y(t)(1 - y(t-\tau)) \tag{6}$$

Apparently, there are 2 equilibrium points $y(t) = 0$ and $y(t) = 1$, and $y = 0$ is unstable. Then linearize Eq. (6) for $y = 1$, and substitute $x(t) = y(t) - 1$, we have

$$\dot{x}(t) = -\frac{1}{2}x(t-\tau) \tag{7}$$

which is amenable to the Lambert-W function, whose principle eigenvalue is given as

$$\lambda_0 = \frac{1}{\tau}W_0\left(-\frac{1}{2}\tau\right) \tag{8}$$

Thus, the local stability of $x(t)$ and $y(t)$ depends on the time delay $\tau$. In fact, the bifurcation occurs at $\tau = \pi$. In the scenario of uncertainty analysis, we assume a symmetric beta distribution of $\tau$ centered at 3.15 to quantify our uncertainty. The result should end up with $\lambda_0 \approx 0$.

**Remark 1** In statistics, the standard beta distribution $\xi' \in [0,1]$, while in polynomial chaos, it is defined as $\xi \in [-1,1]$. Thus, we need to make transformation $\xi = 2\xi' - 1$ for PC.

Correspondingly, we have $\tau = 3.15 + 0.5\xi$, where $\xi' \sim Beta(2,2)$ is the standard statistical beta distribution. The main reason we use beta rather than Gaussian distribution is because of the compact support in the former.

To study the impact of random $\tau$ on $\lambda_0$ and the local stability, we implemented MC, LHD and PCE. For MC, we draw 10,000 samples from the distribution of $\tau$, and evaluate $\lambda_0$ at each realization. LHD is one kind of optimized MC, and similar accuracy was obtained with only half of the computational cost, as summarized in Table . In the PC setting, $\lambda_0$ was represented as

$$\lambda_0 \approx \sum_{i=0}^{P} c_i\phi_i(\xi) \tag{9}$$

where $\phi_i(\xi)$ is Jacobi chaos polynomials. To obtain the coefficient $c_i$ here, we adopt the Stochastic Galerkin method. Here, a Galerkin projection is used to minimize the error of the truncated expansion and the resulting set of equations can be solved to obtain the expansion coefficients. Provided sufficient smoothness conditions are met, PC estimates of uncertainty converge exponentially with the order of the expansion. As indicated in Eq. (10),

$$\langle \lambda_0, \phi_j(\xi) \rangle = \langle \sum_{i=0}^{P} c_i\phi_i(\xi), \phi_j(\xi) \rangle \quad \text{for } j = 0,1,2 \dots, P \tag{10}$$

where $\langle \cdot, \cdot \rangle$ represents the inner product operator with respect to $\xi$. Due to the orthogonal property, $\langle \phi_i(\xi), \phi_j(\xi) \rangle = \Delta_{mn}\delta_{ij}$. This leads to $c_j = \frac{\langle \lambda_0, \phi_j(\xi) \rangle}{\langle \phi_i(\xi), \phi_j(\xi) \rangle}$, and $\langle \lambda_0, \phi_j(\xi) \rangle$ is numerically evaluated via Gaussian quadrature. The first two moments of $\lambda_0$ are thus characterized by the coefficients as $E[\lambda_0] = c_0$, $\sigma[\lambda_0] = \sum_{i=1}^{P} c_i^2$.

Table 1 Comparison of deterministic method, MC, LHD and PC

| Method | $E[\lambda_0]$ | $\sigma[\lambda_0]$ | Sample Number |
|---|---|---|---|
| Deterministic | $6.0387 \times 10^{-4}$ | 0.0000 | 1 |

| MC | $-9.3869 \times 10^{-4}$ | 0.0164 | $1 \times 10^4$ |
| --- | --- | --- | --- |
| LHD | $-1.1000 \times 10^{-3}$ | 0.0164 | $5 \times 10^3$ |
| PC | $-1.1000 \times 10^{-3}$ | 0.0190 | 20 |

The comparison of MC, LHD and PC is summarized in Table 1. Notably, under uncertainty, negative $E[\lambda_0]$ is derived, which is in sharp contrast to the positive value in the deterministic case. Thus suggests that the average dynamic behaviors under uncertainty could be drastically distinct from that under nominal and/or deterministic conditions. This further necessitates the uncertainty quantification for robust process control and monitoring. Also note that, PC registered the accuracy as MC, but only with a fraction of the computational overhead.

3.2 Delayed vibration equation

The relative vibration between cutting tool and the workpiece is a normal phenomenon associated with metal cutting or machining processes. As indicated in Figure 1 (a), the cutting tool continuously removes material from a rotating workpiece in a face turning process. One notorious detriment in such process is the fierce relative vibration, often characterized by violent tool vibration, premature tool wear, loud noise, and inferior part surface. This phenomenon is also referred to as regenerative chatter or self-excited vibrations in the machining process, triggered by time delay.

This self-excited vibration is usually modeled as a second-order DDE in the feed direction, which contributes most to the surface formation [22], as in Eq. (11)

$$\ddot{y}(t) + 2\zeta\omega_n\,\dot{y}(t) + \omega_n^2 y(t) = -kb\big(f_0 - y(t) + y(t-\tau)\big) \qquad (11)$$

Here, $y$ denotes the displacement of the cutting tool, $\zeta$ the damping ratio, and $\omega_n$ the natural frequency of the relative vibration. This is a forced-vibration model, as on the right-hand side, the cutting force is expressed as proportional to the instantaneous cutting area, and $k$ is the force coefficient. The cutting area is determined by depth of cut $b$ and actual feed rate $f_1 = f_0 - y(t) + y(t-\tau)$. As shown in Figure 1 (b), the deviation of actual feed $f_1$ from $f_0$ is due to the tool vibration, or the displacement of the cutting tool at the current as well as previous revolution. The time delay $\tau$ here is the time period of one revolution.



Figure 1 Schematic diagram of facing turning: (a) cutting tool continuously remove material from the rotating workpiece; (b) regenerative chatter model

Predictive models have been explored to generate stability lobe chart or stability regions for a variety of combination of speeds (i.e., $\tau$) and cutting depths $b$ to eschew the costly and time-consuming trial and error alternatives [23] Thereby, mapping the area of stability as a function of process parameters, i.e., depth of cut and spindle speed, can guide the process design to avoid those

detrimental effects of chatter and elevate the production rate.

In our previous work, we investigated stability of Eq. (11) using temporal finite element method [16]. Specifically, we divided the time delay or the time period of one revolution into $M$ elements, and approximate $y(t)$ in each time interval $T = \frac{\tau}{M}$ as $y(t) = \sum_{i=1}^{4} a_{ji}^n S_i(\sigma)$, where $\sigma \in [0, T]$ is the local time within the $j^{th}$ element in $n^{th}$ revolution and $S_i(\sigma)$ is the cubic Hermite polynomial, which possesses the property of $C^0$ and $C^1$ continuity. Substitute $y(t)$ into Eq. (11) and invoke deterministic Galerkin approach, we have

$$Na^n = Pa^{n-1} + N^{-1}Q \qquad (12)$$

where the matrix N, P and Q are derived from the two Galerkin projections on each element, thus of size 2M×2M, as indicated in our previous work [16]. This has reduced the stability analysis to a finite-dimensional Floquet transition matrix problem, see [16] for details. The stability of this equation requires $\lambda_G$<1, where $\lambda_G$ denotes the maximum absolute eigenvalue of the transition matrix $G = N^{-1}P$. This is equivalent to that the principal eigenvalue of its characteristic equation satisfies $Re(\lambda_0)$<0, as $\lambda_0 = \frac{\max(Re(\log \lambda_G))}{\tau}$. The stability chart is indicated in Figure 2 (a). Here, the shaded area represents the unstable cutting, whilst the rest is for stable cutting. Clearly, under this deterministic condition, there is a stark boundary between these two different dynamic behaviors. In the scenario of stability, the system will eventually return to the equilibrium under small perturbation.



Figure 2: (a) The boundary between stability and instability when $\zeta$ is deterministic with $\zeta = 0.02$. (b)(c)(d) The stability derived from PC, MC, LHS respectively, and the color bar stands for the probability of eigenvalue getting zero.

In the case of uncertainty, without loss of generality, we assume standard beta distribution $\xi' \sim \text{Beta}(2,2)$, and in statistics $\xi' \in [0,1]$. The same transformation $\xi = 2\xi' - 1 \in [-1,1]$ is used for PC. Then we have $\zeta = 0.02\xi + 0.01$, which is centered at 0.02, representing the typical values of damping ratio in machining. At each parameter setting $(\tau, b)$, we derived $\lambda_G(\zeta)$ from TFEM using random samples of $\zeta$, and adopted PC to represent the uncertain $\lambda_G(\zeta)$

$$\lambda_G(\zeta) \approx \sum_{i=0}^{p} c_i \phi_i(\xi) \tag{13}$$

We follow the same stochastic Galerkin procedure as explained in case study 1 to solve the coefficients here. The comparison of stability chart obtained using MC, LHD and PC is shown in Figure 2 (b)-(d). Here, the stability chart under uncertainty presents the probability that the process is stable. Thus, industrial practitioners can avoid the area with low probability of stability to optimize process design. We note that while the three methods generate similar probabilistic stability chart, PC is far more efficient, as only 5 samples are used in PC compared to 1000 in MC.

## 4. Discussion and Conclusions

In this present study, we investigated stability of delay complex systems under uncertainty. While the stability of delay systems in deterministic scenarios has been extensively studied, very few research works explored the stability of dynamic behaviors under uncertainty. We applied PC to quantify the uncertainty propagation in such systems in two case studies, for one and second order DDEs. For the first order DDE, the stability is determined upon the principle eigenvalue of Lambert W function, and a stochastic Galerkin framework is applied directly on the principal eigenvalue expression. The result showed that uncertain parameter leads to negative mean of the principal eigenvalues, suggesting stability compared to the positive principal eigenvalues in deterministic case. For the second order DDE, there is no Lambert W function available, or it is limited only to specific forms of DDE. Thus, TFEM is used to approximate the eigenvalue. Again, the Galerkin algorithm is applied upon the TFEM that connects the uncertain parameter and the eigenvalue. In both cases, PC achieves the same accuracy as the robust MC method, but only at a fraction of the computational cost. Overall, our study suggests that the stability could be sensitive to parameters variation, particularly for the process dynamics at the bifurcation point. Thus, special treatment is needed for the robust control or system optimization.

## References

[1] W. Just, A. Pelster, M. Schanz, and E. Schöll, *Delayed complex systems: an overview*. The Royal Society, 2010.

[2] M. Lakshmanan and D. V. Senthilkumar, *Dynamics of nonlinear time-delay systems*. Springer Science & Business Media, 2011.

[3] J. Pittner and M. A. Simaan, *Tandem cold metal rolling mill control: using practical advanced methods*. Springer Science & Business Media, 2010.

[4] J. Li, Y. Kuang, and B. Li, "Analysis of IVGTT glucose-insulin interaction models with time delay," *Discrete and Continuous Dynamical Systems Series B*, vol. 1, no. 1, pp. 103–124, 2001.

[5] Y. Suzuki, M. Lu, E. Ben-Jacob, and J. N. Onuchic, "Periodic, Quasi-periodic and Chaotic Dynamics in Simple Gene Elements with Time Delays," *Scientific Reports*, vol. 6, p. 21037, Feb. 2016.

[6] S. Boccaletti, "The synchronized dynamics of complex systems," *Monograph series on*

*nonlinear science and complexity*, vol. 6, pp. 1–239, 2008.

[7]   J. Nilsson, B. Bernhardsson, and B. Wittenmark, "Stochastic analysis and control of real-time systems with random time delays," *Automatica*, vol. 34, no. 1, pp. 57–64, 1998.

[8]   H. Rezaei, S.-M. Jung, and T. M. Rassias, "Laplace transform and Hyers–Ulam stability of linear differential equations," *Journal of Mathematical Analysis and Applications*, vol. 403, no. 1, pp. 244–251, 2013.

[9]   F. H. Clarke, Y. S. Ledyaev, and R. J. Stern, "Asymptotic stability and smooth Lyapunov functions," *Journal of differential Equations*, vol. 149, no. 1, pp. 69–114, 1998.

[10]  P. Kokotović, H. K. Khalil, and J. O'reilly, *Singular perturbation methods in control: analysis and design*. SIAM, 1999.

[11]  S. Yi, P. Nelson, and A. Ulsoy, "Delay differential equations via the matrix Lambert W function and bifurcation analysis: application to machine tool chatter," *Mathematical Biosciences and Engineering*, vol. 4, no. 2, p. 355, 2007.

[12]  T. Insperger and G. Stépán, *Semi-discretization for Time-delay Systems: Stability and Engineering Applications*, vol. 178. Springer Science & Business Media, 2011.

[13]  C. P. Vyasarayani, "Galerkin approximations for stability of delay differential equations with time periodic delays," *Journal of Computational and Nonlinear Dynamics NOVEMBER*, vol. 10, pp. 061008-1, 2015.

[14]  C. Hwang and Y.-C. Cheng, "A note on the use of the Lambert W function in the stability analysis of time-delay systems," *Automatica*, vol. 41, no. 11, pp. 1979–1985, 2005.

[15]  S. Yi, P. Nelson, and A. Ulsoy, "Delay differential equations via the matrix Lambert W function and bifurcation analysis: application to machine tool chatter," *Mathematical Biosciences and Engineering*, vol. 4, no. 2, p. 355, 2007.

[16]  P. V. Bayly, J. E. Halley, B. P. Mann, and M. A. Davies, "Stability of interrupted cutting by temporal finite element analysis," *Journal of Manufacturing Science and Engineering*, vol. 125, no. 2, pp. 220–225, 2003.

[17]  M. Stein, "Large sample properties of simulations using Latin hypercube sampling," *Technometrics*, vol. 29, no. 2, pp. 143–151, 1987.

[18]  A. Desai and S. Sarkar, "Analysis of a nonlinear aeroelastic system with parametric uncertainties using polynomial chaos expansion," *Mathematical Problems in Engineering*, vol. 2010, 2010.

[19]  D. Xiu and G. E. Karniadakis, "The Wiener--Askey polynomial chaos for stochastic differential equations," *SIAM journal on scientific computing*, vol. 24, no. 2, pp. 619–644, 2002.

[20]  P. Turchin, *Complex population dynamics: a theoretical/empirical synthesis*, vol. 35. Princeton University Press, 2003.

[21]  G. E. Hutchinson, "Circular causal systems in ecology," *Annals of the New York Academy of Sciences*, vol. 50, no. 1, pp. 221–246, 1948.

[22]  C. Kan, C. Cheng, and H. Yang, "Heterogeneous recurrence monitoring of dynamic transients in ultraprecision machining processes," *Journal of Manufacturing Systems*, vol. 41, no. C, 2016.

[23]  F. A. Khasawneh, B. P. Mann, T. Insperger, and G. Stépán, "Increased stability of low-speed turning through a distributed force and continuous delay model," *Journal of Computational and Nonlinear Dynamics*, vol. 4, no. 4, p. 041003, 2009.

# Spatiotemporal Modeling of Nonlinear Dynamics in Cardiac Electrical Activities

Rui Zhu[1] and Hui Yang[2]

[1] Complex System Monitoring, Modeling and Analysis Laboratory
The Pennsylvania State University, University Park
rzz45@psu.edu
[2] Complex System Monitoring, Modeling and Analysis Laboratory
The Pennsylvania State University, University Park
huy25@psu.edu

## Abstract

As the leading cause of death, cardiac diseases account for more than 30% of mortalities in the US. It is urgent to improve the early detection of life-threatening cardiac events by monitoring cardiac electrical activities. Body Sensor Networks (BSNs) have emerged as a key technology for monitoring spatiotemporal dynamics of complex systems. Electrocardiogram (ECG) sensing is increasingly integrated into BSNs for monitoring highly nonlinear and nonstationary cardiac activities. Current ECG systems deploy 224 sensors which are approximately uniformly distributed on the body surface. Large groups of ECG sensors provide rich information on spatiotemporal cardiac dynamics that are not fully available in traditional 3-lead or 12-lead ECGs. However, the approximately uniform distribution of sensors considers little about the non-uniform potential distribution on the heart surface. In this paper, we propose to model spatiotemporal dynamics of cardiac electrical activities with the information recorded by a parsimonious set of sensors. Notably, these sensors are optimally distributed on the body surface to capture a complete picture of spatiotemporal cardiac dynamics. Experimental results show that the spatiotemporal modeling effectively and efficiently detects nonlinear dynamics of cardiac electrical activities and demonstrate that the parsimonious set of sensors can achieve the performance of current ECG sensing systems.

# Heterogeneous Recurrence Analysis of Vectorcardiogram Signals

Ruimin Chen[1], Farhad Imani[1] and Hui Yang[1]

[1] Complex System Monitoring, Modeling and Analysis Laboratory
The Pennsylvania State University
rxc91@psu.edu

## Abstract

Recurrence is one of the most common phenomena in nature. During the past few decades, the theory of recurrence analysis has been studied and made meaningful progress in different domains (i.e. manufacturing, finance, and biology). Typically, the recurrence plot and recurrence quantification analysis became a useful tool for visualizing and quantifying the recurrence behaviors in complex systems. However, traditional recurrence methods are based on homogeneous recurrences, which defines all recurrence and non-recurrence states with black and white colors in recurrence plots respectively. While dynamic system monitoring concerns more about the abnormal recurrences and variations between the recurrence states. As a result, the recurrence analysis cannot capture all nonstationarity and stochastic irregularity in nonlinear dynamics. For example, most of the traditional recurrence analysis on Vectorcardiogram signals assumes the homogeneous recurrence behavior. While heterogeneous recurrences are more related to the variations of recurrence states in terms of state properties and the evolving dynamics. Very little work has been done to study the heterogeneous recurrence behaviors for monitoring and control the Vectorcardiogram signals. There is an urgent need for developing the new tool for characterization and quantification of heterogeneous recurrences on Vectorcardiogram. The objective of this paper is to utilize a novel heterogeneous recurrence analysis approach to study the hidden pattern in Vectorcardiogram. For this purpose, first, we reconstruct the state space and recursively partition it into a hierarchical structure. Next, we extract the quantifiers by fractal representation. Then, we select the important quantifiers by a nonparametric method and utilize the multivariate monitoring method to detect anomaly behaviors in Vectorcardiogram. Experimental results show that the proposed approach not only captures heterogeneous recurrence patterns in the fractal representation but also successfully monitors the changes in the dynamics of a Vectorcardiogram.

# Extraction and classification of convectively coupled equatorial waves through eigendecomposition of Koopman operators

Joanna Slawinska[1] and Dimitrios Giannakis[2]

[1] Department of Physics, University of Wisconsin-Milwaukee, slawinsk@uwm.edu
[2] Courant Institute of Mathematica Sciences, New York University, dimitris@cims.nyu.edu

## Abstract

We present recently developed data-driven technique for dynamical systems and its application for spatiotemporal pattern retrieval of propagating coherent structures in the Earth's tropical atmosphere.

We study spatiotemporal patterns of convective organization using a recently developed technique for feature extraction and mode decomposition of spatiotemporal data generated by ergodic dynamical systems (Giannakis 2017). The method relies on constructing low-dimensional representations (feature maps) of spatiotemporal signals using eigenfunctions of the Koopman operator. This operator is estimated from time-ordered unprocessed data through a Galerkin scheme applied to basis functions computed via the diffusion maps algorithm. Koopman operators are a class of operators in dynamical systems theory that govern the temporal evolution of observables. They have the remarkable property of being linear even if the underlying dynamics is nonlinear, and provide, through their spectral decomposition, natural ways of extracting intrinsic coherent patterns and performing statistical predictions.

We apply this approach to brightness temperature data from the CLAUS archive and extract a multiscale hierarchy of spatiotemporal patterns on timescales spanning years to days, including dominant intraseasonal mode of tropical variability (MJO; Madden and Julian, 1972) but also traveling waves on temporal and spatial scales characteristic of convectively coupled equatorial waves (CCEWs; Kiladis et al. 2009). In particular, we examine if the activity of these coherent structures is modulated by low-frequency atmospheric and oceanic variability. We discuss various properties of waves in our hierarchy of modes, focusing in particular on their across-scale interactions and temporal evolution. As an extension of this work, we discuss the deterministic and stochastic aspects of the variability of these modes.

## References

[1] Giannakis, 2017: Data-Driven Spectral Decomposition and Forecasting of Ergodic Dynamical Systems. *Appl. Comput. Harmon. Anal.*, 2017.

[2] R. Madden and P. Julian, 1972: Description of global-scale circulation cells in the tropics with a 40-50 day period. *J. Atmos. Sci.*, **29**, 1109-1123.

[3] G. N. Kiladis, M. C. Wheeler, P. T. Haertel, K. H. Straub, and P. E. Roundy, 2009: Convectively coupled equatorial waves. *Rev. Geophys.*, **47**, 2.

# Investigating Inherent Timescales of Variability in Blazars Using Structure Functions

Giridhar Nandikotkur

School of Natural Sciences
Fairleigh Dickinson University, Teaneck NJ
giri@fdu.edu

## Abstract

Active Galactic Nuclei (AGN) are galaxies with a supermassive black hole at their center, and jets emanating perpendicular to the plane of rotation of galaxy in both directions. The jets consisting of blobs of plasma moving at relativistic speeds, radiate across a broad range of wavelengths from radio through-X-ray to high-energy gamma-rays. NASA's X-ray satellite, Rossi X-ray Timing Explorer (RXTE), observed the AGN (named Mrk421) numerous times during its 16 years of operation from 1996-2012. The X-ray data are available at various sampling-intervals. Active Galactic Nuclei occasionally flare up when large blobs of plasma are ejected into the jet stream. These flares are followed by a prolonged quiescent-states with medium to low activity. The timescales over which flux (energy received per second per unit area of the detector) varies can be tied to the size of the emitting region of the emission. One of the interesting questions surrounding the X-ray emission is if the quiescent state has an intrinsic time-scale of variability. The light curves (flux vs. time) have been analyzed using structure functions to investigate the presence of inherent timescales of variability and to identify characteristics of variability in flux. We have also extracted the fractal dimension of the light curve at various timescales in order to characterize the statistical nature of the physical process causing the variability. The results of the analysis will be discussed in the context of underlying physical mechanisms responsible for emission of X-rays.

**Part III**

# Thursday April 12th 3:50pm-5:05pm Contributed Talks 3: Data-Driven Approaches (Symposium Hall)

# A Quantitative Image Analysis Method for Nuclei in Subcellular Microscopic Image

Mohammad A. Al-Mamun[1] and Shakeeb R. Chowdhury[2]

[1]Department of Epidemiology of Microbial Diseases, Yale School of Public Health, New Haven, CT, USA and [2]University of Liberal Arts Bangladesh, Dhaka, Bangladesh

## Abstract

Detecting nucleus and analyzing it from microscopic images with different quantitative methods direct us to efficient prognosis and treatments for the complex disease like cancer. In this work, we segmented and classified different nuclei shapes from hundreds of subcellular microscopic nuclei image obtained from a public database called Protein Atlas. We used marker controlled watershed segmentation method to segment the nuclei and classified them using three classifying techniques: support vector machine (SVM), random forest and convolutional neural network (CNN). The CNN outperforms while compared to the other classification techniques. Overall, this method presents a novel tool for segmenting and classifying the subcellular nuclei images.

## 1.Introduction

The nucleus is one of the most prominent cellular organelles, yet surprisingly little has been done to detect and classify them based on their shape, size and texture. Major challenges of subcellular nuclei analysis come from the complexity of the dynamic intracellular environment and include different sizes and shapes and varying internal structure of nuclei other structures. Previously proposed methods are not adequate as most of the them are domain and problem specific which warrant an improved techniques nuclei segmentation and classification for general use [1].

## 2.Materials and Method

We used an improved version of marker controlled watershed segmentation proposed in [1] and included the flooding process used before applying the watershed algorithm. The output of segmentation method provided us individual nucleus image with 150×150 resolution. After that we labelled the image based on the shape: normal (>80% circular) and abnormal (<80% circular). Then, we extracted 21 features: area, centroid, convexarea, ecentricity, majoraxislength, minoraxislength, orientation, solidity, epithelial form factor, perimeter, density index; texture features: mean, variance, skewness, kurtosis, angular second moment, entropy, inverse different moment, correlation, and contrast. These features were fed into two classifiers: SVM and Random forest. For CNN, each nuclei image was fed to the input layer of a sequential network.

## 3.Results and Discussion

For evaluating the classifiers performance, we validate normal and abnormal categories by a domain expert (Table 1) and calculate the classification accuracy based on the extracted features for each segmented nucleus (Table 2). One of the limitation of this study is that the z-slice stacks was not available in the current database to get the projection of whole nucleus shape. The results show, that CNN outperforms in classifying the nuclei. In future, we will improve the classification accuracy by using the CAFEE CNN architecture.

Table 1. Description of segmentation results

| Total samples | 199 |
|---|---|
| Total nuclei | 3356 |
| Perfectly segmented nuclei | 2833 |
| Under-segmented | 112 |
| Over-segmented | 411 |

Table 2. Classification accuracy

| Classifiers | Accuracy= [*]TP+TN/TP+FP+FN+TN |
|---|---|
| SVM | 81.1 |
| RandomForest | 84.7 |
| CNN | 86.2 |

## References

[1] Al-Mamun, Mohammad A., et al. "A quantitative image analysis for the cellular cytoskeleton during in vitro tumor growth." *Expert Systems with Applications* 92 (2018): 39-51.

# Optimal Clusterings

M.T. Bassett, B.P. Newton, J.C. Schlessinger,
W.R. Pulleyblank, P.K. Kuiper, S.A. Lynch, R.E. Miller,
S.T. Morse, J.D. Pleuss, T.B. Russell, J.W. Roginski

Network Science Center, Department of Mathematical Sciences
United States Military Academy, West Point, NY

## Extended Abstract

Clustering is a fundamental problem encountered in data analytics and unsupervised learning. The basic problem is to partition records which share certain attributes into sets (clusters) such that records which are "similar are in the same set and records which are "different are in different sets. For example, health care fraud detection systems often have a first stage wherein provider billing records are grouped into clusters of similar records. These clusters are then analyzed to find statistical outliers, which are provided to investigators as targets for possible audits.

Many clustering heuristics have been proposed, but there is no broadly accepted metric for measuring the quality of solutions produced by these methods. We analyze the so-called *Condorcet metric* which computes penalties for pairs of records in the same cluster which are not sufficiently similar and pairs of records in different clusters which are not sufficiently different. A clustering for which the sum of these penalties over all pairs of records is minimized is defined to be optimum.

We apply a number of standard clustering methods to a variety of datasets and compute the Condorcet metric for the resulting clusterings. This enables us to rank them in order of this metric. We also consider the problem of finding a mathematically optimal clustering of the datasets according to this metric, by formulating the problem as a mixed integer programming problem. This enables us to measure how close the clusterings produced by various methods are to the theoretical optimum for several different data sets.

The clustering problem can be formulated as a problem on a network whose nodes correspond to the records and for which every pair of nodes is joined by an edge with a weight corresponding to the difference between the records. We analyze the effect of "sparsifying the network by removing selected edges. By using a variation of the Delaunay triangulation from computational geometry we are able to reduce the number of edges in the network from $O(n^2)$ to $O(n)$. Not only does this give us a much smaller network, but it also results in construction of clusterings which avoid some of the problems encountered by application of the Condorcet metric to all possible pairs of record.

In addition to developing the theory and methods of this approach, we present computational results on a variety of data sets, including arrest data, demographic data and state political data.

# Is space a word, too?

Jake Ryland Williams[1] and Giovanni C. Santia[2]

College of Computing and Informatics, Drexel University
[1]jw3477@drexel.edu; [2]gs495@drexel.edu

**Abstract**

Common definitions of "word" describe either a speech sound/sound-series or segment of written discourse appearing between space and/or punctuation. They *must* communicate meaning without being divisible into smaller units capable of independent use. But what about cues left between, often neither spoken nor written? Before zero was formalized as a number, thousands of years of mathematics proceeded using placeholders. Similarly, irrational and transcendental numbers were first approached proximally, through geometric relationships and equations of their accepted, rational counterparts. So, the number concept was successively broken and extended from discretization, to ratio, to quantity, to its final logical completion by numbers poorly described as "imaginary" or "complex". These concepts emerged as intellectual curiosities or useful tools, notably before their acceptance into formalism. Considering this analogy, we study the potential inclusion of some useful linguistic cues into a common status with words. We observe evidence through a *proximal* analysis of textual objects accepted as words, finding possibility for space et al. to share this status.

Our approach to the question answering if space is a word focuses on Mandelbrot's modification of Zipf's law for word frequencies in conjunction with Simon's model for how language is generated. For a text of $R$ distinct words: $w_1, \cdots, w_R$; ranked, $r = 1, \cdots, R$, descending, according to their frequencies of occurrence: $f(w_1) \geq \cdots f(w_R)$, Zipf's law is a power-law distribution: $\hat{f}_Z(w_r) \propto r^{-\gamma}$; $\gamma \geq 0$, which Mandelbrot refined by a horizontal translation, $k$: $\hat{f}_{ZM}(w_r) \propto (r + k)^{-\gamma}$; $\gamma \geq 0, k > -1$. Mandelbrot rooted the existence of $k$ both empirically and in the theoretical optimization of communication against its cost of transmission. However, for Zipf's law or Mandelbrot's refinement a word's rank intuitively describes the number of words of similar or more-severe frequency. So, accepting a Mandelbrot translation of $k \approx 2$, we might posit the existence of two unresolved words of exceptional frequency. Natural candidates might thus be space and punctuation, etc., whose place in Zipf's law we investigate, finding evidence.

Our empirical work points towards the possibility of space's inclusion in Zipf's law for words, but in an unexpected way. Unlike a simple $f \propto r^{-1}$ relationship, the Mandelbrot refinement that space fits into entails outlier behavior at small ranks, with $k < 0$. We thus find a connection with the language generation model proposed by Simon. Though put aside for its crudeness, it re-emerged as the preferential attachment mechanism by Barabasi for complex networks nearly 50 years later. Unlike the networks case, Simon's model for words entails the generation of exceptionally high-frequency words through a functional form in close agreement with Mandelbrot's $k$ and space's empirical frequencies: $\hat{f}_S(w_r) \propto (r - \theta)^{-\theta}$; $\theta \in (0, 1)$. While this form is more restrictive than Mandelbrot's, we ask if Simon's model acts as a foundational mode for language generation.

# Emergence of synchrony and chimera in data-driven models of brain dynamics

Kanika Bansal[1,2,3], Timothy D. Verstynen[4], Jean M. Vettel[1,5,6], and Sarah F. Muldoon[3,7]

[1] Human Sciences, U.S. Army Research Laboratory, Aberdeen Proving Grounds, MD
[2] Department of Biomedical Engineering, Columbia University, New York, NY
[3] Department of Mathematics, University at Buffalo, SUNY, Buffalo, NY
[4] Department of Psychology, Carnegie Mellon University, Pittsburgh, PA
[5] Department of Bioengineering, University of Pennsylvania, Philadelphia, PA
[6] Department of Psychological and Brain Sciences, University of California, Santa Barbara, CA
[7] CDSE Program, University at Buffalo, SUNY, Buffalo, NY

## Abstract

The human brain is a complex dynamical system which functions as regional neuronal populations interact to produce spatiotemporal patterns of coherent activity. Often, one can observe the formation of separate domains of coherence and incoherence – a state referred to as 'chimera' in the complex systems framework. We constructed data-driven models of brain dynamics and studied the emergent patterns of synchrony and chimera that are driven by the underlying structural connectivity of the brain. Our results indicate that different regions of the brain are structurally constrained to produce different patterns, in a way that is predictive of their functional (cognitive) role.

## Methods and main results

To study the patterns of human brain activity, we built virtual brain network models based on the observed structural connectivity obtained from diffusion weighted imaging data across a cohort of thirty subjects. For each subject, we modeled regional brain dynamics using Wilson-Cowan oscillators, coupled through each individual's structural connectivity. We then sequentially studied the effects of applying regional activation across brain regions and across the cohort of subjects. As the activation spreads, we analyzed synchronized and de-synchronized populations within the brain by calculating network synchrony. We observed three emergent states, namely coherent, chimera, and metastable. Emergence of these states was found to follow the network connectivity of the activated region. Here, we discuss the importance of these states in the cognitive functioning of the human brain. Further, we describe how the chimera framework can be used to aid our understanding of the structural organization of the human brain and individual cognitive variability.

## Acknowledgments

# Inter-Session Reproducibility of Brain Functional Network Structure in Resting and Task States

Johan Nakuci[1], Javier O. Garcia[2,3], Nick Wasylyshyn[2,3], James C. Elliot[4],
Matthew Cieslak[4], Barry Giesbrecht[4], Scott T. Grafton[4], Jean M. Vettel[2,3,4], Sarah F. Muldoon[1,5]

[1] Neuroscience Program, University at Buffalo, Buffalo, NY
[2] U.S. Army Research Laboratory, Aberdeen Proving Ground, MD
[3] Department of Bioengineering, University of Pennsylvania, Philadelphia, PA
[4] Department of Psychological and Brain Sciences, University of California, Santa Barbara, CA
[5] Department of Mathematics and CDSE Program, University at Buffalo, Buffalo, NY

## Abstract

Network analysis has provided new and important insights into the function of complex systems such as the brain. In the field of neuroscience, both structural and functional human brain networks have been constructed from diffusion MRI, functional MRI (fMRI) and electro/magnetoencephalography (E/MEG), and the topology of these networks has provided new insight into brain function. However, these studies are often limited to analysis of a single recording or scanning session for each subject, or only compared between two subsequent sessions, raising questions about inter-session reproducibility. Given the prevalence of inter-subject variability in network statistics calculated across multiple subjects for data from a single session, it is important to characterize the reproducibility of network statistics across sessions to assess the sensitivity of network measures to between subject variation. Here, we investigate the reproducibility of network statistics measured in functional brain networks by comparing inter-subject and inter-session variability across multiple measures of network structure. We find that intra-subject reproducibility is higher than inter-subject reproducibility, suggesting sensitivity to individual differences in brain network structure. Further, intra-subject reproducibility depends on the frequency band of interest and modality of recording. Our findings suggest that one must be cautious when interpreting variability in measures of resting- and task-state functional connectivity.

## Methods

We study functional brain networks derived from 20 subjects with 8 sessions of rest and task combined fMRI-EEG recordings. The fMRI networks were constructed using Wavelet-Coherence across five frequency bins from 0.06-0.1Hz. Wavelet-Coherence creates networks without having to accommodate negative values as seen in traditional correlations analysis. For EEG data, networks were constructed using the de-biased weighted Phase Lag Index (dwPLI) across five frequency bands: $\delta$ (1-3Hz), $\theta$ (3-7Hz), $\alpha$ (8-13Hz), $\beta$ (15-30Hz), and $\gamma$ (30-50Hz). The dwPLI is an appropriate measure for assessing functional connectivity in EEG data given it is less susceptible to volume conduction. For the resulting networks, we calculated the Clustering Coefficient, Path Length, Small-World Topology, Synchronizability, Assortivity, Hierarchy, and Global- and Cost-Efficiency. The Interclass Correlation Coefficient (ICC) was used to assess inter- and intra-subject variability.

**Part IV**

# Thursday April 12th 3:50pm-5:05pm Contributed Talks 4: Network Dynamics (Tree House)

# Strategic Topology Switching for Security of Multi-Agent Systems

Yanbing Mao[1], Emrah Akyol[1] and Ziang Zhang[1]

[1] Department of Electrical and Computer Engineering, Binghamton University–SUNY,
Binghamton, NY, 13902 USA.
Email: {ymao3, eakyol, zhangzi}@binghamton.edu.

## Abstract

This paper analyses strategic topology switching for multi-agent systems under the class of "zero-dynamics" attack. We first study the detectability of aforementioned attacks for under switching topology. Based on this analysis, we then propose a strategic topology-switching algorithm that optimally changes topology and prevents the entire class of zero-dynamic attacks. The proposed approach outperforms prior work in the sense that i) we do not impose any constraints on the size of the attacker (misbehaving agents) set, ii) one monitoring output is sufficient.

## 1 Introduction

Security concerns regarding networked cyber-physical systems pose an existential threat to their wide-deployment, see e.g., such as Stuxnet malware attack and Maroochy Shire Council Sewage control incident [1]. The "networked" aspect exacerbates the difficulty of detecting and preventing aforementioned attacks since centralized measurement (sensing) and control are not feasible for such large-scale systems [2], and hence requires the development of decentralized approaches which are inherently prone to attacks. Recently, a special class of stealthy attacks, namely the "zero-dynamics" have gained interest. Here, the attacker's goal is two fold: i) not being detected by the system by keeping the monitoring outputs unaltered (hence the name "stealthy"), ii) manipulating the system to accept false data (e.g., aggregation results), which are significantly different from the actual data, see e.g., [3]. While developing defense strategies for such zero-dynamics attacks have recently gained interest [2, 4, 5, 7], the space of solutions is yet to be thoroughly explored. Significant drawbacks of prior work are that they constrain the connectivity of network topology and the size of the misbehaving-agent set [2, 4, 5] or require the knowledge of attack beginning time at the defender side for attack detection [2, 4, 5, 7]. The main objective of this work is to remove such constraints and unrealistic assumptions by utilizing a new approach for attack detection: intentional topology switching.[1] However, before using the dynamic topologies to reveal zero-dynamics attack, the question that *whether the dynamic changes can destroy the system stability in the absence of attacks* must be investigated.

The impact of network topology on the stability of networked systems has been recently studied. For example, Mao et al. [11] find that the second-order consensus can be achieved under

---

[1]Topology changes in a multi-agent network have traditionally been considered a problem dictated by nature, to be mitigated by the designer. Here, we intentionally change the topology to detect the zero-dynamics type of attacks.

certain dwell time of topology-switching signals; for the first-order multi-agent systems studied in [8], stability under switching topology shows that the dynamic connected topologies do not undermine the agents' ability of reaching consensus. We use this stability result of [8] to derive the strategic topology-switching algorithm that reveals zero-dynamics attacks.

## 2 Preliminaries and Problem Formulation

### 2.1 Preliminaries

#### 2.1.1 Notation

For a set $\mathcal{V}$, $|\mathcal{V}|$ denotes the cardinality (i.e., size) of the set. In addition, for a set $\mathcal{K} \subseteq \mathcal{V}$, $\mathcal{V}\backslash\mathcal{K}$ denotes the complement set of $\mathcal{K}$ with respect to $\mathcal{V}$. $\mathbb{R}^n$ and $\mathbb{R}^{m \times n}$ denote the set of $n$-dimensional real vectors and the set of $m \times n$-dimensional real matrices, respectively. Let $\mathbb{C}$ denote the set of complex number. $\mathbb{N}$ represents the set of the natural numbers and $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$. Let $I$ be the identity matrix with compatible dimension. $\mathbf{1}_n \in \mathbb{R}^n$ and $\mathbf{0}_n \in \mathbb{R}^n$ denote the vector with all ones and the vector with all zeros, respectively. The superscript '$\top$' stands for matrix transpose.

#### 2.1.2 Graph Theory

The interaction among $n$ agents is modeled by an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{1, 2, \cdots, n\}$ denoted the set of $n$ agents and $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$ is the set of edges of the graph $\mathcal{G}$. An undirected edge in $\mathcal{G}$ is denoted by $a_{ij} = (i, j) \in \mathcal{E}$, where $a_{ij} = a_{ji} = 1$ if agents $i$ and $j$ interact with each other, and $a_{ij} = a_{ji} = 0$ otherwise. Assume that there are no self-loops, i.e., for $\forall i \in \mathcal{V}$, $a_{ii} \notin \mathcal{E}$. A path is a sequence of connected edges in a graph. A graph is a connected graph if there is a path between every pair of vertices.

**Lemma 1** *[9] If the undirected graph $\mathcal{G}$ is connected, then its Laplacian $\mathcal{L} \in \mathbb{R}^{n \times n}$ has a simple zero eigenvalue (with eigenvector $\mathbf{1}_n$) and all its other eigenvalues are positive and real.*

**Lemma 2** *[10] The Laplacian of a path graph $P_n$ has the eigenvalues as*

$$\lambda_k = 2 - 2\cos(\frac{(k-1)\pi}{n}), k = 1, \cdots, n. \tag{1}$$

### 2.2 Problem Formulation

For simplicity, let $\mathcal{K} = \{1, 2, \cdots\} \subseteq \mathcal{V} = \{1, \cdots, n\}$ denote the set of misbehaving agents, and $\mathcal{M} = \{1, 2, \cdots\} \subseteq \mathcal{K}$ denote the set of outputs used by monitors. The considered multi-agent system with its outputs under attack can be described by

$$\dot{x}_i(t) = \sum_{i=1}^{n} a_{ij}^{\sigma(t)}(x_j(t) - x_i(t)) + g_i(t)\mathrm{u}_i(t - \kappa), i \in \mathcal{V} \tag{2a}$$

$$y_i(t) = x_i(t) + g_i(t)\mathrm{u}_i(t - \kappa), i \in \mathcal{M} \tag{2b}$$

43

where $x_i(t) \in \mathbb{R}$ is the agent state, $\sigma(t) \in \mathfrak{S}$ is the topology switching signal under attack; $g_i(t)$ is agent $i$'s attack signal; $a_{ij}^{\sigma(t)}$ is the element of the coupling matrix that describes the activated $\sigma^{\text{th}}(t)$ topology of undirected communication network; $\mathrm{u}_i(t - \kappa)$ is the unit step function: $\mathrm{u}_i(t - \kappa) = 1$ if $t \geq \kappa$ and $i \in \mathcal{K}$, and $\mathrm{u}_i(t - \kappa) = 0$ otherwise; $\kappa \geq 0$ is the attack-beginning time.

The agent under attack is usually named the misbehaving agent.

**Definition 1** *[5] Consider multi-agent system (2). An agent $i \in \mathcal{V}$ is misbehaving if there exists a time $t$ such that $g_i(t)\,\mathrm{u}_i(t - \kappa) \neq 0$.*

We next make the following assumptions on the attacker and defender.

**Assumption 1** *The considered attacker:*

- *has the knowledge of topology-switching sequences, including the switching times and the activated topologies at each switching times;*

- *can modify the initial conditions arbitrarily if the attack-beginning time is the initial time;*

- *can attack the multi-agent system infinitely over infinite time.*

**Assumption 2** *The considered defender has no knowledge of the attack-beginning time nor the misbehaving agents.*

The following lemma states that arbitrarily topology switching among connected topologies does not undermine the agents' ability of reaching consensus in the absence of attacks.

**Lemma 3** *[8] For any arbitrary switching signal $\sigma(t)$ of connected topologies, the solution of the multi-agent system (2) in the absence of attacks globally asymptotically converges to Ave(x(0)) (i.e., average-consensus is reached).*

## 3   Detectability of Zero-Dynamics Attack

This section will show the advantages of strategic topology switching in revealing zero-dynamics attack.

### 3.1   Zero-Dynamics Attack

**Definition 2** *(Zero-Dynamics Attack [2]) Consider the following two systems:*

$$\begin{cases} \dot{p}(t) = Ap(t) + Bg(t) \\ \bar{y}(t) = Cp(t) + Dg(t) \end{cases} \tag{3}$$

$$\begin{cases} \dot{q}(t) = Aq(t) \\ \tilde{y}(t) = Cq(t) \end{cases} \tag{4}$$

*where $p(t), q(t) \in \mathbb{R}^n$, $\bar{y}(t), \tilde{y}(t) \in \mathbb{R}^m$ with $m \leq n$, $g(t) \in \mathbb{R}^o$ with $o \leq n$, $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times o}$, $C \in \mathbb{R}^{m \times n}$ and $D \in \mathbb{R}^{m \times o}$. The attack signal $g(t) = g(\kappa)e^{\lambda(t-\kappa)}$, $t \geq \kappa$ with $\kappa \geq 0$, is a zero-dynamics attack if $g(\kappa) \neq \mathbf{0}_o$ and $\lambda \in \mathbb{C}$ satisfy*

$$\begin{bmatrix} \lambda I - A & B \\ -C & D \end{bmatrix} \begin{bmatrix} p(\kappa) - q(\kappa) \\ -g(\kappa) \end{bmatrix} = \begin{bmatrix} \mathbf{0}_n \\ \mathbf{0}_m \end{bmatrix}. \tag{5}$$

*The state and output of system (3) satisfy*

$$\bar{y}(t) = \tilde{y}(t), \forall t \geq 0 \tag{6}$$

$$p(t) = q(t) + (p(\kappa) - q(\kappa)) e^{\lambda(t-\kappa)}, \forall t \geq \kappa. \tag{7}$$

**Remark 1** *The resulted state (7) shows that through choosing the parameter $\lambda$, the attacker can achieve its attack objectives:*

- $\mathrm{Re}(\lambda) > 0$*: unstable system;*

- $\mathrm{Re}(\lambda) = 0, \mathrm{Im}(\lambda) \neq 0$*: oscillating;*

- $\mathrm{Re}(\lambda) < 0$*: modifying the steady-state value.*

*The resulted output (6) means the undetectable property of zero-dynamics attack.*

### 3.2 Strategy on Switching Topologies

The following strategy on switching topologies can reveal zero-dynamics attack without constraint on the size of misbehaving-agent set.

**Theorem 1** *Consider the multi-agent system under attack (2). If the strategy satisfies*

**(r1)** *after the attack begins at $\kappa \geq 0$, the first topology-switching signal is time-dependent, i.e, $\sigma(t_k)$ where $k = \underset{r \in \mathbb{N}}{\arg\min} \{t_r > \kappa\}$, is time-dependent;*

**(r2)** *the difference of Laplacian matrices of the two consecutive topologies activated at switching times $t_k$ and $t_{k+1}$ has distinct eigenvalues,*

*then without constraint on the misbehaving-agent set and with only one output for monitor, the strategy can reveal zero-dynamics attack.*

**Proof 1** *The proof can be finished by contradiction. We assume the multi-agent system (2) has zero-dynamics attack. The multi-agent system (2) in the absence of attacks can be written as*

$$\dot{\tilde{x}}_i(t) = \sum_{i=1}^{n} a_{ij}^{\hat{\sigma}(t)} (\tilde{x}_j(t) - \tilde{x}_i(t)), i \in \mathcal{V} \tag{8a}$$

$$\tilde{y}_i(t) = \tilde{x}_i(t), i \in \mathcal{M} \tag{8b}$$

*where $\hat{\sigma}(t)$ is the topology switching signal in the absence of attacks. Let $\breve{x}_i(t) = x_i(t) - \tilde{x}_i(t)$, $\mathcal{L}_1$ and $\mathcal{L}_2$ denote Laplacian matrices.*

*The rest of proof follows the same line of reasoning as the proof of Theorem 1 given in [12]. Following nearly the same steps, we have*

$$g(\kappa) - \mathcal{L}_1 \breve{x}(\kappa) + \lambda \breve{x}(\kappa) = \mathbf{0}_n, \tag{9}$$

$$g(\kappa) - \mathcal{L}_2 \breve{x}(\kappa) + \lambda \breve{x}(\kappa) = \mathbf{0}_n. \tag{10}$$

*It is readily seen from (9) and (10) that $(\mathcal{L}_1 - \mathcal{L}_1)\breve{x}(\kappa) = \mathbf{0}_n$. The strategy (r2) in Theorem 1 implies $\mathcal{L}_1 - \mathcal{L}_1$ has distinct eigenvalues. By Lemma 1 we can conclude that $\mathcal{L}_1 - \mathcal{L}_1$ has properties: (i) zero is one of its eigenvalues with multiplicity one; (ii) the eigenvector that corresponds to the eigenvalue zero is $\mathbf{1}_n$. Thus, the solution of $(\mathcal{L}_1 - \mathcal{L}_1)\breve{x}(\kappa) = \mathbf{0}_n$ is $\breve{x}(\kappa) = \mathbf{0}_n$. Substituting $\breve{x}(\kappa) = \mathbf{0}_n$ into either (9) or (10) yields the same result as $g(\kappa) = \mathbf{0}_n$. Therefore, by Definition 2 that there is no zero-dynamics attack, so a contradiction occurs, which completes the proof.*

### 3.3 Strategic Topology Switching Algorithm

Based on Lemma 1 and Theorem 1, the proposed strategic topology-switching algorithm is described by Algorithm 1.

---

**Algorithm 1:** Strategic Topology-Switching Algorithm

---

**Input:** Dwell time $\tau$, topology set $\mathfrak{S}$ that includes at least two connected topologies that the difference of their Laplacian matrices has distinct eigenvalues.

1 **while** $x(t_{k-1}) \neq Ave(x(0))$ **do**

2      Switch the connected topology of network (2a) from $\sigma(t_{k-1})$ to $\sigma(t_k)$ at $t_k$, such that $\sigma(t_k) \neq \sigma(t_{k-1})$, and $\mathcal{L}_{\sigma(t_k)} - \mathcal{L}_{\sigma(t_{k-1})}$ has distinct eigenvalues;

3      Run the multi-agent system (2) until the time $t_{k+1} = t_k + \tau$;

4      Update the topology-switching time: $t_{k-1} \leftarrow t_k$;

5      Update the topology-switching time: $t_k \leftarrow t_{k+1}$.

6 **end**

---

**Theorem 2** *Consider the multi-agent system under attack (2). If the topology-switching signal is generated by Algorithm 1, then the following properties hold.*

**(i)** *In the absence of attacks, the agents can achieve the average consensus by Lemma 1.*

**(ii)** *Without constraint on the misbehaving-agent set and with only one monitoring output, Algorithm 1 is able to reveal zero-dynamics attack.*

**Proof 2** *The topology set $\mathfrak{S}$ provided to Algorithm 1 includes at least two topologies topology, and at each topology-switching time, Algorithm 1 considers only connected topologies, hence by Lemma 3 we can conclude the property (i).*

*Lines 3–5 implies that all the topology-switching signals generated by Algorithm 1 are time-dependent, thus (r1) in Theorem 1 is satisfied. Line 2 in Algorithm 1 means the difference of Laplacian matrices of every two consecutive topologies has distinct eigenvalues, thus (r2) in Theorem 1 is satisfied. Therefore, by Theorem 1 we can conclude the property (ii).*

**Remark 2** *As $0 \leq \frac{(k-1)\pi}{n} < \pi, \forall k = 1, \cdots, n$, Lemma 2 implies that the Laplacian matrix of a path graph has distinct eigenvalues. Therefore, we can conclude (r2) in Theorem 1 can be easily satisfied if the graph generated by the difference of two graphs is a path graph.*

Table 1: Candidate Topologies

| Index $\sigma(t)$ | $\mathbf{a}_{12}^{\sigma(t)}$ | $\mathbf{a}_{13}^{\sigma(t)}$ | $\mathbf{a}_{14}^{\sigma(t)}$ | $\mathbf{a}_{23}^{\sigma(t)}$ | $\mathbf{a}_{24}^{\sigma(t)}$ | $\mathbf{a}_{34}^{\sigma(t)}$ |
|---|---|---|---|---|---|---|
| 1* | 1 | 0 | 0 | 1 | 0 | 1 |
| 2* | 1 | 1 | 0 | 0 | 1 | 1 |
| 3* | 1 | 1 | 0 | 0 | 1 | 0 |

## 4 Simulation

The simulations on a multi-agent system with $n = 4$ agents will be presented to demonstrate the effectiveness of the proposed strategic topology-switching algorithm. In the simulation setting, the initial states are chosen as $x(0) = [1, 2, 3, 4]^\top$. For simplicity, let the attack-beginning time is the initial time, i.e., $\kappa = t_0 = 0$, which means the attacker can modify the initial conditions arbitrarily.

To convincingly verify the effectiveness of Algorithm 1 in revealing zero-dynamics attack, we consider the extremely bad situation:

- all the agent are misbehaving, i.e., $\mathcal{K} = \{1, 2, 3, 4\}$,

- only one output is available to monitor, let $\mathcal{M} = \{1\}$.

### 4.1 Zero-Dynamics Attack Design

First, we considered the topology set $\mathfrak{S} = \{2^*, 3^*\}$ where the representations of $2^*$ and $3^*$ are given in Table II. Obviously, the set $\mathfrak{S} = \{2^*, 3^*\}$ does not satisfy the strategy (r2) in Theorem 1. Thus, the attacker can easily design a zero-dynamics attack such that Algorithm 1, with only one output of position, cannot reveal it.

Let the attacker's goal is to attack the multi-agent system under Algorithm 1 to be unstable, while not being detected. Following the zero-dynamics attack design method (5) in Definition 2, one of its zero-dynamics attack strategies is easily designed as

- choose parameter $\lambda = 1$,

- introduce attack signal: $g(t) = [0, 4e^t, -2e^t, -3e^t]^\top$,

- modify the initial conditions: $x(0) = [1, 3, 2, 3]^\top$.

The trajectories of attack detection signal $r(t) = y_1(t) - \tilde{y}_1(t)$ and system states are given in Figure 1, which shows the attacker's goal of attacking the multi-agent system to be unstable while not being detected by the defender is achieved. Thus, using only one output, the designed zero-dynamics attack is not revealed under the topology set $\mathfrak{S} = \{2^*, 3^*\}$.

### 4.2 Reveal Zero-Dynamics Attack

Now we turn to the topology set $\mathfrak{S} = \{1^*, 2^*\}$ to real the attack. It can be verified that this topology set satisfies (r2) in Theorem 1, hence by Theorem 2 we can conclude that using only one output, the strategic topology-switching algorithm–Algorithm 1–is able to reveal the designed zero-dynamics attack.

Figure 1: Attack detection signal $r(t)$: the attack is not revealed; system states: the multi-agent system under attack is unstable.

The trajectory of attack detection signal $r(t) = y_1(t) - \tilde{y}_1(t)$ is shown in Figure 2, which illustrates that with all the agents being misbehaving, using only one output, Algorithm 1 succeeds in revealing zero-dynamics attack.



Figure 2: Attack-detection signal $r(t)$: using only one output, the designed zero-dynamics attack is revealed.

## 5   Conclusion

This paper studies strategic topology switching for detection of zero-dynamic attacks on multi-agent systems. The propose approach has the following advantages over the prior work:

- the algorithm has no constraint on the set of misbehaving agents, i.e., all the agents are allowed to be misbehaving;

48

- only one output is enough for monitor;

- the defender is not required to know the attack-beginning time.

## References

[1] A. A. Cárdenas, S. Amin, and S. Sastry. Research challenges for the security of control systems. *In HotSec*, pp. 1–6, 2008.

[2] F. Pasqualetti, F. Dörfler, and F. Bullo. Attack detection and identification in cyber-physical systems. *IEEE Transactions on Automatic Control*, 58(11): 2715–2729, 2013.

[3] B. Przydatek, D. Song, and A. Perrig. SIA: Secure information aggregation in sensor networks. *In Proceedings of the 1st international conference on Embedded networked sensor systems*, pp. 255–265, ACM, 2003.

[4] S. Sundaram and C. N. Hadjicostis. Distributed function calculation via linear iterative strategies in the presence of malicious agents. *IEEE Transactions on Automatic Control*, 56(7): 1495-1508, 2011.

[5] F. Pasqualetti, A. Bicchi, and F. Bullo. Consensus computation in unreliable networks: A system theoretic approach. *IEEE Transactions on Automatic Control*, 57(1): 90-104, 2012.

[6] A. Teixeira, D. Pérez, H. Sandberg, and K. H. Johansson. Attack models and scenarios for networked control systems. *Proceedings of the 1st international conference on High Confidence Networked Systems*, pp. 55–64, ACM, 2012.

[7] A. Teixeira, I. Shames, H. Sandberg, and K. H. Johansson. Revealing stealthy attacks in control systems. *In Communication, Control, and Computing (Allerton), 2012 50th Annual Allerton Conference on*, pp. 1806–1813, 2012.

[8] R. Olfati-Saber and R. M. Murray. Consensus problems in networks of agents with switching topology and time-delays. *IEEE Transactions on Automatic Control*, 49(9): 1520–1533, 2004.

[9] C. Godsil and G. Royle. *Algebraic Graph Theory*. New York: Springer-Verlag, 2001.

[10] D. Spielman. Spectral graph theory (lecture 5 rings, paths, and cayley graphs). *Online lecture notes. http://www.cs.yale.edu/homes/spielman/561/*, September 16, 2014.

[11] Y. Mao, E. Akyol and Z. Zhang. Strategic topology switching for security-Part I: consensus & switching times. *arXiv preprint arXiv:1711.11183*, 2017.

[12] Y. Mao, E. Akyol and Z. Zhang. Strategic topology switching for security-Part II: detection & switching topologies. *arXiv preprint arXiv:1711.11181*, 2017.

# Target Control in Logical Models Using the Domain of Influence of Nodes

**Gang Yang** [1,*]**, Jorge G. T. Zañudo**[1,3,4] **and Réka Albert** [1,2]

[1]*Department of Physics, Pennsylvania State University, University Park, Pennsylvania 16802, USA*
[2]*Department of Biology, Pennsylvania State University, University Park, Pennsylvania 16802, USA*
[3]*Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, Massachusetts 02215, USA*
[4]*Eli and Edythe L. Broad Institute of MIT and Harvard, 415 Main Street, Cambridge, Massachusetts 02142, USA*

Correspondence*:
Gang Yang
gzy105@psu.edu

## ABSTRACT

Dynamical models of biomolecular networks are successfully used to understand the mechanisms underlying complex diseases and to design therapeutic strategies. Network control, and its special case of target control, is a promising avenue toward developing disease therapies. In target control it is assumed that a small subset of nodes is most relevant to the system's state and the goal is to drive the target nodes into their desired states. An example of target control would be driving a cell to commit to apoptosis (programmed cell death). From the experimental perspective, gene knockout, pharmacological inhibition of proteins and providing sustained external signals are among practical intervention techniques. We identify methodologies to use the stabilizing effect of sustained interventions for target control in Boolean network models of biomolecular networks. Specifically, we define the domain of influence of a node (in a certain state) to be the nodes (and their corresponding states) that will be ultimately stabilized by the sustained state of this node regardless of the initial state of the system. We also define the related concept of the logical domain of influence of a node, and develop an algorithm for its identification using an auxiliary network that incorporates the regulatory logic. This way a solution to the target control problem is a set of nodes whose domain of influence can cover the desired target node states. We perform greedy randomized adaptive search in node state space to find such solutions. We apply our strategy to in silico biological network models of real systems to demonstrate its effectiveness.

Keywords: Target Control, Boolean Network, Biological Network, Domain of Influence, Logical Modelling, Network Dynamics

The full text is available at http://biorxiv.org/cgi/content/short/243246v1 .

# From Consensus to Polarization of Opinions in Complex Contagion

Vítor V. Vasconcelos[1] and Flávio L. Pinheiro[2]

[1] Department of Ecology and Evolutionary Biology, Princeton University, Princeton (NJ), USA

[2] The MIT Media Lab – Massachusetts Institute of Technology, Cambridge (MA), USA

The study of how opinions, innovations, behaviors, and knowledge spread has long been a central topic in physical, social, and ecological sciences. In the past, these have been commonly studied through simple contagion processes, processes in which information flows through the contact of two individuals. The inability for simple contagion models to account for the plethora of dynamical patterns observed in the real world, such as the polarization of opinions, has led to a search for additional mechanisms. Recent empirical evidence suggests that different matters spread in different ways, namely, that some require a dependence on the whole neighborhood of an individual to propagate. In that sense, the process of information acquisition requires reinforcement from multiple contact sources. This phenomenon became known as Complex Contagion. Although widely investigated in the literature of cascading effects, complex contagion has only recently received some attention in the context of population dynamics, i.e., when multiple competing opinions co-evolve over time in a population.

Complex Contagion has been commonly modeled by a process of fractional thresholds. This implies that there is a well-defined threshold fraction of neighbors needed for an opinion/idea/innovation to be adopted by an individual. Dynamically, this results in a deterministic process that either percolates through the system or that becomes contained to a few elements of the system. Under this context, it was found that complex contagion spreading is speeded up by clustering of individuals in populations (triangular closures) but that a modular structure of the population can halt the propagation of an opinion. Indeed, the study of opinion dynamics has examined under which conditions a consensus is formed. Typical questions involve the time to consensus and how likely is it for a new opinion to invade a population.

Here, we introduce a new class of complex contagion processes inspired by recent empirical findings in the literature of innovation and knowledge diffusion. [1] We consider different opinions coevolving in a population with potentially asymmetric properties of contagion. We assume that the probability of an opinion to spread to a new individual grows as an arbitrary power of the density of neighbors that already share that same opinion. We explore analytically and computationally the properties of this model in well-mixed and structured populations. Namely, we test well-mixed, homogeneous random, random, scale-free, and modular networks of influence. We show these populations span a dynamical space that exhibits patterns of polarization, consensus, and dominance. We map these patterns to topologically equivalent ones found in the literature of evolutionary games of cooperation. We find that these dynamical properties are robust to different population structures. Finally, we show how modular topologies can create different dynamics and additional dependences on the initial configuration of opinions. Our results are general and of relevance not only for the study of opinions and ideas but also when considering propagation in more abstract networks derived from data, like that of product complexity and product adoption by countries as well as others. [2,3]

[1]     M. Karsai, G. Iniguez, K. Kaski, and J. Kertész, J. R. Soc. Interface **11**, 20140694 (2014).

[2]     C. A. Hidalgo, B. Klinger, A.-L. Barabási, and R. Hausmann, Science (80-. ). **317**, 482 (2007).

[3]     D. Hartmann, M. R. Guevara, C. Jara-Figueroa, M. Aristarán, and C. A. Hidalgo, World Dev. **93**, 75 (2017).

# A Possibilistic Programming Approach to Handle Dynamics of Global Supply Chain Network

Alireza Fazli Khalaf, Yong Wang [1]

Department of Systems Science and Industrial Engineering, Binghamton University, Binghamton, NY 13902, USA

**Abstract**

Global trading nowadays has stimulated the need for designing global supply chains to efficiently satisfy the increasing customer demand across the world. The dynamic nature and complexity of global supply chain networks cause a high degree of uncertainty that can adversely affect all decision levels in long-term planning periods. In this regard, this paper aims to propose an efficient approach for designing a global supply chain network. Unlike previous studies, we consider comprehensive parameters of the global supply chain in the network design and apply an effective possibilistic programming method to cope with uncertainty of the parameters. This method enables decision makers to adjust the level of conservatism in achieving the desired outcome. Numerical examples and the computational results are presented to show applicability and effectiveness of the proposed approach.

## 1. Introduction

In recent decades, the fast growth in globalization has forced companies and corporations to establish their own global supply chain network (GSCN). They have expanded their supply chain network in multiple countries in order to benefit from tariff and trade concessions, cheap labor, low capital subsidies, and reduced logistics costs [1].

A high degree of uncertainty originated from the dynamic and complex nature of supply chain networts (SCN) will greatly influence their overall performance as well as their design [2]. Stochastic programming is a method that can be applied to handle uncertainty of parameters of SCN models. However, in many real situations there is not enough data to determine a reliable statistical distribution of uncertain parameter in the design problem. An alternative approach to stochastic programming is fuzzy set theory. The fuzzy set theory provides a framework to tackle the epistemic uncertainty, including fuzzy coefficients for lack of knowledge as well as flexibility in constraints and the target value of the goals simultaneously. Fuzzy programming and in particular, the possibilistic programming (PP) method employs experts experience and knowledge to determine possibility distribution of uncertain parameters [3]. Globalization considerations in SCN models such as exchange rate of currencies, tariffs and duties, and value-added tax have been an important part of GSCN studies [1]. However, many models are simplified to reduce the complexities by ignoring the dynamic and uncertain nature of globalization parameters.

This paper proposes an efficient multi-echelon GSCN model that integrates both strategic network design decisions and tactical planning decisions to optimize at different decision levels simultaneously and consequently avoid sub-optimalities. To handle different sources of uncertainties that adversely affects reliability of the model outcomes, the PP method will be applied to cope with unavailability or incompleteness and imprecise nature of data. Globalization related parameters and their dynamic nature, which are disregarded in related studies in the literature, will be considered in the model beside demand and supply uncertainties.

---

[1] Corresponding author. Tel: +1 607-777-3845; Fax: +1 607-777-4094. Address: 4400 Vestal Pkwy E, Binghamton, NY 13902-6000, USA. Email: yongwang@binghamton.edu.

The rest of this paper is organized as follows. The related literature is reviewed and commented in Section 2. The GSCN design problem is formulated in Section 3. The PP method is described and applied on the formulated model in Section 4. Illustrative examples and sensitivity analysis on the optimal design results are elaborated in Section 5. Finally, conclusions and future research guidelines are provided in section 6.

## 2. Literature Review

In the area of designing GSCNs, researchers have focused on various topics and issues related to global trade concepts. Bing et al. [4] designed a reverse GSCN for household plastic waste recycling industry. Global parameters such as tax and duties were not considered in the proposed model. [1,5,6] have also studied GSCN design and considered global parameters. However, they are mainly deterministic models. In addition, it is rare in the previous research work to comprehensively handle parameter uncertainties for GSCN design.

In the area of traditional SCN design, many researchers have worked on approaches to tackle uncertainty of parameters. Some researchers have focused on two-stage stochastic programming models [7,8]. As it was previously noted, stochastic programming models has some deficiencies and researchers has eliminated them by application of the PP approach. [2] proposed a bi-objective closed-loop SCN model that minimizes total costs beside maximizing responsiveness of network. They extended a multi-objective PP approach to deal with uncertainty of parameters. Vahdani et al. [9] expanded [2] by considering reliability concepts in the model to achieve robust results under disruptions. Hatefi et al. [10] improved the model of [9] by application of the p-robustness measure to tackle disruption risks. It should be mentioned that many recent SCN design models [11,12] are deterministic and ignore uncertainty of the network. Based on the literature review, this paper is concerned with GSCN design while using the PP method to handle uncertainty of parameters.

## 3. Problem Formulation

The proposed GSCN design model in this paper is an integrated multi-echelon and single product model that optimizes strategic and tactical decisions simultaneously. As it is presented in Figure 1, plants procure needed raw materials from different suppliers to produce final products. Then, products are transferred to customer zones via distribution centers. Notably, to model global trading concepts, tariff related to transporting raw materials and final products between different echelons of the network is considered in the proposed model. Also, based on governmental rules, value added tax is defined for raw material and final products transferred through the network. Currency exchange rates related to different activities such as supply and production are presented in the proposed model. In the extended model it is assumed that network is based on the pull mechanism. The number and location of suppliers and customer zones are fixed and predefined. Demand of customer zones should be fully satisfied and potential location of distribution centers and plants are predefined. Flow of raw materials and products is permitted only in the forward direction and Transportation cost and tariff to transport are measured in USD and other costs are measured in local currency. Transportation, production, procurement and fixed opening costs of facilities and also capacity of facilities and customer zones' demand are all accompanied with possibilistic uncertainty. The main issues to be addressed in the GSCN design model include the determination of raw material and final products flow between consecutive echelons, as well as finding the number and optimal locations for plants and distribution centers. Cost minimization is considered as the main objective of the proposed model.

Figure 1. Structure of the studied GSCN

Indices:

| | |
|---|---|
| $i$ | Index of potential locations of plants $i = \{1,2,\dots,I\}$ |
| $j$ | Index of potential locations of distribution centers $i = \{1,2,\dots,J\}$ |
| $k$ | Index of fixed locations of customer zones $i = \{1,2,\dots,K\}$ |
| $l$ | Index of fixed locations of suppliers $i = \{1,2,\dots,L\}$ |

Parameters:

| | |
|---|---|
| $\widetilde{D}_k$ | Demand of customer zone $k$ |
| $\widetilde{ES}_l, \widetilde{EM}_i, \widetilde{ED}_j$ | Exchange rates to USD from local currency at supplier $l$, plant $i$ and distribution center $j$, respectively |
| $\widetilde{CS}_l, \widetilde{CM}_i, \widetilde{CD}_j$ | Costs of raw material at supplier $l$, manufacturing at plant $i$, and processing at distribution center $j$ |
| $\widetilde{FM}_i, \widetilde{FD}_j$ | Fixed costs of opening plant $i$ and opening distribution center $j$ |
| $\widetilde{PS}_l, \widetilde{PM}_i, \widetilde{PD}_j$ | Capacities of supplier $l$, plant $i$, and distribution center $j$ |
| $\widetilde{RS}_{l,i}, \widetilde{RM}_{i,j}, \widetilde{RD}_{j,k}$ | Tariffs of transporting each unit of raw material from supplier $l$ to plant $i$, each unit of final product from plant $i$ to distribution center $j$, and each unit of final product from distribution center $j$ to customer zone $k$ |
| $SE$ | Units of raw material required for producing one unit of final product (constant value) |
| $\widetilde{TS}_{l,i}, \widetilde{TM}_{i,j}, \widetilde{TD}_{j,k}$ | Transportation costs of each unit of raw material from supplier $l$ to plant $i$, each unit of final product from plant $i$ to distribution center $j$, and each unit of final product from distribution center $j$ to customer zone $k$ |
| $\widetilde{VS}_l, \widetilde{VM}_i, \widetilde{VD}_j$ | Value added tax rates at supplier $l$, plant $i$, and distribution center $j$ |
| $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ | Satisfaction degrees $(0 \leq \alpha \leq 1)$ |

Decision variables:

| | |
|---|---|
| $P_{l,i}, O_{i,j}, U_{j,k}$ | Quantity of raw materials shipped from supplier $l$ to plant $i$, quantity of final products shipped from plant $i$ to distribution center $j$, and quantity of final products shipped from distribution center $j$ to customer zone $k$ |
| $X_i$ | $\begin{cases} 1 & \text{if a production plant is opened at location } i \\ 0 & \text{otherwise} \end{cases}$ |
| $Y_j$ | $\begin{cases} 1 & \text{if a distribution center is opened at location } j \\ 0 & \text{otherwise} \end{cases}$ |

Based on the above notation system, the proposed GSCN design model is presented as follows:

$$
\begin{aligned}
\min \sum_i \widetilde{FM}_i \widetilde{EM}_i X_i + \sum_j \widetilde{FD}_j \widetilde{ED}_j Y_j \\
+ \sum_l \sum_i P_{li} \widetilde{CS}_l \widetilde{ES}_l (1 + \widetilde{VS}_l) + \sum_i \sum_j O_{ij} \widetilde{CM}_i \widetilde{EM}_i (1 + \widetilde{VM}_i) \\
+ \sum_j \sum_k U_{jk} \widetilde{CD}_j \widetilde{ED}_j (1 + \widetilde{VD}_j) + \sum_l \sum_i P_{li} (\widetilde{TS}_{li} + \widetilde{RS}_{li})
\end{aligned}
\tag{1}
$$

$$+ \sum_i \sum_j O_{ij}(\widetilde{TM}_{ij} + \widetilde{RM}_{ij}) + \sum_j \sum_k U_{jk}(\widetilde{TD}_{jk} + \widetilde{RD}_{jk})$$

s.t.

$$\sum_j U_{jk} \geq \widetilde{D}_k \qquad \forall k \tag{2}$$

$$\sum_i O_{ij} = \sum_k U_{jk} \qquad \forall j \tag{3}$$

$$\sum_l P_{li} = SE \sum_j O_{ij} \qquad \forall i \tag{4}$$

$$\sum_j O_{ij} \leq X_i \widetilde{PM}_i \qquad \forall i \tag{5}$$

$$\sum_k U_{jk} \leq Y_j \widetilde{PD}_j \qquad \forall j \tag{6}$$

$$\sum_i P_{li} \leq \widetilde{PS}_l \qquad \forall l \tag{7}$$

$$X_i, Y_j \in \{0,1\} \qquad \forall i,j \tag{8}$$

$$P_{li}, O_{ij}, U_{jk} \geq 0 \qquad \forall i,j,k,l \tag{9}$$

Objective function (1) minimizes the total costs of the network design. The first and second terms of the objective function represent fixed costs of locating plants and distribution centers, respectively. The third term corresponds to the total procurement cost of raw materials from suppliers. The total manufacturing cost of products at plants is represented by the fourth term. The fifth term calculates the total distribution cost at distribution centers. The value added tax and currency exchange rates are considered in procurement, manufacturing, and distribution costs owing to governmental regulations and diversity of location of facilities at different regions of the world. The three remaining terms present the total transportation and tariff costs between different echelons of the network. Constraint (2) assures that the demand of customer zones should be satisfied via different distribution centers. Constraints (3) and (4) guarantee the flow balance at plants and distribution centers, respectively. Constraints (5) through (7) correspond to capacity limitations of plants, distribution centers, and suppliers, respectively. Note that the aforementioned constraints prohibit assignment of flow to unopened facilities at different echelons of the network. Constraints (8) and (9) enforce binary and non-negative restrictions on decision variables.

## 4. Possibilistic Programming and Solution Methodology

PP is used in the proposed model to handle imprecise model parameters. It is introduced to find crisp counterpart of the developed possibilistic mixed integer programming. PP is an appropriate approach to deal with epistemic uncertainty and lack of knowledge about model parameters [13]. In the proposed model, an uncertain parameter has a possibility distribution that denotes the possibility degree of occurrence of the parameter's value. Possibility distributions are determined by using available data and experts' knowledge.

The PP approach employs "an expected value" (EV) and "an expected interval" (EI) of a fuzzy number through its defuzzification process on the original model. In the proposed PP model, uncertain parameters are represented by triangular fuzzy functions [14]. The membership function of a triangular fuzzy number $r$ with three prominent points $r = (r^p, r^m, r^o)$ can be represented as follows:

$$\mu_r(x) = \begin{cases} 0, & if\ x \le r^p \\ f_r(x) = \dfrac{x - r^p}{r^m - r^p}, & if\ r^p < x \le r^m \\ 1, & if\ x = r^m \\ h_r(x) = \dfrac{r^p - x}{r^o - r^m}, & if\ r^m \le x < r^o \\ 0, & if\ x \ge r^o \end{cases} \tag{10}$$

The EV and EI of the triangular fuzzy number based on formulated membership function can be represented as follows [15]:

$$EI(r) = [E_1^r, E_2^r] = [\frac{1}{2}(r^p + r^m), \frac{1}{2}(r^m + r^o)] \tag{11}$$

$$EV(r) = \frac{E_1^r + E_2^r}{2} = \frac{r^p + 2r^m + r^o}{4} \tag{12}$$

Moreover, based on the ranking method [16], for each couple of fuzzy numbers $a$ and $b$, their relationship can be characterized using the following relation:

$$\mu_M(a, b) = \begin{cases} 0, & if\ E_2^a - E_1^b \le 0 \\ \dfrac{E_2^a - E_1^b}{E_2^a - E_1^b - (E_1^a - E_2^b)}, & if\ 0 \in [E_1^a - E_2^b, E_2^a - E_1^b] \\ 1, & if\ E_1^a - E_2^b > 0 \end{cases} \tag{13}$$

Term $a \ge_\alpha b$ is the same as $\mu_M(a, b) \ge \alpha$, which means $a$ is greater than or equal to $b$ at least in degree of $\alpha$. Now, consider the following general optimization model formulation:

$$\begin{aligned} \min\ & Z = c^t x \\ \text{s.t.}\ & a_j x \ge b_j, & j = 1, \dots n \\ & l_j x \le d_j, & j = 1, \dots, n \\ & x \ge 0 \end{aligned} \tag{14}$$

All presented parameters in model (14) are uncertain and it is assumed that they could be modeled with triangular fuzzy membership functions. Based on [14], the minimum satisfaction degree of decision vector $x \in R^n$ in the formulated model is equal to α if $min_{j=1,\dots,n}\{\mu_m(a_j, b_j)\} = \alpha$. Finally, the equivalent crisp counterpart of model (14), with regard to (11) to (13), can be presented as follows:

$$\begin{aligned} \min\ & Z = EV(c^t)x \\ \text{s.t.}\ & [(1 - \alpha)E_2^{a_j} + \alpha E_1^{a_j}]x \ge \alpha E_2^{b_j} + (1 - \alpha)E_1^{b_j}, & j = 1, \dots, n \\ & \left[\alpha E_2^{l_j} + (1 - \alpha)E_1^{l_j}\right]x \le (1 - \alpha)E_2^{d_j} + \alpha E_1^{d_j}, & j = 1, \dots, n \\ & x \ge 0 \end{aligned} \tag{15}$$

Using the explained possibilistic approach, it is possible to consider the satisfaction degree for each constraint and find the crisp counterpart and the optimal solution accordingly. In order to present deterministic model corresponded to the possibilistic model, we use the same notations for the parameters and variables. Three prominent points of fuzzy numbers represent the tainted uncertain parameters in the deterministic model. The crisp counterpart is as follows:

$$\min \sum_i \left( X_i \frac{FM_i^p EM_i^p + 2FM_i^m EM_i^m + FM_i^o EM_i^o}{4} \right)$$

$$+\sum_{j}\left(Y_j\frac{FD_j^p ED_j^p + 2fd_j^m ED_j^m + fd_j^o ED_j^o}{4}\right)$$

$$+\sum_{l}\sum_{i}P_{li}\left[\frac{CS_l^p ES_l^p + 2CS_l^m ES_l^m + CS_l^o ES_l^o + CS_l^p ES_l^p VS_l^p + 2CS_l^m ES_l^m VS_l^m + CS_l^o ES_l^o VS_l^o}{4}\right]$$

$$+\sum_{i}\sum_{j}O_{ij}\left[\frac{CM_i^p EM_i^p + 2CM_i^m EM_i^m + CM_i^o EM_i^o + CM_i^p EM_i^p VM_i^p + 2CM_i^m EM_i^m VM_i^m + CM_i^o EM_i^o VM_i^o}{4}\right]$$

$$+\sum_{j}\sum_{k}U_{jk}\left[\frac{CD_j^p ED_j^p + 2CD_j^m ED_j^m + CD_j^o ED_j^o + CD_j^p ED_j^p VD_j^p + 2CD_j^m ED_j^m VD_j^m + CD_j^o ED_j^o VD_j^o}{4}\right] \qquad (16)$$

$$+\sum_{l}\sum_{i}P_{li}\left[\frac{TS_{li}^p + 2TS_{li}^m + TS_{li}^o + RS_{li}^p + 2RS_{li}^m + RS_{li}^o}{4}\right]$$

$$+\sum_{i}\sum_{j}O_{ij}\left[\frac{TM_{ij}^p + 2TM_{ij}^m + TM_{ij}^o + RM_{ij}^p + 2RM_{ij}^m + RM_{ij}^o}{4}\right]$$

$$+\sum_{j}\sum_{k}U_{jk}\left[\frac{TD_{jk}^p + 2TD_{jk}^m + TD_{jk}^o + RD_{jk}^p + 2RD_{jk}^m + RD_{jk}^o}{4}\right]$$

s.t.

$$\sum_{j}U_{jk} \geq \alpha_1\frac{D_k^o + D_k^m}{2} + (1-\alpha_1)\frac{D_k^p + D_k^m}{2} \qquad \forall k \qquad (17)$$

$$\sum_{i}O_{ij} = \sum_{k}U_{jk} \qquad \forall j \qquad (18)$$

$$\sum_{l}P_{li} = SE\sum_{j}O_{ij} \qquad \forall i \qquad (19)$$

$$\sum_{j}O_{ij} \leq X_i\left[\alpha_2\frac{PM_i^p + PM_i^m}{2} + (1-\alpha_2)\frac{PM_i^m + PM_i^o}{2}\right] \qquad \forall i \qquad (20)$$

$$\sum_{k}U_{jk} \leq Y_j\left[\alpha_3\frac{PD_j^p + PD_j^m}{2} + (1-\alpha_3)\frac{PD_j^m + PD_j^o}{2}\right] \qquad \forall j \qquad (21)$$

$$\sum_{i}P_{li} \leq \alpha_4\frac{PS_l^p + PS_l^m}{2} + (1-\alpha_4)\frac{PS_l^m + PS_l^o}{2} \qquad \forall l \qquad (22)$$

$$X_i, Y_j \in \{0.1\} \qquad \forall i,j \qquad (23)$$

$$P_{li}, O_{ij}, U_{jk} \geq 0 \qquad \forall i,j,k,l \qquad (24)$$

## 5. Experimental Results and Analysis

To evaluate the performance of the proposed model and demonstrate its effectiveness, it is solved based on different scenarios. The obtained results and their corresponding analysis are presented in this section. To evaluate performance of the model, 10 suppliers are considered and their locations are fixed and predetermined. They supply raw materials to plants. There are 10 potential locations for opening plants. Plants produce final products and deliver them to customer zones via distribution centers. There are 10 potential locations for opening distribution centers. Locations of plants and distribution centers can be chosen as a strategic decision. There are 15 customer zones and their locations are predefined and fixed. The sizes of the designed test problem

are presented in Table 1. The prominent points of the uncertain parameters are generated randomly using different uniform distributions functions. To evaluate performance of the developed model, it is coded in General Algebraic Modeling System (GAMS) optimization software.

Table 1. The size of the test problem

| No. of suppliers | No. of plants | No. of distribution centers | No. of customer zones |
|---|---|---|---|
| 10 | 10 | 10 | 15 |

The sensitivity analysis is performed to investigate the effects of the transportation costs between plants and distribution centers on the total cost of the network under different satisfaction degrees. As it is shown in Figure 2, the increase of satisfaction degree, while transportation cost is fixed, results in risk-averse performance of the extended model and accordingly obtained the final solution. A higher satisfaction degree in Equations (20), (21), and (22) enforces the model to consider lower capacity for suppliers, plants and distribution centers. A higher satisfaction degree in equation (17) results in a higher value of demand in the supply chain model. Accordingly, the total costs of transportation and processing increases because of using more raw materials and operating more processes through the network to satisfy the higher customers' demand. Moreover, the figure shows as the transportation costs increase, the objective function value also increases.

In another sensitivity analysis, the total cost of the designed network is investigated under different satisfaction degrees. As Figure 3 illustrates, increasing $\alpha_1$ while holding the other three satisfaction degrees fixed leads to a higher total network cost. Figure 3 also shows that as the satisfaction degrees related to capacity constraints of facilities ($\alpha_2, \alpha_3, \alpha_4$) increase, the total cost of the network increases. As the results show, the model enables decision makers to consider different levels of risk associated with the demand and production capacities.



Figure 2. Sensitivity analysis results of transportation costs

Figure 3. Sensitivity analysis results of satisfaction degrees

## 6. Conclusions and future work

This paper proposes a new possibilistic GSCN model that considers globalization parameters. Possibilistic programming is employed to deal with uncertainties in parameters. The model is established to enable decision makers to adjust the satisfaction degrees of the uncertain parameters. The computational experiments and sensitivity studies show the proposed possibilistic GSCN model can effectively control uncertainty of parameters and increase the quality of decisions. The proposed model can be employed as a reliable and effective decision making tool by international supply chain decision makers.

Future work may involve developing a bi-objective model to minimize the total costs and the lead time in the network simultaneously in the context of GSCN and applying possibilistic programming to handle the uncertainty of parameters. Another potential direction of the further research is to introduce the green concept into the global supply chain network to design an environmentally friendly network while controlling the uncertain parameters.

**References**

[1]     Zhang A, Luo H, Huang GQ. A bi-objective model for supply chain design of dispersed manufacturing in China. International Journal of Production Economics 2013;146:48–58.

[2]     Pishvaee MS, Torabi SA. A possibilistic programming approach for closed-loop supply chain network design under uncertainty. Fuzzy Sets and Systems 2010;161:2668–83.

[3]     Pishvaee MS, Razmi J, Torabi SA. An accelerated Benders decomposition algorithm for sustainable supply chain network design under uncertainty: A case study of medical needle and syringe supply chain. Transportation Research Part E: Logistics Transportation and Review 2014;67:14–38.

[4]     Bing X, Bloemhof-Ruwaard J, Chaabane A, van der Vorst J. Global reverse supply chain redesign for household plastic waste under the emission trading scheme. Journal of Cleaner Production 2015;103:28–39.

[5]     Hammami R, Frein Y. Integration of the profit-split transfer pricing method in the design of global supply chains with a focus on offshoring context. Computers and Industrial Engineering 2014;76:243–52.

[6]     Zhang X, Huang S, Wan Z. Optimal pricing and ordering in global supply chain management with constraints under random demand. Applied Mathematical Modelling 2016;40:10105–30.

[7]     El-Sayed M, Afia N, El-Kharbotly A. A stochastic model for forward–reverse logistics network design under risk. Computers and Industrial Engineering 2010;58:423–31.

[8]     Ramezani M, Bashiri M, Tavakkoli-Moghaddam R. A new multi-objective stochastic model for a forward/reverse logistic network design with responsiveness and quality level. Applied Mathematical Modelling 2013;37:328–44.

[9]     Vahdani B, Tavakkoli-Moghaddam R, Jolai F, Baboli A. Reliable design of a closed loop supply chain network under uncertainty: An interval fuzzy possibilistic chance-constrained model. Engineering Optimization 2013;45:745–65.

[10]    Hatefi SM, Jolai F, Torabi SA, Tavakkoli-Moghaddam R. A credibility-constrained programming for reliable forward–reverse logistics network design under uncertainty and facility disruptions. International Journal of Computer Integrated Manufacturing 2015;28:664–78.

[11]    Eskandarpour M, Dejax P, Péton O. A large neighborhood search heuristic for supply chain network design. Computers and Operations Research 2017;80:23–37.

[12]    Varsei M, Polyakovskiy S. Sustainable supply chain network design: A case of the wine industry in Australia. Omega 2017;66:236–47.

[13]    Naderi MJ, Pishvaee MS, Torabi SA. Applications of Fuzzy Mathematical Programming Approaches in Supply Chain Planning Problems. Fuzzy logic in its 50[th] year 2016, p. 369–402.

[14]    Jimenez M, Arenas M, Bilbao A, Rodriguez MV. Linear programming with fuzzy parameters: An interactive method resolution. European Journal of Operational Research 2007;177:1599–609.

[15]    Heilpern S. The expected value of a fuzzy number. Fuzzy Sets and Systems 1992;47:81–6.

[16]    Jimenez Ma. Ranking fuzzy numbers through the comparison of its expected intervals. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 1996;4:379–88.

**Part V**

# Friday April 13th 11:00am-12:30pm Contributed Talks 5: Social Networks (Symposium Hall)

# Stable and unstable frustrations in aged social networks

Leila Hedayatifar[1,2], Foroogh Hassanibesheli[3], and G. Reza Jafari[2,4]

[1] New England Complex Systems Institute, Cambridge, MA, USA
[2] Department of Physics, Shahid Beheshti University, Tehran, Iran
[3] Department of Physics, Humboldt University, Berlin, Germany
[4] Center for Network Science, Central European University, Budapest, Hungary

## 1 Abstract

In social network dynamics, links can carry weight representing strength of relations. The other important parameter that comes into play is age of links. In presence of memory, individuals remember history of their relations that leads to increase the inclination to change the old ones. This ability develops social concepts such as commitment and allegiance leading to the formation of cultural communities and political groups. Based on Heider Balance theory, dynamics of relations in societies goes towards reducing tension and reaching local or global minimum tension states. We add a temporal kernel function into the conventional differential equation which represents the dynamic of links based on Heider balance theory. Adding this term allows the quality of past relations to contribute to the evolution of the system. Here, relations are considered as positive/negative referring to friendship/animosity, profit/nonprofit, etc. This theory studies dynamics of societies based on triadic configurations in which relations evolve to reduce the number of unbalanced triads. In this regard, assigning a potential energy gives us more quantitative view of the network's dynamic. During aging of links, some nodes get older resulting in the formation a skeleton under the skin of society. Even though network's dynamic gets affected by memory, still the general trend goes towards obtaining stable states with negative energies. The resistance of aged networks against the changes decelerates the evolution of the system and traps it into long-lived glassy states. Despite having some unbalanced triads, when glassy states appear in negative energies, the society is still stable. But when resistance is high, glassy states with positive energies have more chance to be seen. In this situation, frustration goes up and society gets unstable.

Figure 1: In this model, links are positive (solid lines) or negative (dashed lines). Thickness of lines represents age of links. Age of relations increases as they do not change over time.

# Complex Networks of Collaborative Reasoning

Sarah Shugars[1]

[1] Network Science Institute
Northeastern University
shugars.s@husky.neu.edu

Collaborative reasoning is a core facet of human society. In business, politics, and everyday life, individuals with varying opinions, experience, and information attempt to collaborate and make decisions. This reasoning process is more complex than collaborative search [1, 2] because agents each hold their own beliefs and may reject each other's opinions regarding the value of different solutions. The deliberative ideal imagines agents discovering optimal solutions through reasoned exchange, but can that ideal be met, even assuming rational agents? In this work, we show that if agents reason together in good faith, they can reach consensus and identify optimal solutions. We model collective reasoning as a complex network, drawing on scholarship indicating that cognitive processes are networked [3, 4]. We define the solution space as an Nk landscape, initiated as a weighted, signed network, and each agent begins with their own Nk landscape representing their beliefs about the solution space. At each time step, $t$, an agent either provides a reason (e.g. an edge) or evaluates a reason presented to them. An agent moves towards a received opinion if there is a positive cosine similarity between the first eigenvector of the agent's existing beliefs and the first eigenvector of the proposed beliefs. We find that even when agents begin with divergent views rational deliberation results in a shared interpretation of the belief space.



*Figure 1:* Agents Converging on Shared Beliefs

# References

[1] D. Lazer and A. Friedman, "The network structure of exploration and exploitation," *Administrative Science Quarterly*, vol. 52, no. 4, pp. 667–694, 2007.

[2] W. Mason and D. J. Watts, "Collaborative learning in networks," *Proceedings of the National Academy of Sciences*, vol. 109, no. 3, pp. 764–769, 2012.

[3] R. Axelrod, *Structure of decision: The cognitive maps of political elites.* Princeton university press, 1976.

[4] R. J. Shavelson, "Methods for examining representations of a subjectmatter structure in a student's memory," *Journal of Research in Science Teaching*, vol. 11, no. 3, pp. 231–249, 1974.

# On the Impact of Network Connectivity in Colonel Blotto Games

James Headrick, Jenny Liang and Emrah Akyol

ECE, Binghamton University, State University of New York
Email: {jheadri1, jliang31, eakyol}@binghamton.edu

## Abstract

We numerically analyze the networked Colonel Blotto game used in cyber-security. The network connects the cyber-nodes that are vulnerable to attacks to the physical nodes that have significance to the players. We consider the case where there are multiple attackers in a war with a common defender, and focus on the impact of network connectivity on the Nash equilibrium of this game.

The Colonel Blotto game models a scenario in which two players having certain resource levels fight over a finite number of battlefields [1, 2]. The players decide on the amount of resource they deploy on each battlefield in order to maximize their payoff which is defined as the number of battles won. Here, we formulate a three-stage Colonel Blotto game where there are two adversaries fighting a common defender (or two systems defending against a common attacker). We incorporate a network between the cyber-nodes (that are not secure, but not valuable) and the physical nodes (that are connected to cyber-nodes through the network). We numerically compute the Nash equilibrium of the game in various parameter regions. Our preliminary numerical results (shown below) indicate that under some conditions, the players may have an incentive to add battlefields or form a coalition as it improves their expected payoffs. In particular, it might be beneficial for security agencies to increase the number of entities that can be under attack and/or form a coalition with other security agencies to share resources, which increases the overall security.



Figure 1: Payoff versus connectivity

## References

[1] E. Borel and J. Ville, Applications de la theorie des probabilites aux jeux de hasard. J. Gabay, 1938

[2] O. Gross and R. Wagner, A continuous Colonel Blotto game, RAND, 1950.

**Fault Tolerance of a Culper Spy Ring Network:**
**How Decentralization, Optimization, and Resilience Fought in the American Revolution**

**Daniel K. Trembley and Samuel O. Heiserman**
**Department of System Science**
**Binghamton University, Binghamton, New York 13902, USA**

**Abstract**

During the American Revolutionary War, George Washington devised a spy ring to gather information on British movements in and around the besieged city of New York. Small groups of civilians were recruited to observe movements of the British forces. Getting caught as a spy meant imprisonment at best, execution at worst [1]. The spy network that operated in this theater was known as the Culper Spy Ring. The Culper Spy Ring was extremely successful at getting accurate and detailed information from the source to the general command of the Continental Army. This paper describes the topology and model of the Culper Spy Ring and provides an explanation of the Culper Spy Ring's fault and attack tolerant network structure. An inherent and natural topology formed by strong internally connected cliques that are weakly linked to other existing cliques forms a robust structure that can withstand disconnection. This is especially the case when the number of weakly tied nodes remains low and randomly dispersed. The topology of a spy ring network thereby survives as a matter of decentralization among active cliques. The goal of the network, therefore, is to maintain the overall decentralization while allowing the maintenance of connections for the percolation of message passing to be normally distributed.

**Keywords**
Culper Spy Ring, Percolation, Fault, Attack, Tolerance, American Revolution, Network Science

**1. Introduction**
The Culper Spy Ring has gained a mythology in recent years. There are a lot of known facts, but there is a weighty amount of conjecture, half-truths, and just plain made up stories regarding a small group of American patriots. Television shows like *Turn: Washington Spies* or Brian Kilmeade and Don Yeager New York Times Bestseller *George Washington's Secret 6: The Spy Ring That Saved the American Revolution* have brought great celebrity amongst the founders and participants of the Culper Spy Ring. So, the question is asked, "What type of network was the Culper Spy Ring?"

**1.1 A Brief History of the Culper Spy Ring**
The American Revolution brought about a unique need for communicating secret messages over distances to inform General George Washington of the British movements in and around New York City. The secret messages contained information about troop counts, soldier movements, logistics, shipping arrivals and departures, names of loyalist colonists, etc. Originally, the use of regular army personnel was used as couriers for the messages. However, this was found to be a bad idea in communicating secret messages because they were easily spotted, caught, tortured, and hanged. In employing a civil force, the Continental Army could under covert measures deliver secret dispatches over great distances without disruption [1].

The Culper Spy Ring was a small group of individuals who knew each other well and had contacts with other groups of different spy rings. The spy rings all formed connections with each other through an authentication process of secret names, passcodes, and remote meeting locations [1]. An individual in one ring only knew the names of one or two members in another ring. Spy ring members were regular citizenry who socialized in the same local churches, pubs, and schools. The members in the spy ring knew each other well within their social circles of village life. Their communication was face to face within the privacy

of their spaces. The British could not listen in or send a spy into their midst without raising great suspicion. This formed a security layer of protection against getting caught. The British would have to intercept messages between spy ring members (i.e., nodes) and remove the member who were caught in the act of delivery. The other spy ring member waiting for the messages would rarely get caught because the passing of messages was randomly timed and located in the vast wooded landscape of New England. This meant that if the deliverer did not show up at the designated location at the designated time, the mission was aborted until a new schedule was devised.

## 1.2 The Modern-Day Context of the Culper Spy Ring

The idea of the Culper Spy Ring was only introduced into the American historic record around 1929 [1]. The idea that a secret society of spies could elude the historical record for more than 150 years is an incredible story in itself, and only when researchers started analyzing enough journals and diaries of the Culper Spy Ring, did the picture of the United States' first spy ring emerge. It appears that as soon as the Americans had won the Revolutionary War, the need for a spy ring was not required and was summarily disbanded. However, as some historians insist, the Culper Spy Ring possibly evolved into what is now the Central Intelligence Agency or CIA [1] [5].

Using historical facts to recreate a 21$^{st}$ century perspective of the Culper Spy Ring, and relying on the methodologies of network science, the first three questions are asked:
1. Is there a natural network topology that forms in spy rings that makes them inherently attack and fault tolerant?
2. Why was the Culper Spy Ring effective?
3. Could the Culper Spy Ring be effective today in terrorism cells [2]?

The first and second question will be answered in this paper. The third question is far more complicated, but it appears that the modern-day terrorist cells are employing similar network topology to communicate and remain undetected [2] [7].

## 2. Robust Nature of Small Network Clusters: Fault and Attack Tolerance

Providing a network system with high levels of fault and attack tolerance, the system will be able to maintain operation in the face of continuing degradation. The design of the network system will have to factor in fault and attack tolerance. In most networks today, the interconnection of nodes about the network should allow for speed and reliability of the system; however, the tolerance designed into the system that withstands faults and attacks in order to gain more reliability often comes at the expense of speed [3].

The Culper Spy Ring Network (CSRN) topology is a social network based on close interpersonal contact and ideology. The Culper Spy Ring Cluster (CSRC) (see Figure 1) is the collection of individuals who share these attributes. These individuals in the CSRC are linked to each other by one degree of separation. The size of the entire network conforms to small-world and scale-free type networks. In order for the network to function, CSRN must maintain at least one path through the network. Normally, this should resolve to the shortest path method even when many paths are available; however, in a spy ring network, the geodesic distances through the network are best accomplished from a randomized series of choices. That is, the choice of one path should not be repeated a second time through the network. This adds consternation and obfuscation to the attacking system, *the informed agent*, from understanding the exact points needed to disable the network [2] [4].

The expectation of attack and fault tolerance is based on the geometry of tightly jointed edged clusters linked to a different cluster by loosely jointed edges (see Figure 2) [4]. In this case, the network percolation threshold should not exist until many SRCs are joined. Since the nodes are evenly spaced within the SRC, and are joined by one edge between different SRCs, the network will continue to function based on the spatial dimensions, i.e., diameter, and as a function of network decentralization. In other words, the network

does not form *Giant Node Clusters* [5]. This demonstrates that the Culper Spy Ring is not just a random network between SRC, but one developed to orchestrate collective disjoint from other clusters with the maintaining connectivity in spite of percolation loss. It is this that makes the Culper Spy Ring robust. It is attack and fault tolerant by its disjoint simplicity. A node or an entire cluster failure will not cause total network disconnection as long as one percolated route remains.

In the Culper Spy Ring Network, speed is also a factor determined by moving from Start Point A to Start Point B (see Figure 2). The information needs to be actionable, meaning, that if the information arrives too late (Stale) or with too many missing pieces (Data Loss), the information is considered non-actionable. Stale information can rarely be used, and data loss depends on the amount of information that can be reliably retrieved. During the American Revolution, messaging was considered slow. It took many days, weeks, or months for a message to arrive at a final destination. Even though the Culper Spy Ring formed in the mid 1770's, it actually allowed messaging throughputs that could be considered outstanding at the time due to well established road and water ways in and around New York City. The CSR delivered an incredible amount actionable intelligence helping the Continental Army gain the upper hand on the by far superior in numbers and equipment British forces [1].

## 2.1 Culper Ring Network Model

The study of the Culper Spy Ring Network has unique implications in network theory. Strongly connected nodes, with other surrounding strongly connected nodes, are joined together by weakly connected edges (see Figure 2). It is maintaining the weakly connected edges between the Start Node and the End Node that is of prime importance. Because the CSRN is a decentralized network, the inherent morphology of the network structure easily exists that is both overt (strong edge connection) and covert (weak edge connection). By the CSRN historical account, it was extremely adaptive moving in and around opposing forces. There is no record of a CSR member getting caught or outed as a spy [1].

For this model, the Spy Ring Cluster (SRC) is a group of five individuals who know each other through a well-established social network structure: churches, pubs, schools, relations, etc. (see Figure 1). The SRC is considered a clique that is impenetrable by outside influence, and all nodes within the clique are persistent, equidistant, and share the same weight. The SRC is a regular graph of common degree four in which the five individuals are completely connected. The symmetry provides a defensible posture from within itself. Messages passed between members of the clique can be directed randomly. A node that is attacked, thus removed, only effects the possibility of how many places within the SRC the message can move. For complete failure of the SRC, all nodes would have to be disabled, or the SRC would have to be isolated from the entire network.



**Figure 1 – Spy Ring Cluster (SRC) Regular Graph of Degree 4**

## 2.2 Culper Spy Ring Network Topology

The entire design of the CSRN topology is considered small-world, scale-free, and decentralized. It contains only strongly connected SRC to the other SRC by weak connections (see Figure 2). In order for the network to be fault and attack tolerant, a limit needs to be placed on the number of SRCs that would be necessary for achieving secrecy of the networks existence. The number of SRCs was reached in the original Culper Spy Ring in order to provide network obfuscation and is a factor of the distance between the

originating message and the location of the Commander of the Continentals [6]. The network obfuscation comes from the individual SRC knowing each other well, and the SRC to SRC relationships not knowing each other at all. Thereby, the design of the network implements another layer of the fault and attack tolerance [7].



**Figure 2 - Culper Spy Ring Network Diagram**

The next consideration for fault and attack tolerance is distance. The distance between nodes within the SRC is weighted as 1. This distance is considered significant for the freshness of the information. For the model, the individual nodes of the SRC are within the same proximity of each other. The consideration of distance is necessary for the SRC to SRC attachment, but the distance and time it takes to percolate is weighted as 1 to demonstrate a neutral speed attribute of the network (see Figure 3). In the CSRN, there is no associate penalty for either optimizing the speed (i.e., a measure of distance rate of a message from Point A to Point B) or not. It is a matter of not getting caught (i.e., node removal), and maintaining at least one attachment between Point A, Sam Culper, and Point B, General Washington (see Figure 3). This is the reliability of the network and its ability to withstand attacks and faults.



**Figure 3 – CSRN Topology Dimensions**

**2.3 Culper Spy Ring Network Modelling Rules and Assumptions**
The CSRN model rules and assumptions are a set of proposed guidance from the historical record. These rules include the following:
1. Each Spy Ring Cluster is comprised of 5 nodes.
2. Each of the 5 nodes is aligned and unidirectional with each node to each node in the SRC (see Figure 1).
3. For the SRC, only one node is allowed to connect with another SRC node.

The assumptions of the model are as follows:

1. If the network contains only a few paths between the beginning and end point, the network will not be fault or attack tolerant. The network will experience quick disconnection and will effectively be rendered non-operational.
2. If the network contains many paths between the beginning and end point, the network will be very fault and attack tolerant; however, it loses its ability to provide the necessary secrecy of the network. The network must remain in a decentralized state without succumbing to a non-operational state.
3. An optimum number of paths that create a low centrality must be selected to provide connectivity between beginning and end points. This allows the network to remain operational in the event of attack and remain secret.

## 3. Simulation Results of the Culper Spy Ring Network
### 3.1 Algorithm
Networks following the CSRN structure of strongly connected cliques linked together with weak ties were constructed using Python's Network X library. First a number of small fully connected networks (cliques) were generated and linked together in a chain from first to last, each clique connected to each other with just 1 link. Then a 'Start' node was generated and connected by a single link to the first clique in the chain and an 'End' node was likewise generated and connected by a single link to the last clique. These 'Start' and 'End' nodes represented those people who the message started and ended with, the former being Sam Culper and the latter being General George Washington.

With the network structure in place, redundancy was then added to the graphs by forming new single links between the 'Start' and 'End' nodes by randomly choosing nodes within each connecting cluster. This was to increase the robustness of the network to random node failure. Without the redundancy, there would be only one node connecting each of the 'Start' and 'End' nodes to the rest of the network. If either of the 'Start' and 'End' nodes were removed, the path between 'Start' and 'End' would be disconnected and the network would cease to function. Once these fortifying redundancies were put in place, nodes were removed from the graphs one by one until the path between 'Start' and 'End' was fully disconnected. Upon disconnection, it was recorded what proportion of the graph's nodes were removed before failure. The results demonstrated that higher redundancy made the graphs more robust to random node failure. This is shown in figures 4 through 7, as graphs with higher numbers of redundancies were able to lose much higher proportions of their nodes before failure.

### 3.2 Simulation with Few Connected Clusters



**Figure 4**          **Figure 5**

Figure 4 is the frequency histogram of proportion of nodes removed until failure for 1000 simulation runs with 2 redundancies. The number of cliques is set at 100 and nodes per clique is set at 5. Figure 5 is the

frequency histogram of proportion of nodes removed until failure for 1000 simulation runs with 10 redundancies. Number of cliques is set at 100 and nodes per clique is set at 5.

### 3.3 Simulation with Many Connected Clusters



**Figure 6**                    **Figure 7**

Figure 6 is the frequency histogram of proportion of nodes removed until failure for 1000 simulation runs with 25 redundancies. Number of cliques is set at 100 and nodes per clique is set at 5. Figure 7 is the frequency histogram of proportion of nodes removed until failure for 1000 simulation runs with 50 redundancies. Number of cliques is set at 100 and nodes per clique is set at 5.

### 3.4 Simulation with Optimization of Connected Clusters

Having established that higher redundancy enhances graph robustness to random failure, a set of graphs (size 1000) was generated for each level of redundancy and the average link densities were recorded for each set of graphs. The purpose was to measure a second form of robustness to these networks, a robustness from being captured by the British. The strength of these CSRN's are their ability to withstand node failure, and their ability to withstand capture through the decentralization of information flow. With more uniform degree distributions, the graphs decentralize information flows is that no one member of the spy ring is highly connected enough for the British to capture who would have enough information to uncover the spy ring.

In forming the topological structure of these CSRN's, a multi-objective optimization problem emerges from where the goal is to maximize robustness to random failure (maximal redundancy) by simultaneously maximizing the robustness of potential captured by the enemy (minimal variance in degree distribution). This trade-off is that as more redundancy is built into the graphs the higher the deviation of their degree distributions become (see Figure 8). In practice, this means that increasing the redundancy, i.e., manufacturing greater number of paths from 'Start' to 'End', also increases the variation of the degree distribution; thereby, creating certain nodes with more centrality and connectivity that could be more vulnerable to captured.

From our preliminary experiment, it is seen that for graphs (each comprised of 100 groups that each contain 5 nodes) could most often withstand at least 20% node removal before failure with 10 redundancies. And at least 60% node removal is needed before failure with 25 redundancies. An optimal number of redundancies is therefore between 10 and 25, assuming it unlikely that the network will lose a large proportion of its nodes before they are replaced. This range of redundancies is a balance between the two objectives, as it provides solid robustness to random failure, while still maintaining relatively low variation in degree distribution.

**Figure 8 – Redundancy Graph**

## 4. Conclusion

The Culper Spy Ring Network existed long before the mathematical models of network science were formed. The existence of a resilient and robust network cleverly designed over 230 years ago explains many of the naturally occurring principles of small networks. When investigating "real world" topologies of spy rings, it is often in hindsight the information becomes available. This is especially the case of the Al-Qaeda in the 9/11 hijacking plot [7], or the Irish Republican Organizations Bishopsgate bombing [8]. The information from the actual members of the Culper Spy Ring are still coming to light as journals and diaries of the time are still being located and a healthy debate as to who they were. Even as the fate of their imprint on the world's greatest democracy last, the techniques and tools can leave even a darker stain on civilization.

The techniques of the Culper Spy Ring are still being used today [3]. By choosing small networks with an optimal number of connected components, a communication system forms that is inherently fault and attack tolerant in the face of node disconnection [7]. Remembering the saying, "A rebel to somebody is another's freedom fighter," the existence and growth of "terrorist cells" causes great worry among civil society. It is the acknowledgment and understanding of these networks that they can either be used for good or evil. It is up to a people on which side to stand and how far to go. It is the research of this type that will continue to identify, quantify, and adapt to the changing landscape of the spy ring world. It is the ultimate testament to the intuitiveness of the human intelligence that find and exploits using these methods that are as amazing and terrifying today as they were in 1776.

## References

[1] C. River, The Culper Ring: The History and Legacy of the Revolutionary War's Most Famous Spy Ring, CreateSpace Independent Publishing Platform, 2015.

[2] A. Gutfraind, "Optimizing Topological Cascade Resilience Based on the Structure of Terrorist Networks," *PLoS ONE,* vol. 5, no. 11, pp. 1-7, 2010.

[3] W. Najjar and J.-L. Gaudiot, "Network Resilience: A Measure of Network Fault Tolerance," *IEEE Transactions On Computers,* vol. 39, no. 2, pp. 174-181, 1990.

[4] H. Sayama, Introduction to the Modeling and Analysis of Complex Systems, Geneseo: Open SUNY Textbooks, 2015.

[5] B. Kilmeade and D. Yeager, George Washington's Secret Six: The Spy Ring That Saved the American Revolution, New York: Penguin Group, 2014.

[6] R. Albert, H. Jeong and A.-L. Barabasi, "Error and Attack Tolerance of Complex Networks," *Nature,* vol. 409, pp. 379-382; 542, 200.

[7] V. E. Krebs, "Mapping Networks of Terrorist Cells," *CONNECTIONS,* vol. 24, no. 3, pp. 43-52, 2002.

[8] A. W. Dnes and G. Brownlow, "The Formation of Terrorist Groups: An Analysis of Irish Republican Organizations," *Journal of Institutional Economics,* vol. 13, no. 3, pp. 699-723, 2017.

# Centrality Analysis of Temporal and Multiplex Networks using Eigenvectors

Dane Taylor[1]

[1] University at Buffalo, State University of New York (SUNY); danet@buffalo.edu

Numerous centrality measures have been developed to quantify the importances of nodes $\{i, \ldots, N\}$ in complex networks, and many of them are obtained as the leading eigenvector of some *centrality matrix* $\mathbf{C}$, which is size $N \times N$. We extend eigenvector-based centrality to temporal and multiplex networks, which consist of $T > 1$ network layers that represent different types of edges. We propose to couple the layers' centrality matrices $\{\mathbf{C}^{(t)}\}$ together in a supracentrality matrix [1] (see Fig. 1a), which introduces a coupling strength $\epsilon > 0$ that is a tuning parameter to control the rate at which centralities change across layers (e.g., across time).

**Joint, Marginal and Conditional Centralities.** Each entry $\mathbb{v}_j(\epsilon)$ in the dominant eigenvector of $\mathbb{C}(\epsilon)$ encodes a "joint centrality" that reflects both the importance of physical node $i = \mathrm{mod}(j, N)$ and layer $t = \lceil j/N \rceil$. It is convenient to represent the length-$NT$ eigenvector $\mathbb{v}(\epsilon)$ by an $N \times T$ matrix so that entry $W_{it} = \mathbb{v}_{i+N(t-1)}(\epsilon)$ encodes the joint centrality of physical node $i$ in layer $t$. In Fig. 1b, we illustrate an example network with $T = 3$ and $N = 4$. We indicate the associated joint centralities for $\epsilon = 0.5$ with a table in Fig. 1c. The shaded regions in Fig. 1c describe two new concepts that we call "marginal node centralities" (MNC) $\{x_i = \sum_t W_{it}\}$ and "marginal layer centralities" (MLC) $\{y_t = \sum_i W_{it}\}$, which provide *uncoupled* centralities. We study how the centrality of each node evolves across the layers by also defining "conditional node centralities" $\{Z_{it} = W_{ij}/y_t\}$, which quantify the importances of physical nodes in layer $t$ relative to other physical nodes in that layer. As examples, we apply our method to study empirical network datasets including the United States Ph.D. exchange in mathematics.



(a) Supracentrality Matrix

$$\mathbb{C}(\epsilon) = \begin{bmatrix} \mathbf{C}^{(1)} & \epsilon^{-1}\mathbf{I} & 0 & \cdots \\ \epsilon^{-1}\mathbf{I} & \mathbf{C}^{(2)} & \epsilon^{-1}\mathbf{I} & \ddots \\ 0 & \epsilon^{-1}\mathbf{I} & \mathbf{C}^{(3)} & \ddots \\ \vdots & \ddots & \ddots & \ddots \end{bmatrix}$$

(b) Temporal Network

(c) Centralities for $\epsilon = 0.5$

Time Layer

| Node Index | 1' | 2' | 3' | MNC |
|---|---|---|---|---|
| 1 | 0.3570 | 0.3392 | 0.1679 | 0.8641 |
| 2 | 0.2516 | 0.3298 | 0.2357 | 0.8171 |
| 3 | 0.2471 | 0.3207 | 0.2312 | 0.7990 |
| 4 | 0.2655 | 0.3579 | 0.2927 | 0.9162 |
| MLC | 1.1213 | 1.3476 | 0.9275 | |

Figure 1: *Example temporal network.* **(a)** Supracentrality matrix $\mathbb{C}(\epsilon)$ for temporal network. **(b)** Intra-layer edges (black lines) encode the network at different instances and inter-layer identity edges (gray lines) couple these layers together with strength $\epsilon^{-1}$. **(c)** For coupling strength $\epsilon = 0.5$, we indicate joint node-layer centralities $\{W_{it}\}$ (white boxes), MNC $\{x_i\}$ (right-most row) and MLC $\{y_t\}$ (bottom row).

## References

[1] D. Taylor, S.A. Meyers, A. Clauset, M.A. Porter and P.J. Mucha, *Multiscale Modeling and Simulation* 15(1), 537–574 (2017).

# NetSciDraw Application

Matthew Dabrowski, Bradley Dreher, Chukwudi Kanu, Jake Lewis, and Eli Shirk

Department of Systems Science and Industrial Engineering, Binghamton University


**Contact Information:** netscidraw5@gmail.com

## Abstract

NetSciDraw application is a part of the NetSciApps project that focuses on students learning through systems thinking. The problem facing many young academic communities is that children are not learning how to think in terms of networks. Nearly every system around us today can be understood as a network structure that consists of nodes and the edges that connect those nodes. Therefore, the goal of this project is to build an application that helps students think in more creative ways rather than in linear paths. NetSciDraw is a simple canvas website application that allows users to draw and express their ideas using these nodes and edges. The primary grades of focus are ranged from K-12 so that the children can begin to think about solving problems with networks early on in their academic careers. However, the team has designed this application such that anybody can use the application (including teachers, parents, college students, etc.) to think and solve problems in more creative ways. The team was able to use the current application Loopy by Nicky Case as a template for the NetSciDraw application. The team has been able to use aspects from Loopy that we thought would best benefit our users while also adding in our own features to enhance the product. At the current state of the application, the users are able to draw nodes (of varying sizes) as well as edges to connect the nodes. The users can also put their thoughts and ideas into these nodes while additionally being allowed to change the color of each node. The NetSciDraw application has received a lot of positive feedback from users, in addition to a few specific requests such as suggestions on how some features can be either added or improved such as a feature that enables labeling/coloring links as well as, the ability to change the size of the nodes created. Additionally, users emphasized needing the ability to save and load diagrams. As a result, in the near future the team wants to focus on adding features that allow the user to adjust the size of already created nodes, change the color of the edges, and also allow the user to save and load their created networks, while also brainstorming other ideas for improvement. The team will continue to ask for feedback from users and truly thinks that the NetSciDraw application can be a revolutionary learning tool that can change the way students and people approach the world today and in the future. It is important to note that all this could not have been done without the guidance from our advisor Professor Hiroki Sayama.

Link to our website: http://coco.binghamton.edu/netscidraw/index.html

**Part VI**

# Friday April 13th 11:00am-12:30pm Contributed Talks 6: Self-Organization & Collective Behavior (Tree House)

# What Do (Biological) Complex Systems Want?

P. Adrian Frazier
Center for the Ecological Study of Perception and Action
University of Connecticut
p.adrian.frazier@gmail.com

Goal directedness is a fundamental concept in psychology, movement science, and the study of organisms as complex systems more broadly. We've gained some insight by studying machines, but this is limited. Machines pursue the goals that they do because they were designed to do so by engineers. They tell us little about the origins of goals, how organisms come to be directed by them, or how they come to have the goals that they have.

Other than not being artifacts, organisms differ from machines thermodynamically. Organisms exist *for* and *because* of the flow of energy and matter exchanged with environments. They are *driven into existence* by this flow. Organisms *self-maintain* the flow of energy and matter, as a flame draws oxygen in to itself as fuel and expels carbon dioxide as waste. But unlike a flame, which burns until nothing remains to burn, an organism adjusts the way it self-maintains, self-maintaining the processes that self-maintain the flow of energy and matter. Mark Bickhard calls this *recursive self-maintenance*, and it is, he argues, what separates life from non-life.

In recognition of this, researchers have searched for dynamic principles in line with the thermodynamics of life. Dynamic stability has attracted the most attention and has had the most success. Much research has been dedicated to studying the relative stability of, for instance, in-phase and anti-phase coordination patterns between limbs, people, and a variety of other couplings. Conceptually, goals seem to be best characterized by dynamic stabilities, especially fixed-point attractors. Behavior perseveres, even when perturbed and faced with challenges, until the goal is achieved. That same behavior returns when goal-conditions are no longer satisfied. But the organism is not always well served by stability. Psychopathological behavior is often marked by functional fixedness—overly warn-in routines and implicit beliefs that never change. The same can be said of self-destructive habits.

Stabilities, in the form of steady states, are often the result of a bottleneck in the dissipation of energy. The cohesive effect of gravity, for instance, slows the dissipation of a star's energy to a crawl, holding it far from thermodynamic equilibrium for billions of years. Our atmosphere aggregates volatile gasses far out of equilibrium with respect to one another, because they do not react with one another very quickly—they exhibit a kinetic barrier. Life, as Nick Lane has argued, specializes to break kinetic barriers. The biosphere has produced numerous innovations, in this regard, by coupling dissipative structures, effectively amplifying its dissipative power and acting as a generator of *variety*. Each innovation breaks and erects kinetic barriers and becomes a basis for further innovations. Variety begets variety. Science illustrates this at the level of human society. Good science opens doors to more science and proceeds by coupling the efforts of diverse casts, applied, empirical, and theoretical.

Returning to goals, we most commonly conceive of goal directed behavior as resulting from perturbations, where the organism seeks to establish or re-establish a steady state. But we might also conceive of it as resulting from frustration, or bottlenecking, where the organism seeks new organized states that generate variety, open the door to a more comprehensive existence. Whatever the origins of goals, they must be tied up in the flow of energy and matter, the far-from-equilibrium thermodynamics, the disequilibrium converting processes where life is made possible.

# Collective Self-Motion of Multiple Benzoquinone Particles at the Air-Water Interface

Tianqi Chen[1], Gayatri Phadke[1], Jennifer E. Satterwhite-Warden[1], Dilip K. Kondepudi[2], James A. Dixon[3,4], and James F. Rusling[1,5,6,7,*]

[1] Department of Chemistry (U-3060), University of Connecticut, 55 North Eagleville Road, Storrs, Connecticut 06269, USA; tianqi.chen@uconn.edu
[2] Department of Chemistry, Wake Forest University, Salem Hall, Box 7486, Winston-Salem, North Carolina 27109, USA
[3] Department of Psychology (U-1020), University of Connecticut, 406 Babbidge Road, Storrs, Connecticut 06269, USA
[4] Center for the Ecological Study of Perception & Action, University of Connecticut, 406 Babbidge Road, Storrs, Connecticut 06269, USA
[5] Institute of Materials Science, University of Connecticut, 97 North Eagleville Road, Storrs, Connecticut 06269, USA
[6] Department of Surgery and Neag Cancer Center, University of Connecticut Health Center, Farmington, Connecticut 06032, USA
[7] School of Chemistry, National University of Ireland, Galway, Ireland
* Correspondence: james.rusling@uconn.edu

## Abstract

Collective self-motion of molecules in the living system lays the foundation of life. However, it has been difficult to study its thermodynamics due to the complexity of living systems. Similar self-motion has been observed in some less complex, non-living systems, with useful implications for understanding the mechanism of motions in the living system. Herein, we studied a system of multiple irregular benzoquinone (BQ) particles at the air-water interface and investigated its responses to temperature change and magnetic field. The particles formed a flock during self-motion. The flock of multiple BQ particles showed thermotaxis and magnetotaxis (when a single magnetic BQ particle was included), demonstrating the perception capability of the particles to the local environment, mimicking the living system.

## Extended Abstract

Self-organized motion is an important property of all known life forms. Many living systems, such as flock of birds, show self-motion, while maintaining certain collective patterns. Such dynamic pattern, as an example of dissipative structure, drives the non-equilibrium system towards equilibrium through constant dissipation of energy and generation of entropy, and lead to motion. Based on this self-organizing property, such systems are called active matter.

Similar self-organized motion has been observed in the non-living system. With significantly less complex models than living systems, these can have useful implications for understanding the self-motion of the living system. Our previous study showed collision and coupling of multiple BQ disks at the air-water interface. Herein, a 13 mm diameter BQ disk was broken into multiple irregular particles and spread over the surface of NaCl solution (25 mL, 100 mM) in a surfactant-free glass petri dish (90 mm diameter, 25 mL volume). A digital camera recorded the motion of particles and observations were made until particles sank or dissolved. The responses of particles to the temperature gradient were studied by placing either a hot probe or a cold probe at the air-water interface without contacting the solution surface. The response of particles to the magnetic field was studied with the presence of a magnetic $BQ\text{-}Fe_3O_4$ particle as the sensor.

Irregular BQ particles formed a flock and moved collectively over time. When a hot probe was positioned at a distance from the flock to create a hot zone, the flock migrated towards the hot zone and followed the hot probe as its position was changed. On introduction of a cold probe, the flock drifted away from the cold zone towards the warmer area of the dish. When a magnet was held above the dish to create a magnetic field, non-magnetic particles moved towards the magnet with the help of one magnetic sensor particle. Therefore the flock of irregular BQ particles exhibited thermotaxis and magnetotaxis in this multi-component self-motion system at the air-water interface.

**Moving from individual to collective complex coordination dynamics:**
**Theoretical and methodological challenges**
Maurici A. López-Felip[1,2] and James A. Dixon[1]
[1] Center for the Ecological Study of Perception and Action, University of Connecticut
[2] Team Sports Department at Fútbol Club Barcelona, Barça Innovation Hub

Human solo-rhythmic coordination has been mostly studied in finger wagging or pendulum swinging tasks. Nevertheless, moving to inter-personal coordination has not been easy and several issues have challenged the elaboration of a solid approach. Hence, the goal of our work is to review theoretical and methodological issues to bridge insights from intra- to inter-personal coordination in tasks as complex as team sports (i.e., soccer) that are not putatively oscillatory.

**Theoretical issues** – rhythmic coordination is a generic emergent property that brings internal order to natural systems. What is the origin of such order? While intra-individual rhythmic coordination could, in principle, be achieved through a *neural mechanism*, the coupling of two humans synchronizing during a soccer game is predominantly *perceptual*. However, both cases have properties of far from equilibrium systems found across all scales of nature, including "simple systems" that can be entirely physical (e.g., Rayleigh-Bénard convection) or chemical (e.g., BZ reaction). Thus, testing these generic dynamical principles should be conducted by using a model of rhythmic behavior of sufficient generality that fits both, the intra and inter-personal interactions.

**Methodological issues** – tests via HKB model resulted in similar patterns for both intra- and inter-personal coordination. However, when research moves beyond simple dyads coordination breaks down. As an alternative, the Kuramoto model (cluster amplitude, $\acute{r}$ where high synchronization = 1) has been applied to several different studies showing consistency with previous research. This is a mean-field based measure where oscillators are equally weighted with sinusoidal coupling. This becomes an issue because 1) it provides an average synchronization value rather than a precise arithmetic difference between phases of different individual oscillators; and 2) lacks an explanation for how synchrony occurs spontaneously, viz., lacking the potential function by which the HKB was linked abstractly to generic physical processes that it defined.

**Complex coordination in team sports** – synchrony in soccer and other team sports has been limited to relative phase measures from statistical points such as centroids of a team or mean phase of all player's displacements using the Kuramoto. However, a main limitation of current models is that collective behavior is context independent. That is, players of a team can be highly synchronized without this corresponding to a meaningful coordination dynamics relevant to the context of the game. Considering theoretical and methodological issues layout above, our work tests a set of new variables for determining the synchronous coordination of teammates during soccer as a function of target goal and action mode (offense vs. defense). The work extends the multivariate cluster phase analysis method while employing related behavioral dynamics variables that capture players' movement changes relative to global and local goals (i.e., goal angle, heading error, heading direction). When clustering these to the Kuramoto model, the output provided a measure of the degree of coupling from those variables that at the individual level have meaning in the state space and that at the collective level have meaning when individuals share a common goal. Overall the model captured self-organizing dynamics of collective team sports, revealing differences in synchrony as a function of game mode (defense, offence), distance from goal, and local vs. global goal contexts, not captured in previous models. Implications of our model as a preliminary and important step to bridge the gap from individual to collective self-organized behavior will be discussed.

# Self-Organization in Fluid-Solid Interacting Systems Far From Thermodynamic Equilibrium

Bong Jae Chung
Department of Bioengineering, George Mason University, Fairfax, VA.
Email: bchung5@gmu.edu


Ashwin Vaidya
Department of Mathematical Sciences, Montclair State University, Montclair, NJ.
Email: vaidyaa@montclair.edu

The interaction of fluids with solids is an age old problem and has given rise to several interesting mathematical and physical problems on pattern formation, stability and bifurcations. Pattern formation is seen to be a consequence of thermodynamic disequilibrium in the system and lends itself to mechanical and thermodynamic arguments, among which is the well known principle of Maximum Entropy Production (MEP). This theory has proven to be effective in certain contexts although its overarching effectiveness as a universal principle remains to be established. The terminal orientation of a rigid body in a fluid is a relatively simple example of a dissipative system out of thermodynamic equilibrium and serves as a perfect testing ground for the validity of the MEP principle.

A body interacting with fluid generates flow around it resulting in dissipative losses. Typically, dynamical equations have been employed in deriving the equilibrium states of such immersed bodies in fluids, but they are far too complex and become analytically intractable when inertial effects come into play. At that stage, our only recourse is to rely on numerical techniques which can be computationally heavy and time consuming.

Our previous calculations [1] reveal that the MEP principle is a reliable tool to help predict the equilibrium orientation of highly symmetric bodies such as cylinders and spheroids, near thermodynamic equilibrium. In recent work [2], we expand our analysis to examine bodies with less symmetry (for instance, a half-cylinder in a flow) and Reynolds numbers (inertial parameter) substantially greater than zero which is far from thermodynamic equilibrium. Experiments and numerical studies indicate that symmetry-breaking and inertia have a nuanced effect on the MEP principle, transforming it to a Min-Max of Entropy Production.

We are currently extending our work to understand pattern formation in time-dependent systems and for flows past deformable bodies. Our collective results on these problems allow us to gain valuable insight into the self-organizing principles in dissipative systems.

## References
1.      B.J. Chung, McDermid, K.  and **A. Vaidya**, On the affordances of the MaxEP principle, *European Physical Journal  B: Condensed Matter and Complex Systems*, 87, 2014.

2.      B. Chung, <u>B. Ortega</u>, **A. Vaidya**, Entropy Production in a Fluid-Solid System Far From Thermodynamic Equilibrium, *European Physical Journal E: Soft Matter and Biological Physics*, 40: 105, 2017.

# Contagious Yawning in Virtual Reality: The Influence of Social Presence Within and Behind the Scenes

Andrew C. Gallup[1], Nicola C. Anderson[2], Daniil Vasilyev[2], and Alan Kingstone[2]

[1]Department of Social and Behavioral Sciences
State University of New York Polytechnic Institute
a.c.gallup@gmail.com
[2]Department of Psychology
University of British Columbia
nccanderson@gmail.com
daniil.vasilyev@gmail.com
alan.kingstone@gmail.com

**Abstract**

Contagious yawning (CY) is a well-documented phenomenon in humans and some other animals. It has been hypothesized that CY functions in coordinating behavior and promoting collective vigilance across group members, yet the factors influencing the propagation of this response remain largely unknown. Stemming from our earlier work showing that social presence diminishes CY in the laboratory, we conducted four experiments investigating the effects of social presence on CY in virtual reality (VR). We first assessed how actual social presence in the testing environment altered CY in VR, and then how both implied and actual social presence within the simulated environment modified this response. We show that, similar to a traditional laboratory setting, having a researcher present during testing significantly inhibited CY in VR, even though participants were completely immersed in a virtual environment and unable to see the researcher. Unlike our previous study, however, manipulating the social presence in VR (i.e., embedding recording devices and avatars within the simulation) did not affect CY frequency. These data provide additional evidence that social presence is a powerful deterrent of CY in humans, which warrants further investigation. These findings also have important applications for the use of VR in psychological science.

# Teasing apart environmental and social influences on the movement of individuals in group living species - proof of concept

Maggie Wisniewska[1], Nicole Dykstra[1], Lisa O'Bryan[1], Simon Garnier[1], Guy Cowlishaw[2], Andrew J. King[3], and Gareth Russell[1]

[1] Federated Department of Biological Sciences
New Jersey Institute of Technology & Rutgers University - Newark
mw298@njit.edu
[2] Institute of Zoology
Zoological Society of London
[3] Department of Biosciences, College of Science
Swansea University

Understanding how animal groups move across landscapes has been a long-standing interest in behavioral ecology, but much remains unknown about how animals integrate information about their physical environment and social dynamics. Our goal was to develop a spatially explicit, statistical framework to test for, and separate, the physical and social influences on an animal's movement. To that end, we modified an established resource selection framework, specifically conditional logistic model, which requires data on an animal's movement, movement of its neighbors, and the environment. We incorporated the social component by converting the positions of neighbors into a distance map, similar to the environmental layer. We fit this model to simulated group movement data from a three-zone model of animal aggregation. The key aspect of this simulation, preference for intermediate distance from neighbors, was reflected in the preference function from our fit [Fig.1]. We then fit our model to movement data from domestic goats, and landscape imagery, and demonstrated species-specific preference for distance to neighbors, the herders, and vegetation [Fig.2]. Our work provides a new tool for examining the patterns of habitat preferences among individuals with respect to social impact on such preferences. Application of this model may shine new light on behavioral processes leading to population distribution across space and time.



Fig. 1. Graph of the best model fit, expressed as log-likelihood, with respect to the distance preference of a focal animal to its nearest neighbors (intermediated distance is preferred), plotted over the range of such distances observed in the data



Fig. 2. Graph of the best model fit, expressed as log-likelihood, with respect to the distance preference of a focal goat to its group members (close distance is preferred), to herders (far distance is preferred), and to vegetation (either close or far distance is preferred), plotted over the range of such distances observed in the data

**Part VII**

# Friday April 13th 2:45-4:00pm Contributed Talks 7: Social Dynamics (Symposium Hall)

# Is Systems Thinking a competitive strategy?

Mark Sellers, Andreas Pape, Hiroki Sayama

State University of New York at Binghamton

*Abstract*

Advocates of systems thinking contend that we would make better decisions in our complex world if we used systems thinking. This proposition is difficult to test because systems thinking is itself a complex collection of skills. This experiment explores several strategies available to any real-world systems thinker to detect and adapt to a non-linear, unpredictable, complex system.

The traditional El Farol bar game developed by Brian Arthur, as implemented in NetLogo, is used to provide a complex test environment. In the original game, 100 agents compete for 70 seats at the bar with each one using its own unique set of randomized strategies. All patrons fail if the bar is overcrowded. Stay-at-home agents fail if it is not crowded. Into this environment, three types of agents are introduced to compete with the original set of agents. The first is a 'rational actor' using a fixed strategy loosely derived from prospect theory in economics with an option to buy insurance. The second is an original El Farol agent modified with an adaptive genetic algorithm to evolve its set of random strategies. The third is a 'systems observer' with no prior knowledge of the system that develops a unique strategy by only observing the behavior of the other agents.

The results show that typically, the highest rewarded agents are original 'dumb' random agents. This is not surprising since statistically, some of these agents will exhibit a large deviation from the sample mean. Unfortunately, in a winner-take-all environment, these 'winners' appear to have the best strategies. The adaptive agents tend to converge on the sample mean – eliminating both extreme winners and losers. They are painfully average and boring – but they never fail catastrophically. However, large populations of adapters converging on a common average solution will destabilize the system – they are only boringly average if they are rare. Rational actors only do well with cheap insurance and do very poorly with expensive or no insurance. This is a hard refutation of rational strategies in a complex world. The 'system observer' was remarkably effective when it had broad access to the entire set of agents but that success reduced as the observed sample decreased. This agent shows the most promise as a 'systems thinker' but further work is required to find the optimal access to the other agents. In the real world, access to all of the other agents may be too costly in time or resources to be a competitive strategy.

Though not conclusive, there are several useful strategic behaviors we can observe in the complex market environment created by the El Farol model. 1) The top and bottom performers generally used unfounded, seat-of-the-pants strategies. 2) The presence of too many adaptive agents in the system destabilized the system. 3) Rational actors are only successful with access to inexpensive insurance. 4) Individual systems observers were successful, but need broad, perhaps costly, access.

1. Arthur, W. B. (1994). Bounded rationality and inductive behavior (the El Farol problem). *American Economic Review*, *84*(2), 406-411.
2. Wilensky, U. & Rand, W. (2015). Introduction to Agent-Based Modeling: Modeling Natural, Social and Engineered Complex Systems with NetLogo. Cambridge, MA. MIT Press.
3. Rand, W., Wilensky, U. (2007). NetLogo El Farol Extension 3 model. http://ccl.northwestern.edu/netlogo/models/ElFarolExtension3. Center for Connected Learning and Computer-Based Modeling, Northwestern Institute on Complex Systems, Northwestern University, Evanston, IL.
4. Wilensky, U. (1999). NetLogo. http://ccl.northwestern.edu/netlogo/. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL.

# Dynamics of the Opioid Crisis in the United States

Brennan Klein[1]   Michael Cavanagh[2]   Ginetta Salvalaggio[3]

Kay Strong[4]   Ange-Lionel Toba[5]

[1]Network Science Institute, Northeastern University, Boston, USA

[2]Coaching Psychology Unit, University of Sydney, Sydney, Australia

[3]Department of Family Medicine, University of Alberta, Edmonton, Canada

[4]Department of Economics, Baldwin Wallace University, Berea, USA

[5]Department of Engineering Management and System Engineering, Old Dominion University, Norfolk, USA

Opioid-related deaths in the U.S. have seen a steady year-over-year increase since 1999 with some classifying the situation as an epidemic. Case and Deaton (2017) provide a well-researched description of the rise in mortality and morbidity in the 21st century. The researchers attribute the phenomenon of cumulative disadvantage—a progressive worsening of employment, marriage, and health opportunities for some social groups. Hollingsworth et al. (2017) use fixed effects modeling to test the relationship between opioid-related deaths and emergency room visits relative to state and county level unemployment rates. They find evidence of a statistically significant positive relationship. The current research question asks: Do employment losses in a neighboring node $j$ lead to a *spread of despair*, hence, increased opioid deaths in the target node $i$? Using opioid-related deaths and employment per 100,000 persons at the county level, this study attempts to capture the interaction effects through a novel measure of despair, $D_i(t)$. Using a (weighted, directed) network where the nodes are census tracts connected if residents of one tract *commute* to another tract for work (as in Nelson & Rae, 2016), we study the propagation of opioid overdose deaths as a function of employment losses. We hypothesize that despair arising from employment losses spreads via network proximity. We found that the distribution of Pearson's correlation coefficients between the $D_i(t)$ of a county and its number of opioid deaths above what is expected based on the national average was positive (p=0.0013). This positive correlation supports the hypothesis that the economic despair in neighboring counties is associated with a given county's opioid-related deaths in the next year. One implication of this study is the importance of using group dynamics, i.e. interaction effects, to model social problems that involving multiple systems—psycho-social-behavioral and economic. By extension, resolution of such challenges necessitates a complex systems approach.

Figure 1: (a) U.S. Commuter Network, colored by community membership; (b) Distribution of Pearson's correlation coefficients (blue) between $D_i(t)$ and unexpected deaths in census tract $i$, compared to what we would expect if there were no correlation (orange).

## References

Case, Anne and Angus Deaton. 2017. Mortality and Morbidity in the 21st Century. Brookings Institute. https://www.theatlantic.com/health/archive/2017/07/how-job-loss-can-lead-to-drug-use/534087/

Hollingsworth, Alex, Christopher J. Ruhm and Kosali Simon. 2017. Macroeconomic Conditions and Opioid Abuse. NBER Working Paper 23192. http://www.nber.org/papers/w23192.pdf

Nelson G Dash, Rae A (2016) An Economic Geography of the United States: From Commutes to Megaregions. PLoS ONE 11(11): e0166083. https://doi.org/10.1371/journal.pone.0166083

1

# The influence of the circadian and ultradian rhythms to human mobility: empirical evidences from location-based check-ins

Hugo S. Barbosa[1], Marcos Oliveira[2], Diogo Pacheco[3], Ronaldo Menezes[3], and Gourab Ghoshal[1]

[1]Department of Physics and Astronomy, University of Rochester, Rochester-NY, USA
[2]GESIS – Leibniz Institute for the Social Sciences, Cologne, Germany
[3]BioComplex Lab, School of Computing, Florida Institute of Technology, Melbourne-FL, USA

## Abstract

In spite of the inherent complexity of the decision-making processes governing our traveling behavior, human trajectories exhibit different levels of regularities at many spatiotemporal scales. External factors such as work schedules and social coordination have an impact on our mobility regularities, which are likely to influence the uncertainty on the whereabouts of people. In this paper, we explore the time-frequency components of the temporal variation in human mobility predictability, in particular the rhythms of predictability in the frequency domain. Our results reveal not only circadian cycles, but also suggest that ultradian rhythms such as the circasemidian cycle (12h) play an important role in our mobility patterns.

## 1   Introduction

A better understanding of the mechanisms governing human traveling behavior is crucial to a variety of domains such as epidemic modeling [1, 2], traffic management [3] and national security [4], to name but a few. When it comes to individual-level mobility, human trajectories have been shown to exhibit regularities at multiple spatiotemporal scales, despite the inherent complexity of human decision-making processes. Indeed, the analysis of large populations via mobile phone data has suggested the possibility of predicting up to 93% of human movement [5]. Such predictability, however, tells us only part of the story, since it neglects spatiotemporal constraints behind mobility regularities.

Several constraints in our daily life, such as work schedules and biological processes, restrict our traveling behavior spatially and temporally. For instance, our internal circadian and ultradian (i.e., less than 24h) rhythms have a direct impact on our activity schedules and therefore on our mobility patterns [6–10]. These constraints are likely to spill over the uncertainty on the whereabouts of people. Indeed, guessing that a person will be at home on a Tuesday at 4am will most likely be a correct guess for most of individuals. Yet, though this uncertainty depends on the time dimension, such dependence has never been statistically characterized.

In this work, we describe the temporal regularities of the theoretical predictabilities of human mobility, and examine their different frequency and time components. Our results suggest that in addition to the daily routines, mobility diversity is also marked by periods of approximately 12h and 6h, which correspond to the second and forth harmonics of our internal circadian

rhythm. These findings suggest that the decision-making processes responsible for our visitation regularities are governed by our internal biological cycles beyond the sleeping and feeding needs, evidenced by predominance of the 12h periods over the 8h and 6h cycles.

## 1.1 Data

We leverage on open data sets of three Location-based Social Network (LBSN) platforms, namely:

- *BrightKite* - Contains approximately 2.5 years of geo-tagged check-ins produced by 51,000 users across 772,000 unique locations worldwide.

- *Gowalla* - Contains more than 6M check-ins by 107,092 users in 1.2M unique locations worldwide over a 20-month period spanning from Feb 2009 to Oct 2010.

- *Weeplace* - Contains more than 7M check-ins produced by more than 15,000 Foursquare users visiting over 1M locations in approximately 50,000 cities worldwide from Nov 2003 to Jun 2011. The data corresponds to the Foursquare users who have provided their data to the Weeplace service (now defunct) in order to create dynamic visualizations of their activities.

## 2 Background

### 2.1 Uncertainty and predictability in human mobility

The places that a person visits over time can be seen as a time series $X = \{x(1), x(2), \ldots, x(T)\}$, where $x(t) \in \{1, 2, \ldots, N\}$. We can examine the potential predictability of this individual by measuring the Shannon entropy $S$ of the random variable $X$, given by:

$$S(p) = -\sum_{i=1}^{N} p(i) \log_2 p(i), \tag{1}$$

where $p$ is a probability distribution of $X$ (e.g., the normalized frequency of location visits). From an information-theoretic perspective, entropy measures the expected number of bits required to optimally encode the outcome of a random variable, reflecting the degree of *uncertainty* about future events. In the context of human mobility, given $N$ distinct locations if $p(i)$ represents the probability of an individual to visit a specific location $i$, the Eq. (1) gives us the uncertainty with respect to the mobility of this person.

The entropy $S$, however, only captures an indirect aspect of the predictability of human dynamics. To have a precise characterization, we measure the probability to correctly predict future locations, given a past series of observations [11]. For this, we can also derive from an entropy estimate the amount of potential predictability of a sequence (e.g., a mobility trajectory). In other words, we can estimate the probability of correctly predicting the future locations of an individual, given a past series of observations. Let us denote this probability as $\Pi$. For instance, an individual with $\Pi = 0.7$ could have their future locations predicted $70\%$ of the time by a *perfect predictive algorithm* whereas $30\%$ of the time their whereabouts are indistinguishable from a random trajectory. Given a particular definition of a mobility entropy $S$ and the number of visited locations $N$,

this quantity can be calculated by inverting the relation

$$S = H(\Pi) + (1 - \Pi) \log_2(N - 1), \tag{2}$$

from which we see that this measure is a function of the number of locations $N$ and the functional form of the entropy $S$. The quantity $H(x)$ is the well-known binary entropy function with general form $H(x) = -x \log_2 x + (1 - x) \log_2 (1 - x)$.

Note that the entropy $S$ neglects the sequence in which events (i.e., visitation) takes place. For this, in addition to the Shannon entropy, often times, other entropy-related measures are also reported in the literature [5, 11, 12]. The most common one is a sequence-correlated entropy rate estimator based on the Lempel-Ziv (LZ) compression algorithm [5, 11–13]. However, the LZ method only offers a reasonable approximation to the true entropy of a sequence for very long strings [13]. Since in our formulation each trajectory can be broken in up to 168 bins (24 hours × 7 days of the week), each individual sequence will have a fraction of the length of the original sequences, which would overestimate the predictability values. Moreover, in our approach the idea of a mobility trajectory as an arbitrarily long sequence of visits is demoted in favor of a *routine-oriented* perspective. Thus, the sequential information—on which the LZ estimator leverages—is not relevant anymore. Nevertheless, all the analyses presented in here were also performed with the LZ-entropy estimator and the results were qualitatively identical as the ones obtained from the Shannon entropy. Therefore, for simplicity we focus our analyses on the Shannon entropy and their estimated predictabilities based on Eqs. (1) and (2).

## 2.2 Wavelet analysis

To describe the temporal regularities in human mobility predictability, we use the Continuous Wavelet Transform (CWT) to examine the time series in the frequency domain. The wavelet transform is frequently used to extract from a time series both their time and frequency components (unlike the Fourier transform that only produces a representation in the frequency domain). The method has a long history of successful applications to a variety of domains such as climate prediction [14], digital image processing [15], and medical imaging [16] to name a few.

The wavelet transform of a discrete sequence $Y = \{y(1), y(2), \ldots, y(N)\}$ having observations with a uniform step $\delta t$ can be defined as the following:

$$W_Y(s, n) = \sqrt{\frac{\delta t}{s}} \sum_{t=1}^{N} y(t) \psi^* \left[ \frac{(t - n)\delta t}{s} \right], \tag{3}$$

where the '$*$' denotes the complex conjugate and $s$ is the wavelet scale. The wavelet transform can be seen as a convolution of $X$ with a translated and scaled version of a wavelet function $\psi(\cdot)$. By varying the scale $s$ and translating over time (i.e., varying $n$), we construct a representation of the amplitude of the different periodic features of $X$ and how they vary with time. In our analyses, we used the Morlet wavelet due to its improved frequency resolution in comparison with other candidate functions [14, 17]. For brevity, we omit the details of the wavelet transform and refer the interested reader to the literature [18, 19].

## 3 Results

Our activity routines are characterized by temporal regularities with time and frequency components. Here we measure the entropy and predictability values for individual discrete-time bins, treating the visits occurring within each bin as an independent sequence. We define the bins to be one-hour-long windows, since the majority of our activities or visits tend to be at least one hour. As we are also interested on the daily predictability changes, we include the weekdays as a second-level in this binning scheme. Each bin $t$ represents the activities occurring within each hour of the week (e.g., Monday/9h-10h, Thursday/15h-16h), for which we compute the entropy and estimate the predictabilities values $\Pi_t^u$ for each person $u$, using Eq. (2) and Eq. (1).

We first analyze the aggregated patterns of the entire population for each data set. For this, we calculated the sample mean of the predictabilities $\Pi_t$ across all users within each time bin $t$. In Fig. 1, we can see the average $\Pi_t$ with a remarkable 24h periodic patterns in all three data sets, having high predictability periods during evenings and nighttime hours (peaking around 4-5am), and secondary peaks between noon and 5pm.

Though all data sets are from location-based social network platforms, the predictability amplitude in the Weeplace data is much larger than what we observed in the other data sets. A possible cause is the fact that the Weeplace data was provided by Foursquare users who were interested in visualizing their check-in history. It is reasonable to believe that these users were, on average, more active (in terms of their check-in history) than a regular LBSN user. In fact, the median number of check-ins of BrightKite and Gowalla users were 11 and 25 respectively, whereas for the Weeplace this figure was 329 check-ins. Also, a median Weeplace user has visited 131 single locations while for BrightKite and Gowalla these numbers were 8 and 19 unique locations respectively.

With an information-theoretic measure instead of a simpler quantity such as the relative location frequencies or the pure location diversity, we have encapsulated in a single quantity both the location diversity and their relative frequencies. To decompose the mobility regularities into their different time-frequency components, now we perform different instances of time-frequency analyses, leveraging on on the continuous wavelet transform of population and individual-level predictability data. For simplicity, we examine here the rescaled and centered version of the time series $\hat{\Pi}_t = (\Pi_t - \mu_{\Pi_t})/\sigma_{\Pi_t}$, where $\mu_{\Pi_t}$ and $\sigma_{\Pi_t}$ denote respectively the mean and the standard deviation of $\Pi_t$.



Figure 1: **Time-dependent predictabilities** The average predictability $\Pi_t$ across all users exhibits daily peaks (4-5am) and secondary peaks (12-5pm) of predictability throughout the time series.

As shown in Fig. 2a, $\hat{\Pi}_t$ exhibits remarkably similar curves in all three datasets, including most of their main and secondary components. Though the data sets are from different sources, the three data sets seem to capture the same regularities. Such finding suggests that the temporal variation of the mobility diversity is activity-independent and therefore is likely to be a characteristic manifestation of the underlying human dynamics.

To analyze these regularities, we evaluated the wavelet transform of $\hat{\Pi}_t$ (Fig. 2b,d,e) and found



Figure 2: **Wavelet analysis** – (**a**) The rescaled predictabilities $\hat{\Pi}_t$ reveal a remarkably strong agreement of all three dataset both in the time and frequency domains. (**b**) The wavelet spectrum reveals that the circadian period (approximately 24h) is the most prominent component of the predictability regularity while the second most-pronounced frequency is the 12h period. Additionally, the third strongest component is a period of approximately 6h. A closer inspection of the power spectra (**c**) shows that these three main components can also be observed in the Fourier spectrum. (**d**) In the BrightKite data, the 12h period is even more pronounced than what we observed in the Weeplace data, whereas the 6h period, although present in the spectrum, is not statistically significant. (**e**) Conversely, the 6h frequency is statistically significant in the Gowalla data.

a strong spectral agreement in all data sets. Our finding revealed major peaks of high predictability during the nighttime and valleys around noon. Minor peaks are also observed in the afternoon during the weekdays. The wavelet transform represents the power level of both the periods of the regularities ($y$ axis) and the temporal localization when they happened throughout the time series ($x$ axis). Moreover, a 95% significance region is represented as the areas delimited by the solid black lines whereas the uncertainty region due to boundary effects is under the hashed area. For details on the computation of the significance and uncertainty regions we refer the interested reader to Ref. [19].

For the Weeplace data, the wavelet analysis (Fig. 2b) reveals that the circadian period (approximately 24h) is the most prominent component of the predictability regularity, depicted as the yellowish contoured area, although for Monday and Sunday the area is under the uncertainty region (hashed area) accounting for boundary effects. More surprising, however, is the fact that the second strongest component is not the 8h period, as we would expect—given the working schedules and the sleeping cycles—but rather a 12h period (i.e., the circasemidian period) corresponding to the second harmonic of the circadian rhythm. This period manifests in the plot as a reddish contoured area. The third strongest component is centered approximately around the 6h regime during the day. To describe these components, we also estimated the global wavelet spectrum, which gives us the true power spectrum of the time series. A closer inspection of the global power spectrum (Fig. 2c) shows that these three main components can also be observed in the Fourier spectrum.

In the BrightKite data (Fig. 2d), the 12h period is even more pronounced than what we observed in the Weeplace data, whereas the 6h period, although present in the spectrum, is not statistically significant. Conversely, in the Gowalla data (Fig. 2e), the 12h frequency is not sufficiently pronounced whereas the fourth harmonic (6h) is indeed statistically significant. These results suggest that human traveling behavior could be influenced by internal biological processes way beyond the sleeping and feeding necessities, evidenced by the presence of the 12h rhythm.

However, an alternative hypothesis to the influence of the 12h rhythm is that these periods are in fact rooted on population-level inhomogeneities on the activity routines. For instance, the 12h component could be explained by the same 24h rhythms if a large proportion of the users had their 24h schedules offset by 12h. To test for this hypothesis we performed the wavelet analysis on individual-level data. Since the 95% significance region is a function of the lag-1 autocorrelation $\alpha$ [19], we computed the $\alpha$ correlation for each individual user. Our analyses indicated that the $\alpha$ correlation distribution exhibits a bell-shaped curve with mean $\mu_\alpha \approx 0.58 \pm 0.1$. We then selected users in each dataset with $\alpha \approx 0.58$ and performed the wavelet analyses on their predictability data. Fig. 3a depicts the normalized predictabilities $\hat{\Pi}_t^u$ of one typical user from each dataset.

At an individual-level, the predictability curves are not as smooth as in the population data and neither are the frequency components. The wavelet analyses (Fig. 3b), however, reveal the 12h and 24h cycles, even though the 12h one is not sufficiently strong to satisfy the 95% significance threshold, likely due to data sparsity. Nevertheless, a closer loot at the global wavelet spectrum (Fig. 3c) makes it evident that the 12h cycle is manifested in the data of all three users, being observed in both the wavelet and Fourier spectra.

**Limitations of the study**   Our analyses were performed on three LBSNs datasets, which are likely to oversample specific subsets of the users' locations (e.g., social places). Moreover, particular segments of the population are likely to be overrepresented in these datasets such as young

Figure 3: **The wavelet power spectra of typical users from each dataset –** (**a**) The rescaled predictabilities $\hat{\Pi}_t$ of individual users are not as smooth as what we observed in the population data. (**b**) Nevertheless, the wavelet spectra suggests the presence of the 12h rhythms. (**c**) This evidence is further corroborated by analyses of the Fourier and wavelet power spectra.

adults and/or technology enthusiasts. Further analyses could also account for other sources of population-level inhomogeneities such differences in culture or working schedules.

## 4 Conclusions

In this paper we explored the time-frequency components of the cyclic variations observed in the uncertainty and predictability-levels of human visitation patterns, which are signature of the synchronized nature of many of our activities (e.g., work schedules) [20], convoluted with socially-instituted divisions of the time (e.g., hours and calendars). However, our results indicate that in addition to the daily routines, mobility diversity is also marked by periods of approximately 12h and 6h. Moreover, the absence of a marked 8h cycle – the duration of typical working and sleeping hours – in detriment of a prominent 12h period indicates that the decision-making processes responsible for our visitation patterns are being influenced by our internal biological cycles beyond our sleeping and feeding necessities as well as other socially-motivated routines. Finally, our findings might represent an important advancement to our understanding of the fundamental processes governing human mobility behavior.

## Acknowledgments

# References

[1] D. Balcan et al. "Multiscale mobility networks and the spatial spreading of infectious diseases." In: *Proceedings of the National Academy of Sciences of the United States of America* 106.51 (Dec. 2009), pp. 21484–9. ISSN: 1091-6490.

[2] H. Barbosa et al. "Human mobility: Models and applications". In: *Physics Reports* (2018).

[3] P. Wang et al. "Understanding road usage patterns in urban areas." In: *Scientific reports* 2 (Jan. 2012), p. 1001. ISSN: 2045-2322.

[4] R. Laxhammar. *Conformal Anomaly Detection: Detecting Abnormal Trajectories in Surveillance Applications*. 2014, p. 171. ISBN: 978-91-981474-2-1.

[5] C. Song et al. "Limits of predictability in human mobility." In: *Science (New York, N.Y.)* 327.5968 (Feb. 2010), pp. 1018–21. ISSN: 1095-9203.

[6] M. Stupfel and A. Pavely. "Ultradian, circahoral and circadian structures in endothermic vertebrates and humans". In: *Comparative Biochemistry and Physiology – Part A: Physiology* 96.1 (1990), pp. 1–11. ISSN: 03009629.

[7] F. A. Scheer et al. "Plasticity of the intrinsic period of the human circadian timing system". In: *PLoS ONE* 2.8 (2007). ISSN: 19326203.

[8] J. L. Toole et al. "Coupling human mobility and social ties". In: *Journal of The Royal Society Interface* 12.105 (2015), pp. 20141128–20141128. ISSN: 1742-5689. eprint: `1502.00690`.

[9] C. M. Schneider et al. "Unravelling daily human mobility motifs." In: *Journal of The Royal Society Interface* 10.84 (2013), p. 20130246. ISSN: 1742-5662.

[10] S. Hasan et al. "Spatiotemporal Patterns of Urban Human Mobility". In: *Journal of Statistical Physics* 151 (2012), pp. 304–318. ISSN: 0022-4715.

[11] C. Krumme et al. "The predictability of consumer visitation patterns." In: *Scientific reports* 3 (Jan. 2013), p. 1645. ISSN: 2045-2322.

[12] E. L. Ikanovic and A. Mollgaard. "An alternative approach to the limits of predictability in human mobility". In: *EPJ Data Science* 6.1 (2017). ISSN: 21931127.

[13] A. Lempel and J. Ziv. "On the Complexity of Finite Sequences". In: *IEEE Transactions on Information Theory* 22.1 (Jan. 1976), pp. 75–81. ISSN: 0018-9448.

[14] D. Nalley et al. "Inter-annual to inter-decadal streamflow variability in Quebec and Ontario in relation to dominant large-scale climate indices". In: *Journal of Hydrology* 536 (2016), pp. 426–446. ISSN: 00221694.

[15] J.-P. Antoine et al. "Image analysis with two-dimensional continuous wavelet transform". In: *Signal processing* 31.3 (1993), pp. 241–272.

[16] P. S. Addison. *The illustrated wavelet transform handbook: introductory theory and applications in science, engineering, medicine and finance*. CRC press, 2017.

[17] X. Mi et al. "The use of the Mexican Hat and the Morlet wavelets for detection of ecological patterns". In: *Plant Ecology* 179.1 (July 2005), pp. 1–19. ISSN: 1385-0237.

[18] R. D. Wallen. *The illustrated wavelet transform handbook*. Vol. 38. 4. 2004, p. 298. ISBN: 9780750306928.

[19] C. Torrence and G. P. Compo. "A Practical Guide to Wavelet Analysis". In: *Bulletin of the American Meteorological Society* 79.1 (Jan. 1998), pp. 61–78. ISSN: 0003-0007.

[20] A. J. Morales et al. "Global patterns of synchronization in human communications". In: *Journal of The Royal Society Interface* 14.128 (Mar. 2017), p. 20161048. ISSN: 1742-5689.

# Multi-scale Patterns of Self-Identification

Chandler Squires[1], Nikhil Kunapuli[1], Alfredo Morales[1], Yaneer Bar-Yam[1]

[1] New England Complex Systems Institute

Our work focuses on understanding, measuring, visualizing, and modeling the multi-scale spatial properties of culture. We pay particular attention to *self-identification*: people's overt assertion of the local, regional, and national communities for which they feel a sense of belonging.

Quantifying cultural traits remains a difficult and ill-defined task, so we restrict our attention to a proxy for culture that is more well-defined: usage of Twitter hashtags. One virtue of this choice, beyond narrowing the scope of the problem, is that the data is plentiful. However, this choice does come at a price: hashtags only represent a facet of culture, not culture as a whole, due to peculiarities such as self-selection bias among users, brevity of hashtags, and emphasis on topics such as food, location, and politics. Despite these shortcomings, hashtags are still an excellent proxy for the task of inferring self-identification, since some of the most common hashtags are city, region, and country names.

We chose to analyze hashtags coming from France, Spain, and Portugal during the month of October 2014, dividing the countries into 180 $1° \times 1°$ degree squares that we call "sites". For each site, we composed a feature vector describing its culture, with each of the 65,335 features corresponding to a different hashtag: the value of a feature was taken as the total number of distinct users per day of the corresponding hashtag. In order to accentuate features related to self-identification, we multiplied each feature by the logarithm of its *inverse site frequency*, i.e. the inverse of the proportion of sites for which it had non-zero count. To reduce the dimensionality of this description, we used principal component analysis (PCA) to derive a final description of vastly reduced size, with only 8 features. For each site, its values for the new features, called *principal components*, are weighted sums of its values for the original features. The results are shown in Fig. 1.

To understand the origins of these patterns, we compared our results to those emerging from a well-known model of cultural dynamics, the Axelrod model. The Axelrod model provides a reasonable null hypothesis for cultural evolution, assuming that all cultural evolution is due to the imitation of culturally similar neighbors. The principal cultural components from the real data and the Axelrod model both show a downward trend in the average size of cultural clusters, but the real data shows a sharper decrease. This suggests a real-world cultural mechanism not found in the Axelrod model, favoring regional self-identification.

The importance of our work is both methodological and scientific. We find an appropriate use for the abundance of Twitter data, and use a novel combination of data transformations to measure and visualize the data. Furthermore, we compare the results to a standard model to show similarities and departures, which are suggestive of socially consequential trends in cultural dynamics. Our approach provides an inlet into a plethora of future work, including further analysis of spatial statistics of culture, including time series analysis, new models of cultural dynamics, and comparison of the Twitter data set to other proxies for culture.



Figure 1: **Principal Cultural Components** Panel A shows the combination of components 1 (red channel), 2 (green channel), and 3 (blue channel). Panels B-I show the strengths of components 1-8, respectively, centered and scaled to have a maximum absolute value of 1.

**Part VIII**

# Friday April 13th 2:45-4:00pm Contributed Talks 8: Information (Tree House)

# Modes of Information Flow
## *extended abstract*

Ryan James[3][4], Bahti Zakirov[1][2], Blanca Daniela Masante[3][4], and
James P. Crutchfield[3][4]

[1]*City University of New York, College of Staten Island, Department of Engineering Science and Physics*
[2]*City University of New York, Macaulay Honors College*
[3]*University of California, Davis, Department of Physics*
[4]*University of California, Davis, Complexity Sciences Center*

### Abstract

Information flow is a highly useful concept for understanding the behavior of systems. There have been numerous attempts to quantify information flow, but there exists confusion about the meaning of these measures. We consider two common, though flawed measures of information flow, time delayed mutual information and transfer entropy, and demonstrate that it is erroneous to conflate the results given by these tools with what one is to intuitively believe constitutes information flow. We separate information flow into three modalities of shared, intrinsic, and conditional. In this context, we demonstrate that time delayed mutual information and transfer entropy actually turn out to provide combinations of *shared*, *conditional*, and *intrinsic* information flow, and that a third measure is needed to fully be able to disaggregate the types of information flow that exist within a system. We then propose a new measure, *intrinsic transfer entropy*, which utilizes intrinsic conditional mutual information from information theoretic cryptography. This provides the first concrete method of separating information flow into its *intrinsic*, *conditional*, and *shared* components. We apply intrinsic transfer entropy to a variety of systems to demonstrate its usefulness.

## 1 Overview

We propose that information flow between two time series, $X$ and $Y$, can take three forms. The first, intrinsic dependence, is when the behavior of the $X$ time series directly drives the $Y$ time series. The second, shared dependence, is when the behavior of the $Y$ time series can be inferred from the prior behavior of either time series, perhaps due to a common driver. The third, conditional dependence, occurs when the $X$ time series is pairwise independent of the $Y$ time series, but when combined with $Y$ prior behavior becomes predictive. These three types of flow are illustrated in figure 1. The question now is: How do we detect and quantify these three forms of dependence from time series observations?

1

Figure 1: The canonical forms of the three types of dependence that can exist between two variables ($X_{-1}$, $Y_0$) in the context of a third ($Y_{-1}$). The first, intrinsic dependence, exists regardless of the third. The second, shared dependence, exists synchronously with the third. Finally, the third, conditional dependence, exists only when observing the third. From an information theoretic prospective, both intrinsic and shared dependence contribute to the time-delayed mutual information, both intrinsic and conditional dependence contribute to the transfer entropy, but only the intrinsic dependence contributes to the intrinsic transfer entropy.

In order to quantify dependence, information theory has provided three forms of shared information: the mutual information, the conditional mutual information, and the intrinsic mutual information. The *mutual information* [1] quantifies the reduction in uncertainty about variable $Y$ when given variable $X$:

$$\mathrm{I}\left[X : Y\right] = \mathrm{H}\left[Y\right] - \mathrm{H}\left[Y|X\right] = \mathrm{H}\left[X\right] - \mathrm{H}\left[X|Y\right] \tag{1}$$

$$= \sum_{x,y \in \mathcal{X},\mathcal{Y}} p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)} \ . \tag{2}$$

When in the presence of a third variable, $Z$, we can quantify the additional reduction in uncertainty about about variable $Y$ when given $X$, after already having been given $Z$ using the *conditional mutual information* [1]:

$$\mathrm{I}\left[X : Y|Z\right] = \mathrm{H}\left[Y|Z\right] - \mathrm{H}\left[Y|X,Z\right] \tag{3}$$

$$= \sum_{x,y,z \in \mathcal{X},\mathcal{Y},\mathcal{Z}} p(x,y|z) \log_2 \frac{p(x,y|z)}{p(x|z)p(y|z)} \ . \tag{4}$$

Note that is is possible that $\mathrm{I}\left[X : Y|Z\right] > \mathrm{I}\left[X : Y\right]$.

Finally, the *intrinsic mutual information* [11] [12] has been introduced as a bound on the rate at which $X$ and $Y$ can agree upon a secret key:

$$\mathrm{I}\left[X : Y \downarrow Z\right] = \min_{p(\overline{z}|z)} \mathrm{I}\left[X : Y|\overline{Z}\right] \ . \tag{5}$$

The intuition here is that $\mathrm{I}\left[X : Y \downarrow Z\right]$ quantifies the amount of information that $X$ and $Y$ share in spite of any local operation performed on $Z$. Importantly, $\mathrm{I}\left[X : Y \downarrow Z\right]$ is bound from above by both $\mathrm{I}\left[X : Y\right]$ ($p(\overline{z}|z)$ is a constant) and $\mathrm{I}\left[X : Y|Z\right]$ ($p(\overline{z}|z)$ is the identity).

2

We demonstrate how these measures of information can be utilized to quantify information flow.

Historically, information flow had largely been measured utilizing the *time-delayed mutual information*

$$\text{TDMI}_{X \to Y} = \text{I}\left[X_{:0} : Y_0\right] \ . \tag{6}$$

The time-delayed mutual information is sensitive to both intrinsic and shared dependence, as can be seen in Figure 1. While the time-delayed mutual information captures a restricted notion of causality, "...it is designed to ignore static correlations due to the common history or common input signals"[2]. That is, it conflates intrinsic and shared dependence.

In order to cleave the shared dependence from the time-delayed mutual information, Schreiber proposed the *transfer entropy* [2]:

$$\text{T}_{X \to Y} = \text{I}\left[X_{:0} : Y_0 \mid Y_{:0}\right] \ . \tag{7}$$

While Schreiber correctly intuited[1] that "these influences are excluded by appropriate conditioning of transition probabilities." Unfortunately, as evidenced by statements such as "...(time-delayed mutual information) fails to distinguish information that is actually exchanged from shared information due to common history and input signals", he was unaware of the possibly of conditional dependence and the transfer entropy suffers as a result, failing to distinguish intrinsic from conditional dependence.

In order to overcome this weakness of the transfer entropy, we propose the *intrinsic transfer entropy*:

$$\text{IT}_{X \to Y} = \text{I}\left[X_{:0} : Y_0 \downarrow Y_{:0}\right] \ . \tag{8}$$

This measure correctly excises the shared dependence from the time-delayed mutual information, without inducing conditional dependence, resulting in exactly the intrinsic dependence.

Taken together, these measures allow us to quantify intrinsic, shared, and conditional dependence:

- intrinsic dependence $= \text{IT}_{X \to Y}$

- conditional dependence $= \text{T}_{X \to Y} - \text{IT}_{X \to Y}$

- shared dependence $= \text{TDMI}_{X \to Y} - \text{IT}_{X \to Y}$

# References

[1] T. M. Cover and J. A. Thomas, Elements of Information Theory. New York: Wiley, 1991 2

[2] Schreiber, T. (2000). Measuring information transfer. Physical Review Letters, 85(2):461âĂŞ464 3

---

[1]There is no real derivation or proof that transfer entropy does what it was designed to do, namely measure information transfer. It's form was, largely, intuited by Schreiber and justified by behaving plausibly on several examples.

3

[3] J. T. Lizier, M. Prokopenko, and A. Y. Zomaya. The information dynamics of phase transitions in random Boolean networks. In S. Bullock, J. Noble, R. Watson, and M. A. Bedau, editors, Proceedings of the Eleventh International Conference on the Simulation and Synthesis of Living Systems (ALife XI), Winchester, UK, pages 374381. MIT Press, Cambridge, MA, 2008.

[4] O. Kwon and G. Oh. Asymmetric information flow between market index and individual stocks in several stock markets. EPL (Europhysics Letters), 97(2):28007, 2012.

[5] X. R. Wang, J. M. Miller, J. T. Lizier, M. Prokopenko, and L. F. Rossi. Measuring informa- tion storage and transfer in swarms. In T. Lenaerts, M. Giacobini, H. Bersini, P. Bourgine, M. Dorigo, and R. Doursat, editors, Advances in Artificial Life, ECAL 2011: Proceedings of the Eleventh European Conference on the Synthesis and Simulation of Living Systems, pages 838âĂŞ845. MIT Press, 2011.

[6] James, R.G.; Barnett, N.; Crutchfield, J.P. Information flows? A critique of transfer entropies. Phys. Rev. Lett. 2016, 116, 238701.

[7] Villaverde AF, Ross J, MorÃąn F, Banga JR (2014) MIDER: Network Inference with Mutual Information Distance and Entropy Reduction. PLOS ONE 9(5): e96732. https://doi.org/10.1371/journal.pone.0096732

[8] arXiv:0910.4514v2

[9] H. Matsuda, K. Kudo, R. Nakamura, O. Yamakawa, and T. Murata. Mutual information of Ising systems. Int. J. Theor. Phys., 35(4):839âĂŞ845, 1996.

[10] P. L. Williams and R. D. Beer arXiv:1004.2515v1

[11] U. Maurer and S. Wolf, The intrinsic conditional mutual information and perfect secrecy, in Information Theory. 1997. Proceedings., 1997 IEEE International Symposium on. IEEE, 1997, p. 88.

[12] U. M. Maurer, S. Wolf, Unconditionally secure key agreement and intrinsic information, IEEE Trans. Inform. Theory, vol. 45, pp. 499514, Mar. 1999.

[13] Dit python library docs.dit.io

4

# Information transfer in the contractile membrane of slime mold *Physarum polycephalum* during decision-making

Subash K. Ray[1], Gabriele Valentini[2], Abid Haque[1] and Simon Garnier[1]

[1]Federated Department of Biological Sciences,
New Jersey Institute of Technology & Rutgers University – Newark
sr523@njit.edu
[2]School of Earth and Space Exploration,
Arizona State University

## Abstract

Non-neuronal organisms like bacteria, plants, protists, and yeasts, despite lacking a neuron-based system, have demonstrated information processing capabilities that help them thrive in complex and changing environmental conditions. One such remarkable non-neuronal organism is the slime mold *Physarum polycephalum*. It is a large, unicellular, multi-nucleated organism that self-organizes into a complex system of intersecting tubules. It can solve labyrinth mazes, build efficient tubule networks, and make adaptive decisions when faced with complicated trade-offs, such as between food quality and risk, speed and accuracy, and exploration and exploitation. These extraordinary capabilities of slime molds are poorly understood and believed to be encoded in the contraction-relaxation pattern of its membrane. The membrane is composed of multiple rhythmically contracting regions that change their contractile patterns in response to – a) the quality of the local environment, and b) the contractile pattern in the neighboring regions (i.e. physical coupling between the neighboring regions). Such a coupled-oscillator based system could potentially provide a mechanism for information processing and propagation, and lead to the emergence of problem-solving capabilities at the organism level. Here we will present how information exchange occurs between different contractile regions of a slime mold tubule, while the tubule is making a decision between two food sources. We measured membrane contractions at locations along the length of a straight tubule, while two food sources were placed at its terminals (Figure 1). The quality of the two food sources were either identical (10% vs 10% oat-agar food blocks), or with one significantly better than the other (10% vs 2% oat-agar food blocks). The final decision of the tubule was recorded as the food source with higher biomass aggregation. We used transfer entropy to quantify the net information transfer between different contractile regions along the tubule axis. We show that contractile regions near the lower quality food block (in the 10% vs 2% food block condition) drive the contraction pattern of the whole tubule. Whereas, in the identical food choice condition, the contraction patterns were driven by the tubule regions near the two food blocks. With this project, we aim to provide new information on the foundation of decision-making in non-neuronal organisms.
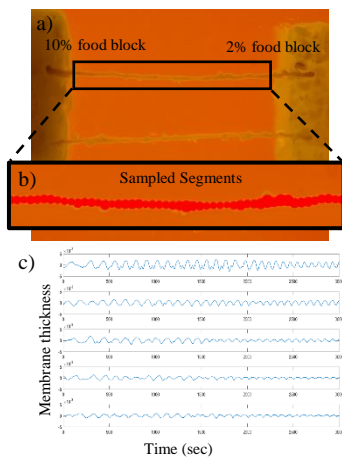
Figure 1. a) Straightened tubules (2cm long) placed between a 10% and 2% oat-agar food blocks. b) Contractions measured at 50 equidistant locations (given by red circles) along the tubule axis. c) Contractions at five evenly spaced locations. Upper curves correspond to locations closer to 10% oat-agar food block, and lower curves closer to 2% oat-agar food block.

# Information flow reveals prediction limits in online social activity

James P. Bagrow[1] and Lewis Mitchell[2]

[1] The University of Vermont      [2] University of Adelaide

**Abstract**   Modern society depends on the flow of information over online social networks, and popular social platforms now generate significant behavioral data. These data offer a wealth of insight to researchers studying human activity. However, it remains unclear what fundamental limits exist when using these data to predict the activities and interests of individuals and whether and to what extent these predictions can be made from an individual's social contacts. By applying tools from information theory to estimate the predictive information within the writings of Twitter uses, we show that approximately 95% of the potential predictive accuracy attainable for an individual is available within the social ties of that individual only, without requiring the individual's data. This information bound holds for any choice of predictive algorithm. Further, distinct temporal and social are visible by measuring the flow of information along social ties, allowing us to better study the dynamics of online activity. Our results have distinct privacy implications: information is so strongly embedded in a social network that in principle one can profile an individual from their available contacts even when the individual chooses to forgo the social platform completely.

**Figure 1:** (**A**) We measure the text entropy rate and ideal predictability to understand how much information about a user's future activity is present in her historical activity. (**B**) To measure information flux and influence we extend this measure to the cross-entropy, capturing how much information about an ego's (green) future is present in an alter's (blue) past. (**C**) These correlated entropies can be related to the predictability $\Pi$, an upper bound on the maximum predictive accuracy of a perfect prediction algorithm ($\Pi = 1$ indicates perfect accuracy and no mistakes, $\Pi = 0$ indicates a complete absence of predictive potential). The distribution of $\Pi$ for the ego (red) is sharply peaked around 0.55; a perfect algorithm has the potential to predict a user's activity with over 50% accuracy. Considerable information on the activity of many (but not all) egos is available in the alter's past (blue), while a random alter generally provides little information (grey; $\Pi < 0.3$ for most pairs).

# References

1. Bagrow, J. P., Liu, X., & Mitchell, L. (2017). Information flow reveals prediction limits in online social activity. arXiv:1708.04575.

2. Bagrow, J. P. & Mitchell, L. (2017) "The quoter model: a paradigmatic model of the social flow of written information." arXiv:1711.00326.

# Quantifying the Causal Structure of Networks

Brennan Klein[1]        Erik Hoel[2]

[1]Network Science Institute, Northeastern University, Boston, MA, USA

[2]Department of Biological Sciences, Columbia University, New York, NY, USA

The causal structure of a network is a function of its directional connections. This structure can actually be quantified and is a function of the *convergence* and the *divergence* in a system (the incoming and outgoing connections, respectively). In 2013, an information-theoretic metric known as *effective information* (*EI*) was shown to encapsulate both of these quantities in systems of logic gates by equating convergence with uncertainty about preceding states and divergence with uncertainty about subsequent states (Hoel, Albantakis, & Tononi, 2013). This metric, originally formulated as a step in the calculation of integrated information (Tononi & Sporns, 2003; Tononi, 2001), draws on information theory to quantify causal influence. In the current work, we show that the *EI* of a network is simply the *noise* present in random walks on the network. This finding allows for the exhaustive classification and quantification of networks' causal structures.

Mathematically, the effective information has been expressed in a number of ways. Here, we introduce a novel definition of the *EI* of a given network as the average Kullback-Leibler (KL) Divergence between the out-weights from $node_i$, $w_{ij_{out}}$, and the stationary in-degree distribution of the network, $P_{in}(k_i)$:

$$\text{EI}_\text{net} = \frac{1}{N} \cdot \sum_i D_{KL}[w_{ij_{out}}||P_{in}(k_i)] \tag{1}$$

Networks with high *EI* contain *more certainty* about the cause and effect relationships between entities in the network, whereas networks with low *EI* contain more noise, thus a random walker at any given node will display *more uncertainty* about future or past nodes. Knowing the *EI* of given a network (and of subgraphs within a network) allows us to extend even further recent advances in the study of *causal emergence* (Hoel, 2017), which is when a "macro-scale" description of a system is a more suitable description of the causal relationships among the components of that system (similar to community detection or coarse-graining in networks). In order to explore causal emergence in networks, we first examined the *EI* of networks with different sizes and structures, one example of which is shown in Figure 1.



Figure 1: By repeatedly simulating Erdős-Rényi (ER) networks of various sizes (from N=10 to N=500) and adjusting the probability, $p$, that any two nodes are connected, we generated hundreds of networks of various sizes. The regime where *EI* is highest is when the average degree is greater than 1.0 and less than $log_2(N)$, which corresponds to the the emergence of a giant component and the point at which every node is likely connected to the giant component.

## References

Hoel, E. P. (2017). When the Map Is Better Than the Territory. *Entropy*, *19*(5), 188.

Hoel, E. P., Albantakis, L., & Tononi, G. (2013). Quantifying causal emergence shows that macro can beat micro. *Proceedings of the National Academy of Sciences*, *110*(49), 19790–5.

Tononi, G. (2001). Information measures for conscious experience. *Archives Italiennes de Biologie*, *139*(4), 367–371.

Tononi, G., & Sporns, O. (2003). Measuring information integration. *BMC Neurosci*, *4*(1), 31.

1

**Part IX**

# Posters

# COLLABORATIVE LEARNING IN COMPLEX ADAPTIVE TEAM STRUCTURE

Muhammad Haris Aziz [a], Summyia Qamar [b]

[a] System Science and Industrial Engineering Department, Binghamton University, New York.

[b] Industrial Engineering Department, University of Engineering and Technology Taxila, Pakistan

**Abstract**

The paper aims at providing pedagogical approach to foster individual and collective learning of students in Complex Adaptive System (CAS). To do so, term-project was designed such that whole class worked As-One team; a complex unit comprising of interacting sub-units. The project was refined through iterative and adaptive process, for which a comparison has been presented between centralized and decentralized decision system. The emergence of system was based on learning from group dynamics. The assessment of project outcome was based on the design and development of innovative product by collaboration and team work of students. Analysis of activities during the whole project indicated that collaboration was initially challenging, however, collective learning which was emerged with the advancement in the project, positively affected individual and team performance. After completion of projects, instructor's evaluation and student survey showed appreciative response regarding project design approach and learning methodology. The As-One strategy fostered collaboration, ownership and intuitive learning among students, to deal with complex situations requiring adaptation and self-organisation.

Keywords: complexity; pedagogical approach; complex adaptive systems; collaboration; emergence; intuitive learning

## 1. Introduction

Complex systems are increasingly influencing our lives and creating challenges to deal with salient traits of complexity-based systems such as change management, self-adaption, unpredictability and emergence. The significance of these traits is broadly acknowledged in different contexts, including sustainability, business management and education [1]. However, the crux of various theories to cope with complexity is the need of "complexivist mindsets", adopted in these contexts. Therefore, to deal with increasingly complex and interwoven world, contemporary education should foster the collaborative and intuitive learning ability among individuals. Theoretical frameworks have been provided to nurture these skills through education systems [1, 2]. But, there remains gap in practical implementation and analysis of concrete experiments. To bridge this gap, this paper aims to provide a complexity-based learning strategy where whole class worked As-One unit and its significance in a higher education engineering course.

## 2. Literature Review

### 2.1. Education as complex system

Research on complex systems highlights learning and education as complex phenomena and regards educational system as Complex Adaptive System (CAS) [3, 4]. Morrison [5] explained CAS as dynamical, emergent and non-linear organization operating in unpredictable environment, where they adapt to micro and macro societal changes through self-organization. CAS are considered emergent due to the property of self-organized interaction among the system agents

while adapting to external changes. Thus, emergence initiates new, unpredictable patterns of system which cannot be represented by individual parts of the system [6]. This property is analogous to the education systems in which self-organized and adaptive evolution of students engender transformations which leads to different forms of the same system. Hence, educational systems, as part of CAS, are considered as proper learning systems [1, 7].

Undergraduate term-project is a specific instance of higher education system. Therefore, it can be regarded as a complex system [2]. However, due to constraints including, fragmented curriculum and limited time, the instructors adopt conventional techniques of student evaluation which hinders collective learning, self-organization and emergence [8].

## 2.2. Why a complexity-based educational strategy

The concept of learning in CAS requires support and guidance of self-organization ability by understanding tendencies and natural behaviour of individuals. Research on complexity management infers that CAS cannot be governed by recurring centralized decision making and control [9, 10]. Therefore, complexity-based educational strategies are required in which students rely on continuous iterative cycles, assessment of results and revision of plans based upon their emerging evidences and collective learning [11].

Importance of collective learning is attributed from the evolution of individual learning and vice versa. This is because interaction of individual students' understanding and knowledge transforms into collective learning which surpasses individual learning and reshapes the whole system [12]. Therefore, ignoring importance of collective learning limits development of students. Another reason to promote complexity based project approach is student heterogeneity because students belong to different ethnicities and have different CGPA. Mainemelis, Boyatzis [13] suggested that adaptive and iterative strategies are appropriate to design projects for different learning styles and skill levels as they emerge from learning processes. Finally, complexity-based strategies contribute to prepare students for a complex and intertwined world as characterized by Barnett and Hallam [14]. Upcoming knowledge is becoming uncertain, unpredictable, and challengeable; therefore, the educational systems should develop adaptivity, self-reliance, collaborative and intuitive learning in students so they can prosper and adapt to uncertain world. Efforts have been made to promote complexity-based education system for better course design [15, 16] and better learning of students at different levels [17-19].

However, the current work summarizes the learning effect in a CAS from centralized decision based collaborative project to decentralized system in successive semester projects, where the heterogeneity of students varied in both experimental setups. In the following sections, the core elements of the project designed through this strategy are defined with key adaptations realized throughout the project development.

## 3. Complexity-based strategy for a term-project

In 2015, complexity-based strategy to design and evaluate an undergraduate term project was adopted to foster collaborative and innovative learning among students. The project was assigned to final year students of Industrial Engineering Department, University of Engineering and Technology, Taxila (Pakistan). Contrary to traditional term projects where three to five students form a group which usually comprise of friends whom they have comfort level, and where each group in a class tries to outperform their fellow beings, this project was unique in a way that the

whole class including the instructor worked As-One team to achieve desired objectives. Analogous to emergent CAS, the project was redesigned in 2017 with improvements and changes in approach, learned by iterative and recurring activities, as mirrored by Davis and Sumara [1] and presented in Figure 1.



Figure 1. Conditions to foster emergence in educational CAS. [1]

### 3.1. Educational and learning goals of project

To foster emergence in the system, learning and educational goals of the project were defined. Learning perspective includes;

(i)     Specialization by facilitating knowledge development through sharing among students, at the same time promoting diversity.

(ii)    Trans-level learning of students by collaboration and decentralized control.

(iii)   Influencing project organization and dynamics through enabling constraints by providing sources of distraction and unpredictability while maintaining the coherence and focus.

The educational objectives of the project were defined by analysing the salient aspects of the course (Computer Integrated Manufacturing) and incorporating them in the project contents, both exogenously (e.g. project design state of art; employability requirements) and endogenously (e.g. course strategic plan; students' skills).

### 3.2. Iteration I: As-One term Project

In 2015, a term-project was designed aimed at engaging students in collaborative dynamics, requiring them to design an innovative tool for Incremental Sheet Forming, based on the specifications provided by the instructor. The project was articulated in stages terminating in an end-product prototype.  In this phase, class was divided into seven teams based on their interests and working compatibility among students. The instructor was the centralized decision body, who defined the narrow project objectives. Whole class was a single team with seven independent but sub-teams including Design, Analysis, Process Planning, Manufacturing, Computer Programming, Project Management (PM) and Concurrent Engineering (CE) Teams [20].

### 3.2.1. Conflicts

After initial excitement faded, the system underwent issues of team collaboration, knowledge transferring and criticism on individuals as well as their teams. Independence in the team structure and functionality created sense of superiority in few teams and affected overall performance of the system. This environment hindered the project progress and led to demotivation of individuals.

### 3.2.2. Emergence of self-adaptable team structure.

To resolve conflicts among teams, enhance collaboration and foster impartiality in system's environment, teams were realigned in such a way that members of CE and PM submerged into other teams and formed integrated nut-bolt structure, shown in Figure 2. This realignment enhanced the learning of instructor to fostered conflict resolution ability and cope with unpredictable scenarios.



(a) Island Team Structure          (b) Integrated (nut-bolt) Team Structure

Figure 2. Evolution in team structure [20]

### 3.2.3. Self-Organization and Trans-Level Leaning.

After self-adaptation of integrated team structure through iterative learning of individuals and instructor, individuals defined their intra-team functionality based on their interests and knowledge which aided in collective learning and improved knowledge sharing process.

### 3.2.4. Project completion and evaluation.

Consequently, in phase 1 of term-project based CAS, the objectives of the system were achieved including product prototype manufactured according to instructor's specifications, as shown in Figure 3, self-organization among individuals, adaptability, conflict management and dealing with complex human behaviour. To encourage team-work and collaboration, the project was evaluated based on collective performance of teams. Qualitative assessment of prototype portrayed the combined effort of teams and their level of learning.



Figure 3. Design and Prototype of tool [20]

### 3.2.5. *Assessment as a 'window onto the system'*

To evaluate this system, a survey was conducted from the participants which showed a positive response and appreciation from students towards the methodology. Few results are shown in Figure 4. This encouraged instructor to continue the evolution of As-One team term-project based CAS, incorporating adaptability to heterogeneous individuals.



j) Response to "Overall how do you rate this learning approach".

h) Response to "Such practice will be beneficial for students in exploring new horizons of knowledge in future".

Figure 4. Survey results [20]

### 3.3. Iteration II: As-One Term Project

The iterative learning and emergent behaviour of CAS allowed redesign of project activities in 2017, whereby knowledge and experience of instructor was transferred to the heterogeneous entities. However, the characteristics of system were found similar i-e unpredictability of human behaviour, conflicting opinions, trans-level learning and external constraints. Overall goal and objectives of the project were same as Iteration 1, except for some methodical alterations.

#### 3.3.1. Decentralized decision system

Based on previous learning of instructor, team structure was pre-defined (i-e integrated nut bolt sub-teams). However, team formation and individual selection was decentralized (i-e individuals selected their team partners, number of sub-teams and role of each sub-team), as group formation is considered critical in individual's acceptance of group activities as well as the success of collaborative learning process [21]. Decentralized decision system provided more autonomy to individuals, thereby, enhanced collaboration and trans-level learning ability. Thus, students self-organized into five sub-teams (CE, PM, Design, Manufacturing, Sales & marketing) and assigned roles and responsibilities, using their cognitive abilities.

#### 3.3.2. Conflicts in product selection

Decentralized decision system at one side offered autonomy, while at other side, challenged selection of a feasible product, formation of sub-teams and distribution of roles. However, teams conducted market surveys for need identification and short-listed few product proposals. With combined consensus, they selected a product (lassi -maker; keeping in view the summer season demand). As the project proceeded, teams faced challenges such as unavailability of machining setup, material procurement issues and lack of engineering information which led to the rejection of product. The project underwent delays due to conflicts in product selection. However, this issue resolved using iterative learning of students in identifying accurate product and its requirements (hybrid pen with dual functionality was selected to be mass produced and sold).

#### 3.3.3. Adaptability to the external constraints

As mentioned above, due to delays in product selection, external pressure on teams increased and aggregated with demotivation. To cope with external environment, there was need of support and

107

adaptability from instructor to accommodate students in this situation. Consequently, instructor allowed students to change the product type with the assessment of economic feasibility and break-even analysis. Henceforth, it sustained CAS by adapting to external constraints. With the consent of instructor, students worked on design and manufacturing of hybrid pen. However, as the system evolved and unpredictable conditions unfolded, students found difficultly in manufacturing of product according to their design and specifications (unavailability of mass-production facilities). Therefore, to complete the project, system elements needed to be flexible and understand the limitations of each other. At that time, instructor allowed to manufacture a single prototype with available facilities.

### 3.3.4. Win-win situation

Consequently, the students succeeded in achieving their goal and developed a prototype, shown in Figure 5. Since the objective of the system was learning, therefore, adaptability and coordination allowed the system to sustain and nourish. Although, decentralized decision system offered many challenges to the students, whereby delay in achieving the goal and compromise on the outcome. However, this system gave autonomy and independence to cope with unpredictable conditions and find feasible solution from constrained environment. Hence, the project completed with a win-win situation by accommodating all entities of the system. (instructor's objective as well as students' goal). The author could not obtain the evaluation of this phase as at that time the department was in transition state from content-based to outcome based education. For this purpose, a pilot study was also made to propose a framework for the department [22].



Figure 5. Design and Prototype of Hybrid pen

### 4. Findings

### 4.1. Time Management Issue

In retrospect, weekly progress reports in both iterations indicated that initial 2/3rd of the time was spend on decision making and conflict resolution while the rest of the work was done in last few weeks. This can be attributed to better collaboration, intuitive learning, clarity of objectives, and time pressure at later stage of the project.

### 4.2. Situational Reaction of Individuals

Evaluation of instructor after mid-session and end-of-session portrayed the reaction of students in different situations. They evaluated low score to the instructor in mid-session where they found difficulty in working in complex team, and uncertainty of achievement. While, after experiencing through the whole project, the end-session evaluation of the instructor was excellent.

### 4.3. Unpredictability and role of external stimuli

The possibility of achieving the goals was found high in the beginning of the project, due to increased motivation and excitement of students. However, as the system constraints unfolded, they imposed hindrance in accomplishment of the marked goals, thus lowering the motivation and

achievement possibility. At that point, external stimuli (guidance and support of instructor) helped in motivating individuals thereby improving the performance of individuals in last stages.

### 4.4. Decentralization leading to increased chaos

It has been observed that in 1st iteration, there was increased adaptability, as well as, the target was achieved up-to-mark. While in 2nd iteration (decentralized decision system), there was increased chaos and conflicts and the target achieved was also below the mark.

### 4.5. Workload Balancing

Workload balancing among sub-teams and individuals is critical. Carefully designed project with thoughtfully defined evaluation criteria could seek equal distribution of workload. Besides work balancing, another observation is that despite executing the project using Concurrent Engineering (CE) technique, there remains uneven distribution of workload over time. Some teams work more at the earlier stages while others virtually remain idle at that time, and some teams work more at the later stages of the project. But then this happens in every project, not all tasks start at once. Still there is window for improvement to engage students throughout the semester.

### 5. Conclusion

Higher educational systems are analogous to CAS as they are emergent, iterative and based on trans-level learning. To cope with actual world of complexity and unpredictability, the educational systems should provide such conditions to foster collaboration, adaptability, self-organization and cognitive abilities in individuals. This research has presented an experimental setup in education system to provide such environment to individuals. Case studies were designed based on term projects where students' attitude was analysed in centralized as well as decentralized decision system. Results of which showed that adopting complexity based pedagogical approaches aid in nurturing of individuals to deal with CAS. Further, the learning of instructor showed that educational systems are also evolutionary and can be improved with combined effort of students and instructors.

To continue study on education complex system, the instructor aims to design a term project in such a way that complete autonomy is given to students where they define their team structures and roles.

### References
1. Davis, B. and D.J. Sumara, *Complexity and education: Inquiries into learning, teaching, and research*2006: Psychology Press.
2. Frei, R. *A complex systems approach to education in Switzerland*. in *ECAL*. 2011.
3. Davis, B. and D. Sumara, *Complexity science and educational action research: toward a pragmatics of transformation.* Educational action research, 2005. **13**(3): p. 453-466.
4. Jacobson, M.J. and U. Wilensky, *Complex systems in education: Scientific and educational importance and implications for the learning sciences.* The Journal of the learning sciences, 2006. **15**(1): p. 11-34.
5. Morrison, K. *Complexity theory and education*. in *APERA Conference, Hong Kong*. 2006.
6. Miller, J.H. and S.E. Page, *Complex adaptive systems: An introduction to computational models of social life*2009: Princeton university press.
7. Newell, C., *The class as a learning entity (complex adaptive system): An idea from complexity science and educational research.* SFU Educational Review, 2008. **1**.

8. Fuite, J. *Network education: understanding the functional organization of a class*. in *Complexity, Science & Society Conference, Liverpool, UK*. 2005.

9. Helbing, D., et al., *Saving human lives: what complexity science and information systems can contribute.* Journal of statistical physics, 2015. **158**(3): p. 735-781.

10. Kempf, K., et al., *Karl G. Kempf.* Chemical Engineering, 2009. **33**: p. 2159-2163.

11. Argyris, C., *Double loop learning in organizations: New forms for turbulent environments.* Journal of Marketing, 1977. **55**: p. 77-93.

12. Van Der Vegt, G.S. and J.S. Bunderson, *Learning and performance in multidisciplinary teams: The importance of collective team identification.* Academy of management Journal, 2005. **48**(3): p. 532-547.

13. Mainemelis, C., R.E. Boyatzis, and D.A. Kolb, *Learning styles and adaptive flexibility: Testing experiential learning theory.* Management learning, 2002. **33**(1): p. 5-33.

14. Barnett, R. and S. Hallam, *Teaching for supercomplexity: A pedagogy for higher education.* Understanding pedagogy and its impact on learning, 1999. **137**.

15. Aron, D.C., *Developing a complex systems perspective for medical education to facilitate the integration of basic science and clinical medicine.* Journal of evaluation in clinical practice, 2017. **23**(2): p. 460-466.

16. Fabricatore, C. and M.X. López, *Complexity-based learning and teaching: a case study in higher education.* Innovations in Education and Teaching International, 2014. **51**(6): p. 618-630.

17. Dubovi, I., et al., *Nursing students learning the pharmacology of diabetes mellitus with complexity-based computerized models: A quasi-experimental study.* Nurse education today, 2018. **61**: p. 175-181.

18. Qian, M. and K.R. Clark, *Game-based Learning and 21st century skills: A review of recent research.* Computers in Human Behavior, 2016. **63**: p. 50-58.

19. Shernoff, D.J., et al., *Student engagement as a function of environmental complexity in high school classrooms.* Learning and Instruction, 2016. **43**: p. 52-60.

20. Qamara, S., et al., *Application of Concurrent Engineering for Collaborative Learning and New Product Design.*

21. Isotani, S., et al., *An ontology engineering approach to the realization of theory-driven group formation.* International Journal of Computer-Supported Collaborative Learning, 2009. **4**(4): p. 445-478.

22. Manzoor, A., et al., *Transformational model for engineering education from content-based to outcome-based education.* International Journal of Continuing Engineering Education and Life Long Learning, 2017. **27**(4): p. 266-286.

**Give Me the BITSTS: A Battery of Integrated Tests of Systems Thinking Skills**
Joe A. Wasserman
West Virginia University

ABSTRACT: Systems thinking, or the ability to understand, make predictions about, and intervene in complex systems, is a crucial skill in various disciplines. To date, most measurement instruments of systems thinking have been focused on a relatively narrow range of systems thinking competencies, time-consuming or difficult to administer and score, and/or domain-dependent. By administering a questionnaire containing five existing systems thinking measurement instruments to 465 participants, this study explored the potential for incorporating existing measurement instruments into a multi-dimensional, easy to administer and score, domain-independent measure of systems thinking. The resulting Battery of Integrated Tests of Systems Thinking Skills (BITSTS) contained 33 items that loaded on six intercorrelated factors in a two-parameter multi-dimensional Item Response Theory model. These items ranged in difficulty from easy to moderately challenging. Although items did not appear to conform to any hypothesized factor structures based on existing typologies of systems thinking skills, item cross-loadings suggested that items were measuring related systems thinking skills, not merely instrument-specific constructs. Although more work is needed to confirm its factor structure and to investigate its validity, the Battery of Integrated Tests of Systems Thinking Skills is a promising instrument for measuring systems thinking skills.

Author Note

Joe A. Wasserman, Department of Communication Studies, West Virginia University.

Correspondence concerning this paper should be addressed to Joe A. Wasserman, Department of Communication Studies, West Virginia University, 108 Armstrong Hall, P.O. Box 6293, Morgantown, WV 26506-6293. E-mail: jowasserman@mix.wvu.edu

# 1. Introduction

Systems thinking—the ability to understand, make predictions about, and intervene in complex systems—has been identified as a crucial skill for learners and practitioners across many disciplines, from sciences and engineering [1] to the social sciences [2]. Because systems are diverse, they must be conceptualized broadly to match: systems are groups of things that are interrelated [3]. Given the importance of systems thinking to so many disciplines, there is a need for measures of systems thinking that are not tied to any particular domain. To date, measurement attempts have been relatively focused, reflecting either a narrow range of systems thinking competencies or expertise in a particular domain. After characterizing systems thinking in greater detail and reviewing existing measures of systems thinking, this study integrates existing, easy-to-score measures into a more complete, domain-independent, multi-dimensional measure of systems thinking skills.

# 2. Review of Literature

Given the breadth of scholarship on systems thinking, the term has been derided as having come "to mean little more than thinking about systems, talking about systems, and acknowledging that systems are important" [4, p. 252]. Nevertheless, more specific typologies of systems thinking have been formalized. These typologies, reviewed below, are useful for researchers and practitioners in that they (a) provide structure to the otherwise potentially nebulous concept of systems thinking and (b) identify a range of more discrete abilities that can all be considered part of systems thinking more broadly.

## 2.1. Systems Thinking: Definitions and Typologies

The diverse literature on systems thinking contains multiple conceptualizations of systems thinking [5], [6]. Given this conceptual diversity, it is productive to synthesize the literature to identify commonalities among systems thinking competencies. Toward this end, we focused on two typologies. The first, Ben-Zvi Assaraf and Orion's [7] list of eight discrete systems thinking skills, was selected for its established influence, having been cited 380 times as of the time of writing [8]. The second, Stave and Hopper's [9] taxonomy of systems thinking, was selected for having drawn on a framework that is well-established in the learning sciences: Anderson and Krathwohl's [10] revision of Bloom's taxonomy of educational objectives.

Ben-Zvi Assaraf and Orion [7] identified eight systems thinking skills within a three-tiered hierarchy. The lowest level of this hierarchy focuses on identifying system components (skill 1), the second on synthesizing these components (skills 2–5), and the third on implementing systems thinking skills (skills 6–8). These skills are as follows.

1. "The ability to identify the components of a system and processes within the system" (p. 523) comprises two distinct skills, (a) identifying the entities that are crucial to a given system and (b) identifying the processes through which these entities change over time.

2. "The ability to identify relationships among the system's components" (p. 523) is distinct from the first skill in that it emphasizes *interrelations* among entities, as opposed to processes that occur within a given entity.

3. "The ability to identify dynamic relationships within the system" (p. 523) entails recognizing not only that entities are interrelated, but that these interrelations cause changes to entities over time.

4. "The ability to organize the systems' components and processes within a framework of relationships" (p. 523) integrates the first three systems thinking skills. Whereas the first three skills emphasize identification, the third focuses on the organization or integration of entities, processes, and interrelations into a holistic framework or system. These first

four systems thinking skills are closely tied to the aforementioned conceptualization of systems as groups of interconnected entities see [3].

5. "The ability to understand the cyclic nature of systems" (p. 523) means recognizing the feedback loops within systems through which entities mutually influence each other over time, rather than proceeding from a starting point to an ending point.

6. "The ability to make generalizations" (p. 523) involves synthesizing understandings of entities, processes, or interrelations within a given system, or of a whole system. These generalizations might include classifying entities or interrelations within a system. Generalizing about a system more holistically might involve characterizing that system in terms of dynamic changes over time, feedback loops among interrelated entities, or temporal delays embedded in these interrelationships and processes.

7. "Understanding the hidden dimensions of the system" (p. 523) involves recognizing that the entities, interrelations, and processes within a system may not be immediately apparent from surface inspection.

8. "Thinking temporally: retrospection and prediction" (p. 523) entails understanding that a given state of a system is a product of past interactions and processes, while future system states will be due to present dynamics.

Stave and Hopper [9] developed a seven-tiered, hierarchical taxonomy of systems thinking that they aligned with Anderson and Krathwohl's [10] revision of Bloom's taxonomy of cognitive processes of remembering, understanding, applying, analyzing, evaluating, and creating. Although their taxonomy overlaps in part with Ben-Zvi Assaraf and Orion [7], it is novel in its alignment with an existing, well-established taxonomy of learning and describes additional systems thinking skills. This taxonomy of systems thinking skills [9]—and their corresponding cognitive processes from Anderson and Krathwohl [10]—are as follows.

1. *Recognizing interconnections* includes "seeing the whole system, understanding how parts relate to and make up wholes, and recognizing emergent properties" (p. 14) and maps onto Ben-Zvi Assaraf and Orion's [7] skills one through three, and is related to conceptualizations of systems as groups of interrelated entities [3]. It more specifically involves recognizing how a system is comprised of its component parts, as well as the way in which those parts are interrelated to form the whole system.

2. *Identifying feedback* includes "recognizing and identifying interconnections and feedback" (p. 14) and maps onto Ben-Zvi Assaraf and Orion's [7] skills two, three, and five. This skill entails recognizing the causal relationships among entities in a system, including indirect chains of causal relationships and feedback loops.

3. *Understanding dynamic behavior* includes "understanding the relationship between feedback and behavior, including delays" (p. 14) and maps onto Ben-Zvi Assaraf and Orion's [7] skills three, five, and six. Beyond recognizing interrelationships, this systems thinking skill involves understanding how system structures—such as feedback loops and delays—produce behaviors that are characteristic of that system.

4. *Differentiating types of variables and flows* includes "understanding the difference between rates and levels" (p. 14) and is not contained in Ben-Zvi Assaraf and Orion [7]. This skill entails quantitatively characterizing entities and interrelationships of a system in terms of stocks and flows. Stocks are quantities of material or information that have accumulated in a particular place, whereas flows are movements of this material or information into, out of, or between stocks [5]. Stocks are characterized by their levels, or how much material or information they currently hold. Flows are characterized by rates,

or the speed at which material or information moves through them.

5. *Using conceptual models* includes "using general systems principles to explain an observation" (p. 14) and maps onto Ben-Zvi Assaraf and Orion's [7] skills six and eight. This skill involves applying general systems concepts to observations of particular systems in order to explain how system structures produced a particular system state.

6. *Creating simulation models* includes "describing connections in mathematical terms and using both qualitative and quantitative variables" (p. 14) and is not contained in Ben-Zvi Assaraf and Orion [7]. This advanced systems thinking skill incorporates all of the prior skills and entails creating simulations of a particular system by formalizing a description of that system's entities and their interrelations. Creating simulations also involves comparing them to an external standard, such as observations of a system, other formal models, or other simulations.

7. *Testing policies* includes "using simulation to test hypotheses and develop policies" (p. 14) and is not contained in Ben-Zvi Assaraf and Orion [7]. This advanced systems thinking skill can involve creating simulations, but otherwise incorporates all of the prior skills. Testing policies entails using simulations to identify potential changes to a system that would produce desirable future states of the system.

## 2.2 Systems Thinking Measurement

A number of diverse measurement instruments of systems thinking that reflect various aspects of the aforementioned systems thinking skills have been developed. To maximize the utility of systems thinking measures, I propose that instruments should meet the following seven criteria: (a) reflect a wide range of systems thinking skills, (b) capture a wide range of systems thinking abilities, (c) not require knowledge of a particular content area—i.e., be domain-independent, (d) be easy to administer, (e) be easy to score, and (f) be as short possible so as to minimize burden on participants while satisfying the preceding criteria. The following review of systems thinking measures focuses on measures that meet most or all of these criteria.

**Stocks and flows.** An early attempt at measuring domain-independent systems thinking ability—and understanding of stocks and flows in particular—involved three open-ended, domain-independent graphing tasks [11]. Each item graphically and textually depicted a stock and flow of some material over time. Responses were elicited as freehand-drawn graphs and scored for correctness. These graphing-response measures of stock-and-flow understanding have been adapted as closed-ended, domain-independent, numeric-response items [12], [13]. In these closed-ended versions of stock-and-flow tasks, items graphically depict one inflow and one outflow, eliciting responses by asking participants to identify the moments of greatest inflow, greatest outflow, greatest stock, and least stock.

**Measures of systems thinking skills in earth sciences.** A number of more domain-dependent, closed- and open-ended measures of systems thinking skills have been developed in the context of primary science education in earth sciences [7] and biology [14]. These measures include closed-ended true-or-false questions about system components, as well as simple, dynamic, and cyclical relationships; open-ended system drawings and concept maps that integrate system components and relationships into more-or-less coherent frameworks; open-ended word association tasks; semi-structured interviews on making generalizations about systems, identifying hidden aspects of systems, and generating explanations and predictions; and a semi-structured repertory grid [7].

**Systems Thinking Assessment.** The Systems Thinking Assessment (STA) is a closed-ended, domain-independent multiple choice test that was iteratively and cyclically developed to

measure middle school students' systems thinking abilities [15]. This 52-item test was developed in Greek and measures identification of system components, reasoning about causal relationships, recognizing how system structures produce system behaviors, and reasoning about flows of matter that include feedback loops [16].

**Complex Systems Concepts Inventory.** The Complex Systems Concepts Inventory (CSCI) is a domain-independent measure of understanding general complex systems concepts that includes both closed-ended multiple choice and open-ended free response items [17]. Complex systems are characterized by individual actors or agents from whose relatively simple interactions emerge complex, macroscopic patterns [18]. As such, the CSCI measures understandings of general systems properties such as emergence, feedback, and self-organization; explaining system behaviors; and making predictions about complex systems.

**Systems Thinking Skills and Systems Thinking Measures**

Based on aforementioned systems thinking skills typologies, the reviewed measures of systems thinking were categorized (Table 1). Each existing measure reflects a subset of systems thinking skills. In combination, these measures reflect the majority of basic to moderate-level systems thinking competencies. By combining these measures, it should be possible to arrive at a multi-dimensional measure of systems thinking skills that captures this range of competencies. The categorizations of measures in Table 1 present competing hypothesized factor structures for the systems thinking measures examined in this study.

Table 1. *Measures of systems thinking categorized by the skills they are hypothesized to measure*

| Ben-Zvi Assaraf and Orion [7] | Stave and Hopper [9] | | | | | | |
|---|---|---|---|---|---|---|---|
| | 1. Recognize interconnections | 2. Identify feedback | 3. Understand dynamic behavior | 4. Differentiate types of variables and flows | 5. Use conceptual models | 6. Create simulation models | 7. Test policies |
| 1. Identify components and processes | STA | STA | STA | STA | | | |
| 2. Identify relationships | GDN, STA | STA | STA | STA | | | |
| 3. Identify dynamic relationships | GDN, STA | STA | STA | STA | | | |
| 4. Organize within framework | | | | | | | |
| 5. Understand cycles | STA | CTQ, STA | STA | CTQ, STA | | | |
| 6. Make generalizations | | | CSCI | | | | |
| 7. Understand hidden dimensions | | | CSCI | | | | |
| 8. Retrospection and prediction | | | CSCI | | | | |
| N/A | | | | DS | | | |

*Note.* See text for acronyms of systems thinking measures.

**3. Method**

**3.1. Participants**

Participants were 465 English-speaking individuals age 18 and over. Of the 452 who

reported their age, ages ranged from 18 to 97 years ($M = 27.1$, $SD = 12.6$). Of the 458 who reported their sex, 189 (41.3%) were male, 251 (54.8%) were female, and 12 (2.6%) were trans or nonbinary. Of the 454 who reported their race, 374 (82.4%) were White, and the next most common race was Asian ($n = 26$, 5.7%). Of the 456 who reported their highest education, 244 (53.5%) had completed high school or some college, 130 (28.5%) had completed a two- or four-year degree, and 81 (17.8%) had obtained a graduate or professional degree.

## 3.2. Procedure

Following approval from the university's Institutional Review Board, participants were recruited from a student email list at a large Mid-Atlantic university, social media, online forums, and online databases of psychology studies so as to obtain participants with a wide range of systems thinking abilities. Participants completed an online questionnaire containing items measuring endorsement of systems concepts, systems thinking skills, and demographics.

## 3.3. Measures

**Systems thinking skills.** A range of systems thinking skills were measured using five existing measures of systems thinking that included both open- and closed-ended items, to be combined into a single Battery of Integrated Tests of Systems Thinking Skills (BITSTS). All measures were discussed in the preceding literature review.

*Department store task.* The DS [13] contains four numeric free-response questions on a graph of flows of individuals entering and exiting a department store over a 30-minute period. Answers within ±1 of the correct response were scored as correct.

*Groundwater System Dynamic Nature Questionnaire.* The GDN [7] contains eight true-or-false (*agree*, *uncertain*, or *disagree*) items on the interrelationships between water and other elements of the hydro-system.

*Cyclic Thinking Questionnaire.* The CTQ [7] contains six true-or-false (*agree*, *uncertain*, or *disagree*) items on the cyclical nature of the stocks and flows of the hydro-system.

*Systems Thinking Assessment.* The STA [15], [16] was translated from Greek into English and contains 29, 4-option multiple choice questions on system components, interrelationships among components, temporal causal dynamics, and feedback loops.

*Complex Systems Concepts Inventory.* The CSCI [17] contains two closed-ended puzzles that ask participants to implement system rules; two closed-ended, 4-option multiple-choice questions on predicting outcomes of system dynamics; and eight open-ended questions on explaining system dynamics. Open-ended responses were scored according to a codebook developed by the first author based on Tullis and Goldstone [17] and guidance by Tullis. After six hours of coder training, the first author and a research assistant coded a random subsample of 50 cases, achieving 100% agreement on one item, 98% agreement with interrater reliabilities of Cohen's $\kappa = .94$, .66, and .00 on three items, and Cohen's $\kappa$ between .70 and .81 on the remaining four items. The two lowest values were deemed acceptable because they were obtained with only one coder disagreement on items that were rarely scored as correct. After resolving remaining disagreements via discussion, the remainder of the data were divided in half and independently scored by the first author and research assistant.

## 3.4 Data Analysis

Questionnaire pages on which participants spent less than one second per item were coded as missing to eliminate poor-quality data due to speeding [19]. Because items were dichotomous (0 = *incorrect*, 1 = *correct*), they were treated as reflective indicators of latent variables in item response theory (IRT) models [20] using diagonally weighted least squares (WLSMV). IRT estimates latent ability based on responses to categorical items and the

likelihood of correct responses to each item as a function of individuals' latent ability level. Measures' unidimensionality was tested via confirmatory factor analysis and inspection of overall model fit and local fit as indicated by items' $R^2$. If unidimensionality was implausible, the dimensionality of each measure in isolation was explored by exploratory factor analysis (EFA). The number of factors to extract was decided by evaluating Eigenvalues, model Chi-square difference tests, and the interpretability of rotated loadings. One-parameter IRT models were compared to less parsimonious two-parameter IRT models via model difference tests and comparison of BIC. For each unidimensional IRT model, a subset of the best-performing items reflecting a wide range of abilities were retained [20]. Next, these reduced unidimensional systems thinking measures were combined and their dimensionality explored via EFA.

## 4. Results

For DS [13], a two-factor, one-parameter model containing all four original items was retained (Table 2: http://tiny.cc/BITSTS18). These factors were positively correlated, $r = .43$ ($p < .0001$). The first factor was interpreted as reflecting the ability to read a graph and as such were excluded from subsequent analyses, the second as systems thinking about relationships between stocks and flows. Negative standardized threshold values indicate that the first factor contained relatively easy items, and positive standardized threshold values indicate that the second factor contained relatively difficult items. For the GDN [7], a one-factor, one-parameter model containing five of the initial eight items was retained (Table 3: http://tiny.cc/BITSTS18). This measure was interpreted as recognizing interconnections and relationships between system entities in the context of the hydro-cycle. Negative item difficulties indicate that these items were relatively easy. For the CTQ [7], a one-factor, two-parameter model containing five of the initial six items was retained (Table 4: http://tiny.cc/BITSTS18). This measure was interpreted as understanding cycles, feedback, and flows in the context of the hydro-cycle. The range of items' difficulties indicates that the measure identifies a range of abilities. For the translated STA [15], [16], 15 items from the original 29 were selected from the retained 24-item, one-factor, two-parameter model (Table 5: http://tiny.cc/BITSTS18). This measure was interpreted as identifying and recognizing system components, relationships and interconnections, understanding temporal causal dynamics in everyday contexts. The final items were selected to maximize the breadth of the test information function by including high-discrimination items with the largest range of difficulties. Nevertheless, negative item difficulties indicate that these items are relatively easy. For the CSCI [17] a two-factor, two-parameter model containing nine of the original 12 items was retained (Table 6: http://tiny.cc/BITSTS18). These factors were not significantly correlated, $r = .23$ ($p = .076$). The first factor was interpreted as systems thinking about complex systems dynamics such as emergence of macroscopic phenomena from (sometimes hidden) micro-level interactions, the second as making predictions about the long-term consequences of local interactions for a whole system. Mostly positive standardized threshold values indicate that the items on both factors were relatively difficult.

To investigate the overall dimensionality of the preceding measures of systems thinking skills, an EFA was estimated using oblique geomin rotation to explore the dimensionality of these items in combination. After dropping one STA and two CSCI items for having empty cells in a bivariate table with two or more other items. A six-factor model was retained (Table 7: http://tiny.cc/BITSTS18). All factors were correlated with at least one other (Table 8: http://tiny.cc/BITSTS18). These items in combination comprised the Battery of Integrated Tests of Systems Thinking Skills (BITSTS). Although each factor's strongest loadings were comprised of items from single measures, all measures except DS exhibited substantial cross-loading.

117

Although cross-loadings suggest an underlying factor structure not attributable to constructs unique to each measure, neither of the potential hypothesized factor structures (see Table 1) appear to be clearly supported.

## 5. Discussion

This study examined the potential for existing measures of systems thinking skills to be shortened and combined into a Battery of Integrated Tests of Systems Thinking Skills (BITSTS) to measure a range of systems thinking skills and abilities. The dimensionality of the DS [13], the GDN [7], the CTQ [7], the translated STA [15], [16], and the CSCI [17] were tested and explored. Measures with items with high information functions ranging in difficulty were combined. Although it was predicted that the factor structure of the BITSTS would reflect typologies of systems thinking skills in the literature, the six-factor structure of the BITSTS appeared to primarily reflect the measures from which items originated, rather than theorized typologies of systems thinking skills. Overall, these results suggest that the BITSTS may hold promise as an instrument for measuring a range of systems thinking skills.

It was proposed that in order to maximize the utility of a measure of systems thinking skill, an instrument should: (a) reflect a wide range of systems thinking skills, (b) capture a wide range of systems thinking abilities, (c) not require knowledge of a particular content area—i.e., be domain-independent, (d) be easy to administer, (e) be easy to score, and (f) be as short possible so as to minimize burden on participants while satisfying the preceding criteria. The BITSTS appears to meet these criteria. The BITSTS's six-factor structure, in which all factors were significantly correlated with one-to-four other factors, suggests that it reflects multiple distinct—but related—systems thinking skills. Additionally, the range of item difficulties indicates that it captures a breadth of ability levels. Although it could be argued that GDN and CTQ reflect domain-specific understandings of the hydro-cycle, the weak and non-significant correlation between factors three and four suggest that they are not intimately related. The BITSTS can be administered as a questionnaire, and the closed-ended items are easy to score. At 33 items, the BITSTS is dramatically shorter than the initial set of systems thinking measures.

## 6. Conclusion

The BITSTS contained items of primarily low difficulty, indicating that it is best at discriminating between individuals of relatively low systems thinking ability and may be limited in its utility for capturing higher abilities. Given the established difficulty of systems thinking [21], measuring low-level systems thinking is still of value. Future research should include additional, more challenging items, such as closed-ended versions of other stock-and-flow tasks [11]. In addition to confirming the factor structure of the BITSTS, it is important to establish its validity. In particular, the BITSTS should be able to predict certain outcomes (e.g., performance on systems dynamics tasks) and be predicted by certain experiences (e.g., taking a class in complex systems, playing complex boardgames). Moreover, it would be important to explore the components of intelligence to which systems thinking are related, and to simultaneously establish that performance on the BITSTS cannot be reduced to general intelligence.

References

[1] National Research Council, *A framework for k-12 science education: Practices, crosscutting concepts, and core ideas*. Washington, D.C.: National Academies Press, 2012.

[2] G. D. Garson, "Computerized simulation in the social sciences: A survey and evaluation," *Simul. Gaming*, vol. 40, no. 2, pp. 267–279, Apr. 2009.

[3]   L. von Bertalanffy, *General system theory: Foundations, development, applications*, Revised edition. New York, NY: George Braziller, 2015.

[4]   J. W. Forrester, "System dynamics, systems thinking, and soft OR," *Syst. Dyn. Rev.*, vol. 10, no. 2–3, pp. 245–256, Jun. 1994.

[5]   D. H. Meadows, *Thinking in systems: A primer*. White River Junction, VT: Chelsea Green Publishing, 2008.

[6]   P. M. Senge, *The fifth discipline: The art and practice of the learning organization*. New York, NY: Doubleday/Currency, 2006.

[7]   O. Ben-Zvi Assaraf and N. Orion, "Development of system thinking skills in the context of earth system education," *J. Res. Sci. Teach.*, vol. 42, no. 5, pp. 518–560, May 2005.

[8]   "Google scholar." [Online]. Available: https://scholar.google.com/scholar?cites=13443129322880353132&as_sdt=5,49&sciodt=0,49&hl=en. [Accessed: 26-Nov-2017].

[9]   K. Stave and M. Hopper, "What constitutes systems thinking? A proposed taxonomy," in *Proceedings of the 25th International Conference of the System Dynamics Society*, Boston, MA, 2007.

[10]  L. W. Anderson and D. R. Krathwohl, *A taxonomy for learning, teaching, and assessing: A revision of Bloom's taxonomy of educational objectives*. New York, NY: Addison Wesley Longman, 2001.

[11]  L. Booth Sweeney and J. D. Sterman, "Bathtub dynamics: Initial results of a systems thinking inventory," *Syst. Dyn. Rev.*, vol. 16, no. 4, pp. 249–286, Dec. 2000.

[12]  G. Ossimitz, "Stock–flow thinking and reading stock–flow-related graphs: An empirical investigation in dynamic thinking abilities," presented at the 2002 System Dynamics Conference, Palermo, Italy, 2002.

[13]  J. D. Sterman, "All models are wrong: Reflections on becoming a systems scientist," *Syst. Dyn. Rev.*, vol. 18, no. 4, pp. 501–531, Dec. 2002.

[14]  O. Ben-Zvi Assaraf, J. Dodick, and J. Tripto, "High school students' understanding of the human body system," *Res. Sci. Educ.*, vol. 43, no. 1, pp. 33–56, Feb. 2013.

[15]  K. C. Constantinide, "Development and validation of a systems thinking assessment instrument for 10-14 year olds," Doctoral dissertation, University of Cyprus, Nicosia, Cyprus, 2015.

[16]  K. Constantinide, M. Michaelides, and C. P. Constantinou, "Development of an instrument to measure children's systems thinking," in *E-Book Proceedings of the ESERA 2013 Conference: Science Education Research For Evidence-based Teaching and Coherence in Learning*, Nicosia, Cyprus, 2014, vol. 11, pp. 13–23.

[17]  J. G. Tullis and R. L. Goldstone, "Instruction in computer modeling can support broad application of complex systems knowledge," *Front. Educ.*, vol. 2, 2017.

[18]  M. Mitchell, *Complexity: A guided tour*. New York, NY: Oxford University Press, 2009.

[19]  D. Wood, P. D. Harms, G. H. Lowman, and J. A. DeSimone, "Response speed and response consistency as mutually validating indicators of data quality in online samples," *Soc. Psychol. Personal. Sci.*, vol. 8, no. 4, pp. 454–464, May 2017.

[20]  T. Raykov, *Item response theory*. Philadelphia, PA: Course booklet. Statistical Horizons, 2017.

[21]  L. Booth Sweeney and J. D. Sterman, "Thinking about systems: Student and teacher conceptions of natural and social systems," *Syst. Dyn. Rev.*, vol. 23, no. 2–3, pp. 285–311, Jun. 2007.

# Resiliency Preparedness: Using Multi-layer Networks to Assist in Disaster Response Processes

Kameron Grubaugh[1] and Jonathan Roginski[2]

[1] Department of Mathematical Sciences
United States Military Academy at West Point
kameron.grubaugh@usma.edu
[2] Network Science Center
United States Military Academy at West Point
jonathan.roginski@usma.edu

## Abstract

Natural disasters strike every year, causing billions of dollars in damages and the loss of thousands of lives across the world. To reduce the effects of these events, Network Science can enhance the ability to inform decision makers about how to most effectively, in terms of both time and resources, spread information through vulnerable portions of a city. A successful information dissemination plan will reduce the amount of time between government officials declaring emergency information and that same information reaching the average citizen in danger. This case study on Newburgh, New York uses a three-layer network to represent the geospatial layout, informal community structures, and church catchments to represent the existing informal relationship network in the city. Analysis shows that using clergy members is the most efficient way to distribute information to a broad populous because in many cities, Newburgh included, a relatively small number churches encapsulate a majority if the churchgoers city-wide. This allows widespread dissemination of information to occur with little governmental resource expenditure. The research also displays the connection between geospatial network structure and the formation of informal communities in cities, which, when leveraged effectively, can also increase the flow of information in a time of crisis.

## 1 Introduction

Networks can describe everything in the environment around us, from social relationships to the spread of illness. The study of these interactions, or Network Science, is an emerging, theory driven field that focuses the dynamic behavior of the relationships the systems contain [1]. The ability to apply this field to anything makes it a highly valuable skill set to all organizations, including the Department of Defense and the President of the United States. In the most recent National Security Strategy, President Trump expressed his interest in maintaining a military that is capable of preventing the growth of violent extremist networks by addressing the underlying conditions that cause such behavior in our current areas of interest [2]. Network Science can aid in policy

recommendations for executives and commanders operating in the complex environments President Trump refers to. Providing a means of estimating which actions the United States forces can take to reap the most benefit across the region, as modeled by a network, is a valued aspect of the decision making process. If a commander must decide between building a school in Village A or establishing clean water in Village B, he or she may not have a way of estimating which will provide the largest step forwards towards the end goal; the use of Network Science can provide insight for him or her to make that decision. This study hopes to evaluate network resiliency; the more resilient the network, the less likely it is to suffer following an attack, such as a natural disaster.

## 2   Related Work

To properly establish a multi-network of any type, one should first understand the basics of network structure and how to measure appropriate statistics. The most basic goal of Network Science is to understand interactions between objects, such as humans, social groups, or the transfer stations of a power grid. This research revolves around the structure of networks and the relevance of the vertices within, with mention of how new interactions, or the deletion thereof, may affect the network as a whole.

Centrality measures gather varying aspects of the importance of individual vertices. The research will highlight two of the primary single-layer network centrality measures that hold significant importance to analyses: betweenness and closeness. Significant to measuring the success of passing information through a network is the betweenness centrality. This centrality measures the number of shortest paths between two nodes that a given vertex lies on, divided by the number of shortest paths possible [3]. This returns the likelihood a given vertex will be on the shortest path, and inherently determine how likely the vertex is to be involved in communication paths between the two nodes. A common use of this measurement is determining which vertex to attack or eliminate to cause the greatest harm to a network. If you remove the vertex on the most shortest paths, the network dynamic must accommodate the change, making it more difficult for information to flow through the network and often cuts entire vertices off from the remainder of the network.

Another one of Freeman's centrality types that I use is the closeness centrality. Indicative of how quickly information is distributed from a given node, the closeness centrality calculates how connected a vertex is in relation to the other vertices in the system [3]. To calculate the closeness centrality of vertex $i$, take $\frac{N-1}{\sum_j d_{i,j}}$ where $N$ is the number of vertices in the network and $d_{i,j}$ is the length of the shortest path between vertex $i$ and any other vertex in the network, $j$. The closer the measurement is to one, the faster information can be distributed to the rest of the network from that source.

As the decades progress, the theoretical foundation of network structure evolves to reflect advancements in the understanding of such. The seminal paper of Paul Erdős and Alfréd Réyni argued that all networks are random in nature [4]. The paper further explains that each graph is random because all edges are equally likely to be selected in a random draw, or that the formation of a random graph is a stochastic, rather than deterministic, process [4]. An example of a random network as Erdős and Réyni imagined is the interstate highway system. There are a given set of vertices, for instance the cities connected by the system, by which the selection of any edge is equally likely. Under this premise, when a new vertex joins the network, its edges form in a random manner. This method does not, however, capture the reality behind network formation.

This would mean that if a new city wished to form along one of the edges, or interstates, it would randomly connect to other cities using new roadways. In practice, cities tend to form along existing roadways because they allow massive flows of economic goods and people between other important cities.

First captured by Albert-László Barabási and Réka Albert, they note this tendency by stating that as more interactions and participants enter a system, the probability that a given vertex has $k$ edges, or interactions, does not change [5]. The power law equation used to find the probability is $P(k \ links) = k^{-\gamma}$ where $\gamma$ where the power law coefficient specific to that network and will always fall between 2 and 4 [5]. In application this means all graphs will self-organize into a scale-free state, leading the joining vertices to interact with already highly influential vertices. An example of a scale-free network according to the Barabási-Albert model is the network of airports in the United States. If a new airport joins the system, the operators will desire it have flights to the most connected airports, such as Chicago and Atlanta. The presence of an edge connecting this new, small airport to the larger, well-connected one will enable the new airport to transport passengers to a vast array of cities across the network. If it connected randomly upon admission to the network, the new airport may only form relationships with other small airports, leaving it on the periphery of the network where it can only appeal to a small number of customers.

Most recently, Duncan Watts and Steven Strogatz further expanded on the idea of network structure by researching the commonly referenced six degrees of separation notion. Their research produced the Watts-Strogatz model which states that network tend to have high clustering (also referred to as cliquishness) and low path lengths due to the close knit nature of naturally forming networks [6]. This model is especially useful with social networks where several smaller groups of friends or business partners are connected by relatively few edges between them. These bridges between modules lead to shorter paths between cliques because there is a streamlined flow of information along the bridge between groups, decreasing the amount of time it takes for this information to travel through the network to the desired targets. The commonly known test of this model is the "six degrees" test, of which the most common examples are of the "six degrees of Kevin Bacon" game and connecting two randomly selected FaceBook members with only six connections. The notion that a randomly selected person in rural Africa can connect with a randomly selected academic in America with only six connections on a consistent basis would seem near impossible under the randomness of the Erdős-Réyni model and highly unlikely under the Barabási-Albert model; the Strogatz-Watts model, however, demonstrates this tested phenomena.

## 3  Methods

The data collection process of this study took place in three phases. First, initial interviews informed the early hypotheses about which sections of the city should be the focus of a resiliency study. Following initial interviews came the formation of the network with a focus on the most vulnerable sector, which is the area between South Street and the southern limit of the municipality. Lastly, a final round of interviews completed the remaining holes in the layers of the network. Each of the steps is equally important in the process of forming a multi-layer network.

The initial interview phase opened up bountiful insight to the informal networks of Newburgh. The initial ideas of using the government as a starting point quickly changed after both of the initial interviews challenged the actual influence of formal governance structures on the network. At this

point, the emphasis shifted to the informal power structure within the operation of Newburgh. This realization came rather early in the interviews, so many of the follow-on questions focused on uncovering information to help form these networks. To help narrow the discussion, the interviewees laid out where the most troubled sectors of the city were. This helped reached a deeper level of analysis and questioning because the focus narrowed to roughly half of the city.

The formation of the geospatial network required several decisions early in the process. One hypothesis presented in an interview suggested that the presence of one-way roads is correlated to the vulnerability of a region. There are other socioeconomic factors that influence this, so the use of one versus two way roads did not seem appropriate for this study without also conducting a study on the correlations between one way roads and the social factors that place people in these areas. This left the researchers with a decision between treating the width of roads in terms of lanes differently as well. This would give main roads even more influence over a network that already values these main thoroughfares heavily as they connect many of the vertices. To avoid overly emphasizing the importance of a block being connected to a main road, weights are not assigned to roads. All roads are treated equally in this study. To form the network, relationships between nodes that are adjacent to each other are drawn. This places slightly more emphasis on blocks that are larger, but in most cases a large block does not have any more neighbors than those that are standard size, so this did not influence the network greatly.

The final phase consisted of developing the informal influence layers of the network with greater detail and establishing layers for comparison to reality. All layers beyond the geospatial one came from these interviews, to include communities that formed based on informal communities and the catchment of churches. A majority of this information came from collaboration between the Newburgh City Planner and the researchers. Key comparisons will come from the church catchments and informal communities, as well as comparing the computational communities that form to those that formed over decades within Newburgh.

## 4 Results

Interpretation of the data is the most crucial aspect in applying the mathematics that lay the foundation for Network Science. A complete analysis of a multi-layer network includes not only the overall assessment, but also an in depth look at each contributing layer.

### 4.1 Geospatial

Transforming the structure of a city into a network is fairly straight forward; treat each city block as a vertex with an edge connecting it to the blocks which it shares a border with. The formation of the Newburgh network used this same technique on a larger scale, shown in Figure 1. The right hand side of Figure 1 shows the vertices scaled with respect to the hub centrality, highlighting how some blocks have more connections than others. While this method may seem elementary at first glance, it best maintains the integrity while attempting to represent the geographic space that is Newburgh. City blocks, with few outliers, are generally the same size. Therefore, a power holder with influence over more city blocks will typically hold influence over a greater geographic space than one with fewer blocks.

Using solely the geographic network, one would expect the most influential blocks to be those most central in the sector being analyzed. This belief held true to an extent. Many of the blocks
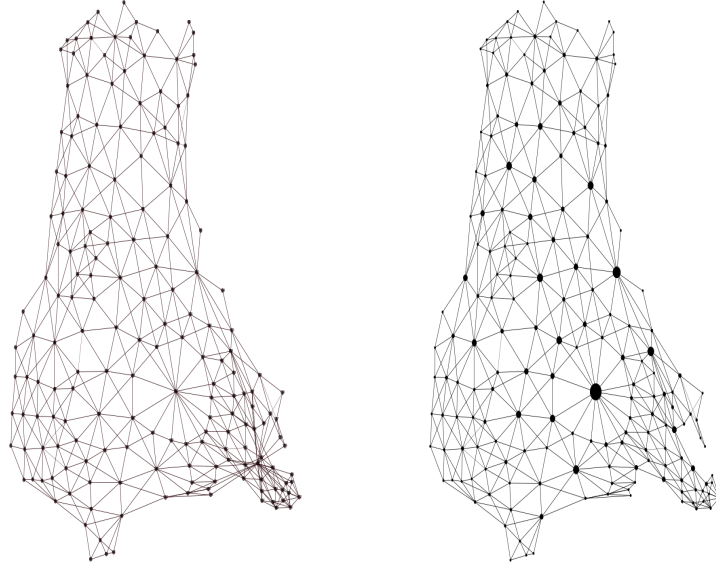
Figure 1: The geospatial network of Newburgh (left) with vertex sizing by betweenness centrality (right).

with the highest closeness and betweenness centrality are disproportionately located to the south east of the geographic center of the sector, slightly towards Block 126, which is the most connected block in the system due to its size. A summary of the centrality measures for the blocks in the top five of at least one of these metrics is shown in Table 1. Looking beyond the sphere of influence for Block 126, similar trends are noticeable with other larger blocks in the western half of the sector. Inherent to the nature of a geographic network is the larger a parcel of land, the greater influence it will have on the system. This layer alone does not provide much insight, but when paired with the following layers it reveals where weaknesses are in the system.

| Vertex | Betweenness Centrality | Betweenness Rank | Closeness Centrality | Closeness Rank |
|--------|------------------------|------------------|----------------------|----------------|
| 126 | 0.2903 | 1 | 0.2442 | 1 |
| 105 | 0.1080 | 8 | 0.2399 | 2 |
| 102 | 0.1185 | 5 | 0.2339 | 3 |
| 101 | 0.0774 | 15 | 0.2327 | 4 |
| 107 | 0.0651 | 20 | 0.2309 | 5 |
| 172 | 0.1788 | 2 | 0.2179 | 16 |
| 168 | 0.1378 | 3 | 0.2189 | 15 |

Table 1: Summary of Closeness and Betweenness Centralities for Highest Ranking Blocks in Each Measure.

## 4.2 Informal Community Structures

Populations of a city tend to form geographic sub-cities, often referred to as communities or boroughs. In Newburgh, these communities formed primarily around ethnicity. During the interview

phase of the project, subjects identified four major communities in the sector emphasized in this study: a predominantly African American region to the north of Broadway, a group of Hispanic households to the south of Broadway, a community of mostly Italians with some Hispanics mixed in the west side of Robinson Avenue, and a middle-class, mostly Caucasian section known as Washington Heights located in the south east region of the sector of study. Interviewees also identified buffer zones where these regions tend to mix together, most notably between the Hispanic and middle-class section. Including the mixture region and areas with no specific community, this makes six separate communities. When allowing mathematics to determine the number of communities in the infrastructure network, the computer determines there are seven unique ones. Altering the resolution for modularity to show just six communities, there are strong similarities between the graphs. A comparison of the graphs is shown in Figure 2.
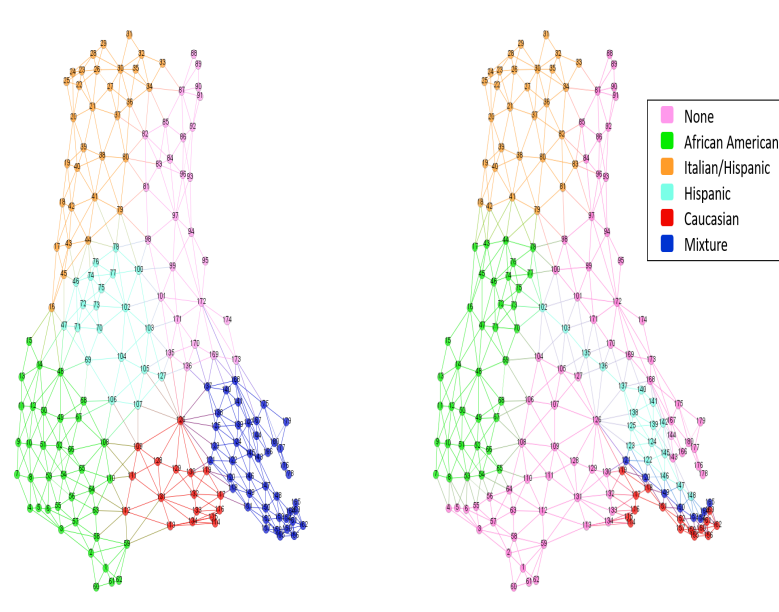


Figure 2: Graph of modularity with resolution 1.3 (left) and a graph of the six informal communities within the sector of study in Newburgh (right).

The striking similarities between the calculated modules of the graph and the real-life informal communities makes it apparent that the simple layout of a city can have an effect on the resilience of the city's network. Additionally, this comparison tells us that the communities are not only tightly grouped, but also that roads serve as the geographic boundary between them in many cases. Washington Heights and the Hispanic neighborhood do not have a major road between them like the other informal community boundaries, but the Washington Heights infrastructure is tightly packed compared to other communities, causing the appearance as a community on the modularity graph.

### 4.3   Church Catchments

Separate communities are often tied together by other factors, such as ministry centers and places of worship. Religious affiliation is closely related to ethnicity, especially with traditional religions

such as Catholicism. In Newburgh, there are five primary churches in the sector of focus: two Catholic, one African Methodist, one non-denominational with Baptist roots, and one Episcopal.
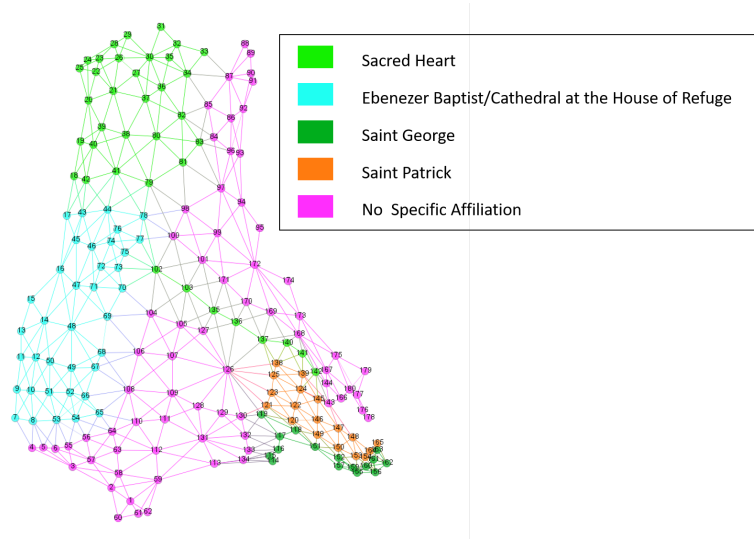


Figure 3: Catchment of the five major churches in the sector as applied to the geospatial network.

Interviews with the local residents identified these churches and confirmed hypotheses on the catchment of such. The African-American population is typically split between the African Methodist church, Ebenezer, and the non-denominational service, The Cathedral at the House of Refuge. Preliminary interviews identified the reverend of the African Methodist church as a potential name for a primary influence holder over the entire network due to his rapport with the African-American community of Newburgh, which led to our emphasis on churches as serving a key role in the system. The Catholic churches catchment is predominantly the Hispanic and Italian regions, with the Italian and western portion of the Hispanic community attending Sacred Heart and the eastern portion of the Hispanic section along with the mixture region attend Saint Patricks. Lastly, the middle-class neighborhood mostly attends the Episcopal service, Saint Georges.

The graph shown in Figure 3 closely resembles that of the communities, prompting the conclusion that the worship leader at their respective church should be considered in the influence structure. Every week, each of these influences has the undivided attention of their congregation for at least forty-five minutes. The position they hold also brings an inherent form of legitimacy with their followers, giving them the ability to spread work quickly on behalf of the city in times of crisis. Churches has numerous forms of informal groups that spread word quickly  prayer chains, worship groups, youth and adult bible studies, weekly dinners, and so on.

## 4.4  Implications

Interpreting a multi-layer network compares to peeling an onion - the removal of each layer reveals something new about the system. The goal of this study was to inform a rapid response plan to reach the widest range of individuals in Newburgh's most vulnerable areas with the lowest amount of time and resources possible. Using all of the included layers to inform this plan, it is clear that the clergy members present the most wide-spread flow of information due to the high worship

rates in the sector of study and the clear catchments, allowing targeting of certain regions within the sector if need be. Second and third order information flow makes the clergy such an appealing option because other power holders, such as business owners, have narrow client bases and can only target those individuals and their close connections.

## 5  Conclusions

A multi-layer network for the southern portion of the city of Newburgh, New York was developed as a tool for assisting in the disaster preparedness of the municipality. Ultimately, the analysis supports that the clergy of the five largest churches in the sector of study are best suited for distribution of information across this vulnerable portion of the city due to the wide range of citizens the reverends, pastors, and ministers can reach with second and third order dissemination of information. The greatest contributions to the discipline come from the creation of the geospatial layer, a developing area in Network Science, and displaying the strong similarities between calculated network structure and the formation of informal communities within a city. These similarities, which are correlated to the church catchments based on ethnic backgrounds and settlement, supports the claim that the physical layout of the city can be used as a tool to strengthen a disaster response plan by informing city officials where to best focus visual efforts of information dissemination.

## Acknowledgments

## References

[1] Lewis, T. (2011). Network Science: Theory and Applications. John Wiley & Sons, 7.

[2] Trump, Donald. The National Security Strategy of the United States of America. Executive Office Of The President Washington DC, 2017.

[3] Freeman, L. C. (1978). Centrality in Social Networks Conceptual Clarification. Social Networks, 1(3), 215-239.

[4] Erdős, P., & Rényi, A. (1960). On Random Graphs. Publ. Math. Inst. Hung. Acad. Sci, 5(1), 17-60.

[5] Barabási, A. L., & Albert, R. (1999). Emergence of Scaling in Random Networks. Science, 286(5439), 509-512.

[6] Watts, D. J., & Strogatz, S. H. (1998). Collective Dynamics of 'Small-World' Networks. Nature, 393(6684), 440.

# Nonlinear network dynamics under perturbations of the underlying graph

## Anca Rădulescu

State University of New York, New Paltz

radulesa@newpaltz.edu

(*joint work with* Sergio Verduzco-Flores, Ariel Pignatelli, Simone Evans)

### Abstract

Recent studies have been using graph theoretical approaches to model complex networks, and how hardwired circuitry relates to the ensemble's dynamic evolution in time. Understanding how configuration reflects on the coupled behavior in a system of dynamic nodes can be of great importance when investigating networks from the natural sciences. However, the effect of connectivity patterns on network dynamics is far from being fully understood.

We investigate the connections between edge configuration and dynamics in simple oriented networks with nonlinear nodes [1] which update in both discrete and continuous time. In discrete time, we use complex quadratic nodes [2,3]. We define extensions of the traditional Julia and Mandelbrot sets, and we study the changes in their topology and fractal behavior in response to changes in the network's adjacencies. In continuous time, we illustrate coupled Wilson-Cowan equations [4]. We use configuration dependent phase spaces and a probabilistic extension of bifurcation diagrams in the parameter space, to investigate the relationship between classes of system architectures and classes of their possible dynamics.

In both cases, we differentiate between the effects on dynamics of altering edge weights, density, and configuration. We show that increasing the number of connections between nodes is not equivalent to strengthening a few connections, and that certain dynamic aspects are robust to the network configuration when the edge density is fixed. Finally, we interpret some of our results in the context of brain networks, synaptic restructuring and neural dynamics in learning networks.

**References:**
[1] A. Rădulescu, Neural network spectral robustness under perturbations of the underlying graph, *Neural Computation* **28** (2015).
[2] A. Rădulescu, A. Pignatelli, Real and complex behavior for networks of coupled logistic maps, *Nonlinear Dynamics* **84** (2016) 2025-2042.
[3] A. Rădulescu, S. Evans, Networks of coupled quadratic nodes. Preprint: arXiv:1712.05953 (2018).
[4] A. Rădulescu, S. Verduzco-Flores, Nonlinear network dynamics under perturbations of the underlying graph, *Chaos* **25** (2015).

# LSTM Based Time-Series Link Prediction

Zeinab S Jalali, Chilukuri K. Mohan

4-206 CST, Department of Electrical Engineering and Computer Science
Syracuse University, Syracuse, NY 13244, USA
zsaghati, mohan@syr.edu

**ABSTRACT**
The problem of predicting future links is an important one relevant to many applications in bioinformatics, e-commerce, social networks and interactions of individuals over the internet. Traditional link prediction methods are often based on modeling and analyzing static snapshots of the network and fail to predict future links in real world networks that are dynamic. In this paper, we propose a new time series link prediction method using recurrent networks based on Long Short-Term Memory (LSTM) units. In the proposed approach, input data passes through the LSTM units, in a chain consisting a series of snapshots at time 1 to time T-1, with information about the presence of various links; predictions are then made about the existence of links at time T. Our experiments show better results compared to other methods (such as dynamic methods that use baseline methods like common neighborhood, Jaccard-coefficient, and Adamic-Adar) in terms of different accuracy measures.

**KEYWORDS**
Link Prediction, Time-Series, Long Short-Term Memory, Recurrent Neural Network.

## 1. INTRODUCTION

Link prediction is an important task in several research areas such as data mining, social networks, network science, and biology [1]. Many prior link prediction methods have been formulated on a static network where they predict hidden or future links given one single snapshot of a network. But real-world networks evolve dynamically, and link prediction needs to be modeled as a sequence of snapshots. The essential problem, referred to as time-series link prediction [2], is to determine which links will appear in the graph at time *T*, given a series of snapshots at time *1* through *T-1*. Some researchers have addressed time-series link prediction by generalizing static link prediction methods [3-7]. Most of these works do not consider periodic patterns and temporal trends of the communication intensities than can be extracted from the time-series information.

Long Short-Term Memory (LSTM) [8] units constituting recurrent neural networks [9] have recently been used successfully in other time-series prediction problems, such as natural scenes parsing [10] and tree-structured sentence representation learning [11]. We propose an LSTM based time-series link prediction model to predict the link occurrences using time series information. We identify a set of features that are needed for the relevant learning task and train an LSTM neural network on the prepared dataset. We compare our method with dynamic algorithms in terms of their prediction performance, using Mean Squared Error (MSE), Precision (Pr), Recall (Re), Accuracy (Ac), F-value (F1) and AUC on six datasets from different sources and domains.

Section 2 defines the problem, reviews background and related works on link prediction and RNNs and presents the datasets. Section 3 introduces a time-series link prediction algorithm based on LSTM. Section 4 presents the experimental study on predicting links in different datasets and discusses the results, and Section 5 summarizes the conclusions and discusses future directions.

## 2. PRELIMINARIES

We begin by describing the problem of time-series link prediction and discussing related work; this is followed by describing the data sets explored in our work.

### 2.1 PROBLEM STATEMENT

Let a *graph series* be a list of undirected graphs *(G1, G2, ..., GT-1)* corresponding to a list of symmetric adjacency matrices *(M1, M2, ..., MT-1)* [3]. Each $M_i$ is an $N \times N$ matrix with elements in {0,1}, where 0 denotes the non-existence of an edge and 1 denotes existence of an edge in *E(Gi)*. The time-series link prediction problem aims to predict occurrence of edges at time *T* given this graph series, specified as an $N \times N$ matrix $G_T$ with each element in {0,1}.

### 2.2 LINK PREDICTION

Time-series link prediction is an extended version of static link prediction where the aim is to infer which new interactions among its members are likely to occur in the near future given a snapshot of a graph [3]. Current approaches can be categorized into four groups: feature-based models, Bayesian probabilistic models, probabilistic relational models and linear algebraic methods [12]. We address link prediction as a binary classification problem, which can be considered as a feature-based link prediction method. One of the earliest feature-based link prediction methods was proposed in [12]. The learning method in their work uses different graph-based similarity measures to compute scores for probable edges. This work was extended in [13] where they added external data outside the scope of graph topology and changed the problem into a binary classification problem. This supervised classification approach became popular since then and has been used in different works in this area [14-15]. All these works are focused on static link prediction, our work uses the idea in these works to address time-series link prediction.

Time series link prediction was first studied in [2] where they introduced the problem and represented an initial effort toward building an integrated time-series link prediction model. Subsequently, temporal information was used with an extension of a local probabilistic model, in [3]. The link prediction problem was formulated as a periodic temporal link prediction in [4] and they tried to find the underlying periodic pattern. Cross-temporal link prediction was proposed in [5] where they predicted the links in different time frames. A new weighted time series link prediction was proposed in [6], in which they use both the similarity methods and the time-series patterns and [7] proposed a time series link prediction method which uses learning automata.

### 2.3 BASELINE METHODS

This section briefly introduces the commonly used link prediction methods for static graphs. Static graph link prediction algorithms take this single graph as input and produce the score matrix *S*. We change the score matrix to adjacency matrix by considering a threshold in a way that if the score for *S(v1, v2)* is greater than the threshold then *M(v1, v2)* will be 1, otherwise, it will be 0. The scores for the baseline methods are computed as follows:

1. Common Neighborhood(CN): The link occurrence score is set to the number of common neighbors between given nodes.
2. Jaccard-Index (JI): The link occurrence score is set to the Jaccard similarity between given nodes
3. Adamic-Ahar (AA): This index refines the simple counting of common neighbors by assigning the less-connected neighbors more weight.
4. Preferential Attachment (PA): The link occurrence score is set to be the product of the

degrees of the involved nodes.

5. Salton (SN): This feature is defined as the number of common neighbors between inputs divided by the square root of the multiplication of the degrees of the nodes.

## 2.4 RECURRENT NETWORKS

Recurrent neural networks (RNNs) are models that are not constrained by acyclic information flow constraints between neurons, and can hence represent relationships between data over time; they address sequential information, and have been applied successfully to several learning problems [10-11]. LSTM networks [17] improve on the performance of recurrent networks, especially for *deep* networks containing multiple layers, successfully tackling the vanishing gradient problem faced by networks when dealing with long data sequences. LSTMs allows for weight adjustments as well as truncation of the gradient by keeping the error flow constant through special units called gates, depending on which information is useful. LSTMs have been widely used in different areas and were able to achieve satisfactory results [18] compared to other methods in areas such as language modeling [19], translation [20], and acoustic modeling of speech [21]. Prior work with recurrent neural networks includes node classification [22], group activity recognition [23], feature learning [24], and prediction of structured sequences [25]. To the best of our knowledge, the RNN approach has not been previously applied to address link prediction using LSTM.

## 2.5 DATA DESCRIPTION

The datasets used in this work are networks from different sources and applications domains. These datasets are: BUP, CEG, UAL, INF, EML and NSC. BUP is a network of political blogs, CEG is biological network, UAL is an airport traffic network, INF is a network of face-to-face contacts in an exhibition, EML is a network of individuals who shared emails and NSC is a co-authorship network. For each dataset, we have 5 snapshots over time. Table 1 summarizes network properties: number of nodes ($|V|$), number of edges ($|E|$), average degree ($<k>$), average clustering coefficient (c), average shortest path length (ASPL), diameter (D) and heterogeneity (H).

### Table 1: Description of network data sets

| Name | $|V|$ | $|E|$ | $<k>$ | C | ASPL | D | H |
|------|------|------|------|------|------|---|------|
| BUP | 105 | 441 | 8.4 | 0.49 | 3.08 | 7 | 1.42 |
| VEG | 297 | 2148 | 14.46 | 0.29 | 2.46 | 5 | 1.80 |
| UAL | 332 | 2126 | 12.81 | 0.63 | 2.74 | 6 | 3.46 |
| INF | 410 | 2765 | 13.49 | 0.46 | 3.63 | 9 | 1.38 |
| EML | 1133 | 5451 | 9.62 | 0.22 | 3.61 | 8 | 1.94 |

## 3. LSTM BASED TIME-SERIES LINK PREDICTION

A list of undirected graphs $(G_1, G_2, ... G_{T-1})$ correspond to a list of symmetric adjacency matrices $(M_1, M_2, ..., M_{T-1})$ which are different snapshots of a network. For each snapshot, we compute several proximity features and convert the link prediction problem into a multi-variable learning problem, addressed using an LSTM network. We model the link prediction problem as a supervised classification task, where each data point corresponds to a pair of vertices in the input graph. We divide our work into four parts: 1) Feature selection; 2) Time-series feature based dataset construction; 3) Applying LSTM network; and 4) Finding evaluation measures.

## 3. 1 FEATURE SELECTION

In link prediction, features should represent some form of proximity between each pair of nodes that are considered as data points [13]. In this study, we use six well-known proximity features: Common Neighbors [26], Jaccard-coefficient [27], Adamic-Adar [28], preferential attachment [29], Katz [12] and Salton [30]. These features are the basis of most link-prediction methods.

## 3.2 TIME-SERIES FEATURE BASED DATASET CONSTRUCTION

After computing features for each data point on all snapshots, we need to change the dataset to a time-series dataset to be used as the input for LSTM network. We represent input features of data-point $(v_1, v_2)$ at step $T$ by $f_1(T), f_2(T), ... , f_n(T)$, and use $M_t(v_1,v_2)$ as the output for the supervised learning problem at step $T$. We build a classification model that can predict the unknown labels of each data-point $(v_i, v_j)$ by having a set of features from $G_1$ to $G_{T-1}$. Converting sequences to pairs of input and output sequences, we frame the supervised learning problem as predicting the link at time $T$, given the computed features at the prior time steps.

## 3.3 APPLYING LSTM NETWORKS

Before training an LSTM network on our input datasets, we normalize the features using the method proposed in [31] to equalize ranges of the features and make them have approximately the same effect in in the computation of similarity. Then we split the prepared dataset into training and testing sets. For each dataset, we use $k$-fold evaluation, where $k$ is equal to 5. We define the LSTM with 50 neurons in the first hidden layer and 1 neuron in the output layer for link prediction. Each model is trained for 50 training epochs with a batch size of 72. After training, we forecast for the entire test dataset.

## 3.4 EVALUATION MEASURES

As we converted the time-series link prediction problem to classifying each data-point into two groups (connected and disconnected), this problem can be considered as a classification problem. So, the classification measures including accuracy(ac), precision(pr), recall(re), F-measure(f1) and root mean square deviation (RMSE) are used. In addition, the AUC metric proposed in [33] is also evaluated. This metric tries to determine whether a random missing link has a higher score than a random non-existent link or not. If the AUC has a value greater than 0.5, it is better than the random link prediction algorithm; and the farther from 0.5, the algorithm is more accurate.

## 4 SIMULATIONS

In this section the behavior of proposed methods compared to dynamic methods on several networks are investigated. We compare our method with the dynamic method proposed in [2] which models the occurrence of each link as an independent time series and builds an autoregressive integrated moving average (ARIMA) model based on the link's past occurrence series and predicts the occurrence of the link in the future based on the intralink dependencies over time. It then combines the results from ARIMA model with the results from static link prediction algorithms to consider the interlink dependency patterns. To implement this method, we first predict the occurrence of each node using ARIMA model then we combine the results with the results from five baseline methods (CN, JI, AA, PA, SN) introduced in section 2.3 and random method (RM). We should mention that as baseline methods are static, we need to reduce the graph

series *(G1, ..., GT)* to a single graph *G(1-T)* with the corresponding adjacency matrix *G(1-T)* [13] and then compute the results.

## 4. EXPERIMENTS
Three sets of experiments, described below, were carried out.

### 4.1 EXPERIMENT I
In this experiment, we compare the proposed method with six dynamic methods using the classification measures including accuracy (AC), precision (PR), recall (RE), F- measure (F1), and Root mean square deviation (RMSE). Table 2 shows the results for different datasets; the best results are shown in boldface. For each dataset, the results shown are averages of 5 runs, with corresponding standard deviations.

*Table 2: Average comparison results for all datasets*

|  | CN | JI | AA | PA | SN | RM | LSTM |
|---|---|---|---|---|---|---|---|
| RMSE | 0.19±0.05 | 0.17±0.07 | 0.16±0.04 | 0.16±0.08 | 0.19±0.04 | 0.58±0.31 | **0.05±0.02** |
| AC | 0.92±0.09 | 0.97±0.02 | 0.84±0.24 | 0.97±0.03 | 0.52±0.38 | 0.59±0.01 | **0.99±0.01** |
| PR | 0.64±0.29 | 0.21±0.16 | 0.39±0.38 | 0.18±0.18 | 0.54±0.39 | 0.48±0.32 | **0.74±0.11** |
| RE | 0.28±0.11 | 0.36±0.09 | 0.35±0.13 | 0.28±0.09 | 0.28±0.08 | 0.08±0.04 | **0.95± 0.1** |
| F1 | 0.16±0.01 | 0.24±0.09 | 0.18±0.08 | 0.15±0.09 | 0.27±0.08 | 0.10±0.06 | **0.83±0.09** |

### 4.2 EXPERIMENT II
In this experiment, we compare the proposed method with dynamic methods using the AUC metric. The results are shown in Table 3.

*Table 3: Comparison of different methods using AUC*

|  | CN | JI | AA | PA | SN | RM | LSTM |
|---|---|---|---|---|---|---|---|
| BUP | 0.53± 0.04 | 0.68± 0.03 | 0.56± 0.05 | 0.71± 0.03 | 0.70± 0.06 | 0.58± 0.02 | **0.97± 0.02** |
| CEG | 0.57± 0.0 | 0.60± 0.01 | 0.57± 0.01 | 0.58± 0.01 | 0.72± 0.01 | 0.53± 0.01 | **0.95± 0.02** |
| INC | 0.52± 0.11 | 0.55± 0.02 | 0.52± 0.11 | 0.53± 0 | 0.62± 0.09 | 0.55± 0.03 | **0.97± 0.03** |
| UAL | 0.83± 0.02 | 0.84± 0.01 | 0.66± 0.03 | 0.77± 0.01 | 0.76± 0.03 | 0.63± 0.01 | **0.75± 0.13** |
| EML | 0.60± 0.04 | 0.66± 0.01 | 0.52± 0.04 | 0.63± 0.02 | **0.73± 0.04** | 0.59± 0.02 | 0.6± 0.11 |
| NSC | 0.56± 0.01 | 0.60± 0.01 | 0.52± 0.01 | 0.59± 0.01 | 0.57± 0.01 | 0.63± 0 | **0.89± 0.01** |

### 4.2.3 EXPERIMENT III
Since the input dataset is imbalanced and the data-points with label equal to 0 outnumber the data-points with label equal to 1, which may lead to bias over this label in the output dataset, we oversample those data points with label equal to 1 using the technique proposed in [32], to improve the performance of the classifier. The comparison of the proposed method before and after oversampling in terms of AUC and Precision is shown in Figure 1.
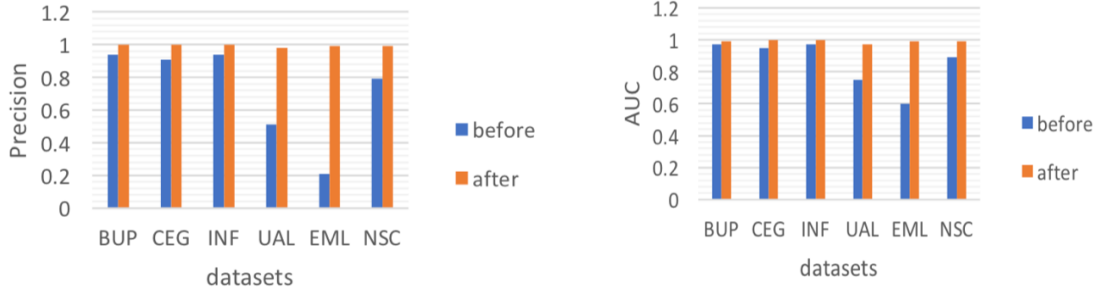
*Figure 1: Effect of oversampling the minority class*

**4.3 DISCUSSION**

The results reported from experiment I show that the proposed method works better than other methods in terms of accuracy, precision, recall and RMSE. In terms of AUC, the results of experiment II show that in most cases (five out of six datasets) the method outperforms other methods. To make our proposed method more accurate, we designed experiment III, in which we oversampled train data to prevent LSTM to be biased toward labeling data as 0. The results show that oversampling significantly improves the results. Figure 2 summarizes all the comparisons in term of precision (percentage of correctly detected links) and AUC after oversampling. The results reported here demonstrate that the proposed time series link prediction method after oversampling is able to achieve an average better accuracy and precision than other methods.



*Figure 2: Comparison of all methods using precision and AUC*

**5. CONCLUSIONS AND FUTURE WORK**

We have presented a new time series link prediction method which uses LSTM to predict the existence or non-existence of each link at future instants by using past data. First, we changed the structure of the graph into a sequence-based structure to be used as the input of recurrent neural network. We also reframed the problem as a supervised learning problem, extracted features and used an LSTM recurrent neural network to solve the link prediction problem. The proposed method was evaluated using AUC metric, accuracy, precision, recall and RMSE on six link prediction networks. Oversampling the minority class was used to address the class imbalance problem.

The experimental results reported here showed that the proposed algorithm is superior to other dynamic algorithms. The better result may be attributed to the learning capability of an LSTM recurrent network to evolve through time; different similarity metrics can be used at different times for future predictions.

In future work, we would like to consider a number of large datasets from different domains to address the link prediction problem. We also plan to consider the problem for streaming graphs.

## 5. REFERENCES

1. Zhao, Peixiang, Charu Aggarwal, and Gewen He. "Link prediction in graph streams." *Data Engineering (ICDE), 2016 IEEE 32nd International Conference on*. IEEE, 2016.
2. Bilgic, Mustafa, Galileo Mark Namata, and Lise Getoor. "Combining collective classification and link prediction." *Data Mining Workshops, 2007. ICDM Workshops 2007. Seventh IEEE International Conference on*. IEEE, 2007.
3. Tylenda, Tomasz, Ralitsa Angelova, and Srikanta Bedathur. "Towards time-aware link prediction in evolving social networks." *Proceedings of the 3rd workshop on social network mining and analysis*. ACM, 2009.
4. Dunlavy, Daniel M., Tamara G. Kolda, and Evrim Acar. "Temporal link prediction using matrix and tensor factorizations." *ACM Transactions on Knowledge Discovery from Data (TKDD)* 5.2 (2011): 10.
5. Oyama, Satoshi, Kohei Hayashi, and Hisashi Kashima. "Cross-temporal link prediction." *Data Mining (ICDM), 2011 IEEE 11th International Conference on*. IEEE, 2011.
6. Huang, Shiping, et al. "Link prediction based on time-varied weight in co-authorship network." *Computer Supported Cooperative Work in Design (CSCWD), Proceedings of the 2014 IEEE 18th International Conference on*. IEEE, 2014.
7. Moradabadi, Behnaz, and Mohammad Reza Meybodi. "A novel time series link prediction method: Learning automata approach." *Physica A: Statistical Mechanics and its Applications* 482 (2017): 422-432.
8. Gers, Felix A., Jürgen Schmidhuber, and Fred Cummins. "Learning to forget: Continual prediction with LSTM." (1999): 850-855.
9. Medsker, L. R., and L. C. Jain. "Recurrent neural networks." *Design and Applications* 5 (2001).
10. Socher, Richard, et al. "Parsing natural scenes and natural language with recursive neural networks." *Proceedings of the 28th international conference on machine learning (ICML-11)*. 2011.
11. Tai, Kai Sheng, Richard Socher, and Christopher D. Manning. "Improved semantic representations from tree-structured long short-term memory networks." *arXiv preprint arXiv:1503.00075*(2015).
12. Liben-Nowell, David, and Jon Kleinberg. "The link-prediction problem for social networks." *Journal of the Association for Information Science and Technology* 58.7 (2007): 1019-1031.
13. Al Hasan, Mohammad, et al. "Link prediction using supervised learning." *SDM06: workshop on link analysis, counter-terrorism and security*. 2006.
14. Doppa, Janardhan Rao, et al. "Chance-constrained programs for link prediction." *NIPS Workshop on Analyzing Networks and Learning with Graphs*. 2009.

15. Wang, Chao, Venu Satuluri, and Srinivasan Parthasarathy. "Local probabilistic models for link prediction." *Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on*. IEEE, 2007.
16. Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." *Neural computation* 9.8 (1997): 1735-1780.
17. Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to sequence learning with neural networks." *Advances in neural information processing systems*. 2014.
18. Graves, Alex, et al. "A novel connectionist system for unconstrained handwriting recognition." *IEEE transactions on pattern analysis and machine intelligence* 31.5 (2009): 855-868.
19. Graves, Alex. "Generating sequences with recurrent neural networks." *arXiv preprint arXiv:1308.0850* (2013).
20. Sak, Haşim, Andrew Senior, and Françoise Beaufays. "Long short-term memory recurrent neural network architectures for large scale acoustic modeling." *Fifteenth Annual Conference of the International Speech Communication Association*. 2014.
21. Xu, Qiongkai, et al. "Collective Vertex Classification Using Recursive Neural Network." *arXiv preprint arXiv:1701.06751*(2017).
22. Deng, Zhiwei, et al. "Structure inference machines: Recurrent neural networks for analyzing relations in group activity recognition." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016.

# Ecological Network Structure and Data Resolution

Kayla R. Sale-Hale[1] and Surendra Hazarie[2]

[1] Ecology & Evolutionary Biology, University of Arizona, kaylasale@email.arizona.edu
[2] Physics, University of Rochester, shazarie@u.rochester.edu

The stability of an ecosystem is influenced by the diversity of the structure of its component networks. We examine two plant-pollinator systems from each of 7 locations in the Canary Islands, resulting in 14 weighted bipartite networks (Figure 1). The unique combination of island pairs and edge weights provide many ways to aggregate the data, each with its own implications. We measure the modularity and Shannon-Jensen divergence of the networks under these different levels of resolution and find that the relationships between the networks change in nontrivial ways, such as the devaluing of physical closeness. This suggests that data resolution must be considered when making any conclusions about ecological stability or conservation.
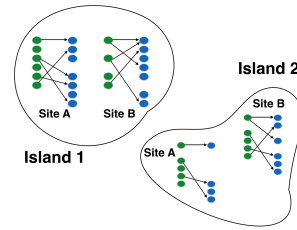
**Canary Island Plant-Pollinator Networks**



Figure 1: Schematic of the island plant-pollinator networks.
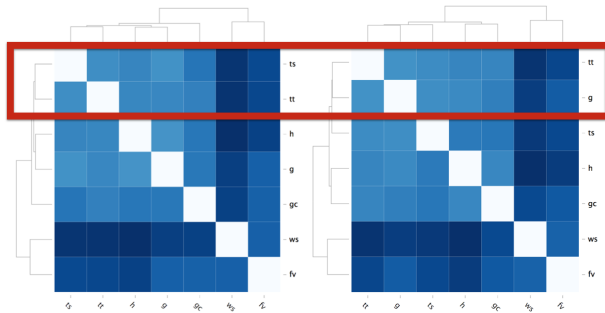
**Shannon-Jensen Divergence**



Figure 2: SJD for the 7 site pairs, at the species level (left) and aggregated into genus (right). Sites "tt" and "ts" both correspond to locations on the island of Tenerife and are grouped together when species level precision is retained, but not when it is discarded.

We use the generalized Louvain algorithm to measure the community structure of the islands across different resolutions. Since these communities indicate functional groups, we can observe how the interpretation of a species' role in the ecosystem changes with resolution. We measure the Shannon-Jensen divergence of the networks in our system, performing hierarchical clustering to group them based on similar information content. We find that the algorithm always clusters same-island sites together, but clustering beyond that is determined by the aggregation method. Usually one would expect islands that are physically close to each other to be the most related, but this is not the case under certain aggregations (Figure 2), further indicating the sensitivity of ecological stability to resolution.

# A Novel Approach to Event-driven Simulation of Spiking Neural Network

Fatemehossadat Miri[1] and J. David Schaffer[2]

[1]Systems Science and Industrial Engineering Department
Watson School of Engineering and Applied Science
Binghamton University
fmiri1@binghamton.edu
[2]Institute for Justice and Well-being
Binghamton University
dschaffe@binghamton.edu

The objective of this research is to explore the computational capabilities of Spiking Neural Networks (SNNs) for machine intelligence as opposed to advancing biological understanding. Artificial Neural Networks (ANNs) are adaptive and powerful tools that process data similar to human brain. SNNs are the third generation of ANNs, which take the goal of making a machine as intelligent as human brain one step closer to reality. SNNs have the same computational capabilities of $2^{nd}$ generation ANNs while require exponentially less neurons[1]. Neuron models and learning models defined for SNNs are more biologically accurate than previous neural network models. SNNs also offer hybrid analog/digital hardware while require low energy as it could re-use its neurons and synapses for detecting different patterns.

To achieve the objective of this research, it is essential to develop a suitable tool and conduct experiments. There are two methods to simulate SNNs and perform experiments, first, time driven and second, event driven simulations. In this paper, a novel approach to perform event-based simulation of SNNs is presented. None of the available SNN simulators offer ALL the required features including Synaptogenesis and Neurogenesis capabilities, event driven logic, parameter flexibility and certain neuron and synapse models that abstract some of the biological details. Hence, an event_driven simulator, Synapse 2.0, was developed in Java from scratch and it is equipped with all the named features.

The novel logic of Synapse 2.0 employs the most two recent Postsynaptic Potential (PSP) waves to predict future events and exact time of spike events. Additionally, information from the two PSPs along with the state of neurons at each event are used to remove or make corrections to previously predicted events. Furthermore, current version of simulator, Synapse 2.0, offers a unique combination of flexibilities for neuron and synapse specifications. It allows for arbitrary connections, excitatory and inhibitory synapses, synaptic parameters and learning paradigms for each synapse and arbitrary neuron parameters for each neuron.

A set of experiments was run on two simulators, Synapse 2.0 and SNNS, which is a time-driven simulator. Results from both simulators confirm that the combination of this new events-driven simulation logic with the novel combination of all the flexibilities for parameters was successful and Synapse 2.0 behaves as expected. Next version of the simulator will be equipped with two processes, Synaptogenesis and Neurogenesis, which will allow SNN to self-evolve its structure during training phase.

## References:

1. Maass, Wolfgang. "Networks of spiking neurons: the third generation of neural network models." *Neural networks* 10, no. 9 (1997): 1659-1671.

# Agent-based model parameter estimation and variable reduction using metaBUS: An application to a collective leadership model

Neil G. MacLaren[1], Yiding Cao[2], Ankita Kulkarni[1], Francis J. Yammarino[1], Michael D. Mumford[3], Shelley D. Dionne[1], Hiroki Sayama[2], Shane Connelly[3], Tyler J. Mulhearn[3], Robert Martin[3], Erin Todd[3], and Frank A. Bosco[4]

[1]School of Management, Binghamton University, State University of New York
[2]Department of Systems Science and Industrial Engineering, Binghamton University, State University of New York
[3]Department of Psychology, University of Oklahoma
[4]Department of Management, Virginia Commonwealth University

## Abstract

This study used a new meta-analysis platform called metaBUS to generate empirically derived parameter estimates for the relationships between variables being considered for inclusion in an agent-based model (ABM). We investigated a variety of methods for gathering large amounts of correlation data from metaBUS and compared two common dimension reduction algorithms, principal component analysis (PCA) and community detection, in order to use the information in metaBUS to empirically reduce the variable set. We used weighted degree centrality to identify the most important node in each community for use in future model development. To our knowledge, this study is the first to use metaBUS to assist in generating parameter values for an ABM and the first study to use a network analysis tools to empirically reduce a variable set in the organizational sciences.

## 1 Introduction

Social science researchers are increasingly turning to simulations and computational modeling to understand complex systems. For example, researchers have recently used agent-based models to explore the dynamics of knowledge transfer in organizations [1] and team shared mental model convergence [2]. Unfortunately, many social science computational models use arbitrary parameter values: although these models can help improve our understanding of the topic area, the field needs to move beyond arbitrary parameter values in order to support more rigorous, testable models.

One way to move beyond arbitrary parameter values for computational models is to simply estimate them from the literature. However, the social science literature is vast and plagued with nomenclature difficulties: a variety of operationalizations exist for the same construct under different names and the same names have been used for operationalizations of different constructs. Furthermore, current methods of estimating initial parameter values from the literature require manual search and analysis; even aided by electronic search tools and computation aids, this kind

of analysis is laborious and error-prone [4]. Despite the difficulties the social science literature contains a large amount of potentially relevant information that could support model development—an efficient, accurate way of developing initial parameter estimates from this literature is needed.

A new database called metaBUS (Meta-Analysis Omnibus, metabus.org) may be able to partially fulfill this requirement. MetaBUS provides access to 1.04 million empirical findings—each represented by variable means, standard deviations, and other summary information—in 23 journals in applied psychology, organizational behavior, and human resources. Published meta-analyses, simulations, and theoretical papers are excluded. MetaBUS was designed to support user-defined meta-analysis, but its ability to search and meta-analytically summarize data from correlation tables in large numbers of published articles has potential for other uses. The purpose of this study was to test whether or not metaBUS could provide initial literature-based parameter estimates for an agent-based model.

## 2  Method

This study uses meta-analysis tools supported by the metaBUS database of empirical findings to search the literature for correlations between each possible pair of variables specified from a conceptual model. We then compare the ability of two procedures, principal component analysis (PCA) and community detection, to use this information to support variable reduction decisions.

### 2.1  MetaBUS Procedure

In order to work with metaBUS the researcher first needs an idea of under what variable names the model constructs might have been measured in the literature—these are then translated into search terms that metaBUS can accept. In this study, the team first developed a conceptual model based on a qualitative literature review [3]. Conceptual variables were then specified with behaviorally anchored rating scales and organized into a code book. The team used the code book definitions and scales to generate a list of variable synonyms and related constructs and the list was translated into metaBUS search term sets using published recommendations [4]. These search term sets were then used pairwise to generate meta-analytic estimations of the correlation between two potential model variables.

Stated more fully, we took the following steps to use metaBUS to generate estimates of the correlation between potential model variables:

1. Select search methodology: MetaBUS has a taxonomic tree of domain constructs and allows the user to search its database using metaBUS defined construct categories. Searches using taxonomic categories will return variables coded in metaBUS as members of that category. Simple text string searches are also permitted, allowing the user to search published variable names directly rather than metaBUS-defined categories of variables.

2. Develop search term sets: Because we were searching a large number of variable pairs, inclusion and exclusion decisions based on topic and level of analysis match were done as much as possible using metaBUS's search tools. We relied on the detailed definitions of constructs both in the original review [3] and in the variable code book derived from it to make decisions about which synonyms to use for a given variable, at what level of analysis it should have been measured, and what potentially related concepts needed to be excluded. Once we

had built an initial list of search term sets, each set was tested against the full list of search term sets to confirm that each set was recovering a reasonable set of published findings.

3. Search: MetaBUS supports pairwise search of its database and returns a list of studies found and meta-analytic summary statistics. At this time, general users are only able to meta-analyze one pair of variables at a time. Therefore, one of the authors (FAB), a metaBUS project lead, compared each search term set—representing a coded variable in the larger study—against itself and every other search term set in the full list directly in the database.

Depending on the topic area, one search method or the other may be preferable, and an additional goal of this study was to investigate the differences between the two search methods. Because of this goal, we conducted the searches before selecting the search methodology.

## 2.2 Principal Component Analysis

MetaBUS searches can provide information about the correlations between constructs in the conceptual model, but the number of potential variables derived from a conceptual model can still be large. PCA and community detection are two methods that are commonly used to reduce dimensions in a system with a large number of variables. PCA is guaranteed to produce uncorrelated variables, can reveal important clustering behavior that might have been hidden in the untransformed data, and produces output with minimal information loss. A drawback of PCA methods, however, is that the individual PCs can be difficult to interpret. We prepared the metaBUS data for PCA by considering each variable to represent both a row and a column in the data matrix; the values in the matrix were the correlation estimates from metaBUS. Each variable was compared against itself in the database, and the result of those searches, appearing on the diagonal, was a correlation estimate of less than 1.0.

## 2.3 Community Detection

Empirical networks typically have a structure in which nodes are densely connected in some regions, called communities, but sparsely connected between those regions. To use community detection for variable reduction, each variable is seen as a node in a network with an undirected edge, the weight of which is proportional to the correlation coefficient, between each variable. Although the output communities may still be somewhat correlated, community detection makes it clear which original variables have been assigned to which new variable set—an important advantage in interpreting the final model. We used the Louvain algorithm [5] with a resolution factor of 0.9 in Gephi 0.9.2 to identify communities.

We used weighted degree centrality, calculated over the whole graph, to select the most important variable within each community. We decided to use weighted degree centrality for several reasons. First, the graph relied on weighted edges; the Louvain algorithm already uses edge weights in its calculations, so a centrality measure that also relied on edge weights seemed appropriate. Second, the graph contained no information on inter-node distance. Third, because the data was derived from metaBUS searches and not directly from nature, the results are likely biased both by the adequacy of our search terms and by what has previously been empirically studied and indexed by metaBUS: centrality measures that relied only on the existence of a connection between two nodes therefore seemed inappropriate. Finally, weighted degree centrality also uniquely identified the most important node in each community without having an extremely skewed distribution.

## 3  Example Study

Previously, several of the authors and their colleagues developed a conceptual model as an attempt to describe the influences and outcomes of leadership at the collective level of analysis [3]. This model was based on an extensive literature review as well as a series of previous laboratory and historiometric studies [6, 7]. As part of research supported by U.S. Army Research Institute grant #W911NF-17-1-0221, the authors are in the process of testing this model through a multi-method study design that involves both laboratory experiments and computational modeling. This work is ongoing, and a total of 60 experimental sessions at two universities are planned. Each session will be analyzed using behaviorally anchored rating scales quantifying 132 variables derived from the conceptual model; data from these experiments will be described using a system of dynamical equations. The team will separately implement the conceptual model in an agent-based simulation. Results from each method will be compared in the final analysis. Future directions for the research project include developing training methodologies for planning and leadership teams based on findings from the current effort.

We used the metaBUS procedures described above to assign search strings to each of the 132 variables identified from the conceptual model. All search sets of both nodes and text strings were tested in the metaBUS Shiny application against the full set of node searches identified for this project in order to refine search terms. Many of the 132 variables deal with specific hypothesized aspects of the collective planning process and have not previously been well studied. Tests of reasonable strings returned no meaningful results, so the final data set considered 93 of the 132 variables (for a total of 4,371 possible meta-analytic correlation estimates). To develop the data sets, the lists of search strings were used to conduct pairwise searches between each variable, including each variable with itself, as represented in metaBUS by its search string. We used the meta-analytic correlation estimate between two variables when one or more studies were found matching the search results; thus, several variables in these results are only represented by one published correlation estimate.

### 3.1  Search Method

All variables were translated into metaBUS search term sets using both the taxonomic category and text search string methods, and each method was used by itself and in combination with the other in all searches. The text string searches found information on 15% of the 4,371 possible correlations; concept taxonomy searches found information on 27% of possible correlations; and the combined text string and concept taxonomy searches found information on 40% of possible correlations. Using the set of searches where text string searches and concept taxonomy searches both returned results, the text string searches were found to be moderately correlated with the taxonomy searches ($r = 0.47$) and both sets of searches were found to be largely similar in their ability to predict the combined search results, with the larger number of taxonomy-based search results having a predictably greater impact on the combined result (text search $R^2 = 0.51$; taxonomy search $R^2 = 0.63$). For this study, moderator analyses would not be used given the large number of searches conducted, so minimizing false positives was prioritized over sample size. Furthermore, we felt it was important to uniquely specify each variable and rely on empirical variable reduction to reduce the set rather than use theory to reduce the set before the searches began. Because each study variable could be uniquely specified using text string searches but not taxonomy searches, and

Table 1: Summary of the variable communities, including the most central variable for each community with several centrality measures. Variable definitions are discussed in [3] and the study code book is available upon request. (O) = outcome variable; (L) = leader variable; (C) = collective variable.

|   | % Nodes | Central Variable | Degree | Weighted Degree | Betweenness |
|---|---------|------------------|--------|-----------------|-------------|
| A | 16.9 | Performance Outcomes (O) | 68 | 14.7 | 0.292 |
| B | 11.3 | Performance Monitoring (C) | 27 | 6.7 | 0.010 |
| C | 11.3 | Game Progress (O) | 48 | 11.9 | 0.063 |
| D | 11.3 | Problem-Solving Capacities (O) | 36 | 12.5 | 0.023 |
| E | 18.3 | Leader-Member Exchange (L) | 47 | 14.6 | 0.047 |
| F | 12.7 | Trust (C) | 42 | 13.3 | 0.031 |
| G | 18.3 | Procedural Justice (L) | 31 | 10.5 | 0.014 |

because in test searches the text string search method seemed to provide the ability to conceptually better match the measured variables from the code book, the text string method was chosen and used in all future analyses.

## 3.2 PCA

The text search results were considered too sparse for meaningful PCA. No variable had less than 40% missing data, and only 17 variables had less than 80% missing data; no information at all was found for 23 variables. The PCA method was not considered further in this study.

## 3.3 Community Detection

We used both Gephi 0.9.2 and the python-louvain package to determine community structure. Although results varied depending on function settings, a common and interpretable result in both software applications was seven communities. These communities seemed to have some stable properties, but community number and membership remains to be fully explored in future analysis. The results of a seven-community analysis in Gephi are presented here.

Table 1 displays the results of the community detection procedure. Communities E and G include the highest percent of original variables (18.3% each). A variable in Community A, Performance Outcomes, has the highest degree, weighted degree, and betweenness centrality for the graph as a whole.
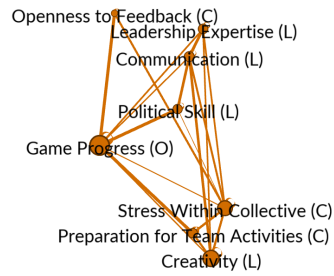
Figure 1 (A) - (G) shows the network structure of each community displayed as its own subgraph. In each graph the size of the node is proportional to the weighted degree centrality of that node; the thickness of each edge is proportional to the correlation between the two nodes. Figure 1 (F) shows the subgraph of the six nodes with the highest weighted degree centrality. One variable from each of six of the communities made up the top six most central nodes in the whole graph; the seventh community's highest centrality node, Performance Monitoring from Community B, had the 19th highest weighted degree centrality.
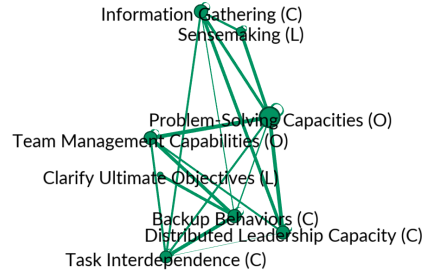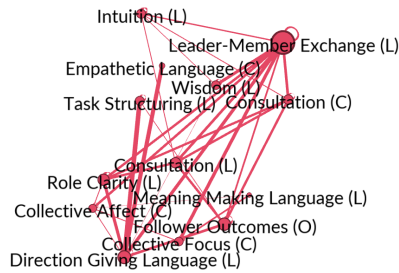
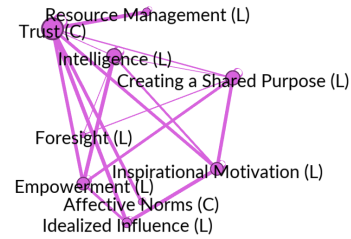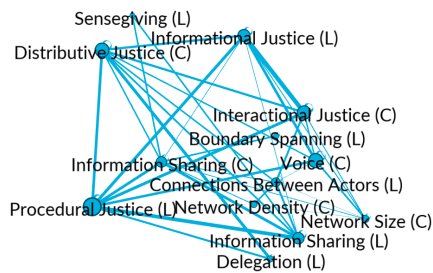(a) Community A      (b) Community B

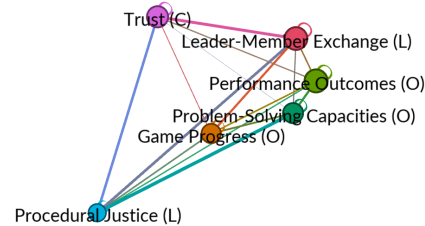(c) Community C      (d) Community D

(e) Community E      (f) Community F

(g) Community G      (h) Central Nodes

Figure 1: Variable community subgraphs (A) – (G) and the most central network nodes by weighted degree centrality (H).

## 4 Discussion

We have introduced a method for initial variable reduction and parameter estimation that investigated the empirical support for correlation between 4,371 pairs of variables. Although further analysis is indicated—and planned—an initial set of clearly defined variables is available for modeling and presented in Table 1. We can define sets of other variables that appear to be closely related to these focal variables and have information about how each focal variable is related to the other. Although these results are preliminary, they are encouraging.

This study has several important limitations. First, for an edge to appear in the graph researchers must have measured the pair of variables in a way that could be summarized in a correlation table, the measurements must have been published in a journal metaBUS indexes, and our search term sets must have found the published correlation in metaBUS. It is not clear to what extent these concerns are biasing our results. We hope to partly address this limitation by comparing the results of the metaBUS analysis with data from our experiments.

Second, the Louvain algorithm is stochastic. Although seven communities is a common result, six communities and five communities are also common results with equivalent modularity index values. While weighted degree centrality is a stable measurement regardless of community delineation, in this data set there is no obvious set of nodes that are the most important nodes in the graph without community detection. Thus, the instability of both the number of communities and the membership of those communities presents a problem in interpreting the results and is another potential avenue for future research.

Despite these limitations, our results appear to both align with previously published analysis and suggest new avenues for future research. For example, Performance Outcomes, an outcome variable in the experiments, has the highest magnitude in all centrality measures for Community A and the graph as a whole. Work unit performance is the most important outcome variable in organizational studies, and it is therefore unsurprising that its betweenness centrality should be an order of magnitude higher than every other node in the graph: it is simply the most commonly studied. However, when edge weights are considered, the variables with the strongest in-community relationships with Performance Outcomes are Problem Definition, defined by behaviors such as making explicit task assignments when the situation is unclear, and Adaptive Performance, the extent to which a collective changes operating procedures when a new problem arises. Within-community relationships such as these are suggestive and, if stable in future robustness analysis, worthy of further investigation.

Another example of both agreement and opportunity is found in Community E. Given the long tradition in leadership studies of differentiated between leader task and relationship behaviors, it is notable that Leader-Member Exchange (LMX), usually considered to reflect the extent to which a leader forms individual relationships with their subordinates, is a close second to Performance Outcomes in terms of weighted degree centrality. In this analysis LMX is more central to the network as a whole than our Follower Outcomes variable, which represents follower trust, satisfaction, and loyalty towards other members of the collective.

## 5 Conclusion

Returning to our research question, we find that although metaBUS was designed to support standard meta-analysis, we can use metaBUS for variable reduction and parameter estimation for a

social science computational model. MetaBUS can return results for some variable relationships, with the percentage of possible results returned varying by search method. Although the high proportion of missing data may preclude the use of PCA for variable reduction, community detection appears to give meaningful results. In this data set we were able to identify variable communities and important variables within those communities despite the large amount of missing data.

By using the community detection process described above we were able to identify seven communities of variables and one central variable within each community. By selecting from among codebook variables, rather than deriving a set of abstract new variables, we are able to retain comparability with the experiment portion of the overall study and use metaBUS to draw meta-analytic correlation estimates. These results will support the development of an agent-based modeling by providing the most important variables to model, their definitions, and a value estimate for the relationships between the variables.

## Acknowledgments

## References

[1] Sheen S Levine and Michael J Prietula. How knowledge transfer impacts performance: A multilevel model of benefits and liabilities. *Organization Science*, 23(6):1748–1766, 2012.

[2] Shelley D Dionne, Hiroki Sayama, Chanyu Hao, and Benjamin James Bush. The role of leadership in shared mental model convergence and team performance improvement: An agent-based computational model. *The Leadership Quarterly*, 21(6):1035–1049, 2010.

[3] Carter Gibson, Tristan McIntosh, Tyler Mulhearn, Shane Connelly, Eric Anthony Day, Francis Yammarino, and Michael D Mumford. Leadership/followership for long-duration exploration missions final report. Technical Report TM-2015-218567, NASA, Houston, TX, 2015.

[4] Frank A Bosco, Krista L Uggerslev, and Piers Steel. MetaBUS as a vehicle for facilitating meta-analysis. *Human Resource Management Review*, 27(1):237–254, 2017.

[5] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10):P10008, 2008.

[6] Francis J. Yammarino, Michael D. Mumford, William B. Vessey, Tamara L. Friedrich, Gregory A. Ruark, and Jason M. Brunner. Collective leadership measurement for the U.S. Army. Study Report 2014-01 Contract No. W91WAW-09-C-0090, U.S. Army Research Institute for the Behavioral and Social Sciences, Fort Belvoir, VA, 2014.

[7] Tamara L. Friedrich, William B. Vessey, Matthew J. Schuelke, Michael D. Mumford, Francis J. Yammarino, and Gregory A. Ruark. Collective leadership and George C. Marshall: A historiometric analysis of career events. *The Leadership Quarterly*, 25:449–467, 2014.

# RESOURCE-PRICE-DEMAND DYNAMICS MODEL FOR CLOUD MANUFACTURING

Summyia Qamar[a], Muhammad Haris Aziz[b], Mohammad T. Khasawneh[b], Chanchal Saha[c]

[a] University of Engineering and Technology, Taxila, Pakistan
[b] Department of System Science and Industrial Engineering, Binghamton University, New York
[c] IBM Corporation Poughkeepsie, New York

## Abstract

This paper captures the dynamics of resource-price-demand in Cloud Manufacturing (CMfg) with the application of Complexity Science. To date, research on pricing in CMfg is mainly focused on the Cloud Service Providers (CSP), ignoring the profit margin of Resource Service Providers (RSP) or manufacturers, which can be achieved through their price adjustment based on Resource Demander (RD) demand. To support this argument, quantitative modelling approach is applied to develop an integrated model of customer demand and manufacturer's resources, in the context of changes brought by the prices. Simulation study of proposed model is performed under different scenarios using real world parameters. The behaviour of system under continuous demand and unlimited resources, as characterized in CMfg system, shows a convergent behaviour. Thus, with continuous evolution of system's components (demand and available resources), the prices of manufacturers should also be adjusted accordingly. Whereas, in restricted demand and resource scenario, the system converges to equilibrium, whereby allowing manufacturers to generate revenue by dynamic pricing of their resources according to converging demand and maximum available resources.

## 1. Introduction

Complexity Science has gained much attention due to its broad spectrum of applications in biology, engineering, environmental and neuro sciences and networking. As a whole, it encompasses complex mathematical and computational modelling approaches to deal with complicated systems which constitute of inter-related components exhibiting independent properties [1]. Since, overall properties of a system cannot be predicted deterministically by aggregating or averaging individual components behaviour, therefore, this leads to emergent behaviour of a system [2].

One such complex system prevailing in temporary manufacturing is Cloud Manufacturing (CMfg), which is characterised by ubiquitous, on-demand access to virtually distributed manufacturing resources [3]. It enables consumers on-demand purchase of manufacturing service from service market. Consumer participation in service transaction is closely linked with the price of service. Consequently, an appropriate price would make a service provider more competitive and assist in gaining high revenue. Therefore, the study of relation between customer demand, manufacturer's resources and pricing requires adoption of complexity science approaches. However, so far literature on pricing in CMfg provides Cloud Service Provider (CSP) based pricing mechanism in which the prices of Resource Service Providers (RSP) of manufacturer are considered fixed. To

abridge this gap, a quantitative modelling approach is adopted to formulate a mathematical model of demand-resource-price dynamics.

## 2. Literature Review

Complex systems science and modelling has gained public attention in 1990's due to its wide spectrum of applicability in topics such as social networks, evolution and environmental changes [4]. However, the core concepts of complex system science, that goes over almost every area of its applications, are emergence and self-organisation [5]. To study this emergent behaviour different studies have been conducted using operations research modelling and simulation approaches, including discrete event modelling, agent-based modelling and system dynamics [4, 6].

The evolutionary nature of complex systems helps them in growth and improvement by learning through their environment [7]. This evolutionary phenomenon has been observed in manufacturing industry which is transforming from production-oriented to service-oriented manufacturing [8]. One such system prevailing in contemporary industrial environment is CMfg [3, 9]. Typically, CMfg is comprised of three types of users in service transaction which includes RSP or manufacturer who owns resources, CSP or third party who is responsible for delivering services and service demander or customer who purchases the services. They are connected with each other through knowledge sharing mechanism [10]. The price of services offered on the cloud are initially provided by different manufacturers. CSP then develops an equilibrium between customer's willingness to pay and the manufacturer's profit margin. Therefore, the profit of manufacturer is primarily dependent on the prices offered by him on the cloud.

Literature on pricing highlights the significance of Dynamic Pricing (DP) strategy for price-sensitive demand. McAfee and Te Velde [11] reported around $500million increase in revenue of American airline industry by the application of DP. DP in cloud environment has been widely studied in literature with different strategies applied to maximize revenue [12-15]. Xu and Li [16] proposed a DP framework for revenue management of Amazon. Ranaldo and Zimeo [17] developed a utility model for cloud service transaction to optimize the utility of CSP in parallel with pricing on cloud. GUO, PENG [18] studies CMfg behaviour under cost uncertainty and considered risk factor of participants. However, their model was based on profit equilibrium under CSP-based strategy. Li and Ma [19] proposed a price Stackelberg game model under probabilistic selling conditions. Cheng, Tao [20] proposed a CMfg resource service utility equilibrium under centralized decision system based on time, revenue and reliability factors. Lin, Yang [21] extended the literature on RS providence by proposing multi centric management and resource allocation in CMfg. Argoneto and Renna [22] worked on capacity sharing with the cost parameters related to logistics and manufacturing. Liu, Zhang [23] proposed a RS sharing model in CMfg to increase utilization of resources, satisfaction rate and utility of enterprise as compared to traditional Network Manufacturing model.

Apart from the contribution of above-mentioned literature in study system dynamics of CMfg, consideration of fixed pricing of RSP impedes their revenue generation which can be achieved through demand-based pricing of their services. Therefore, to support this argument, quantitative modelling approach is applied to study the emergent behaviour of demand and availability of resources under variant behaviour of prices.

## 3. Quantitative Modelling of Complex Cloud Manufacturing System

To represent the dynamics of pricing on the manufacturer's resources and customer demand, a hypothesized CMfg environment is formulated where the manufacturer's resource availability, customer demand and price interaction is studied.

## 3.1.Notations

D=Customer Demand
$\overline{D}$= Demand Constant
R=Available Resources
$\overline{R}$=Resource Constant
P=Price
$\overline{P}$=Price Constant
t=Time period, t=1,2,…T
$\alpha$=Price to demand conversion ratio=$\overline{D}/\overline{P}$
$\beta$=Demand to price conversion ratio=$\overline{P}/\overline{D}$
$\gamma$=Resource to price conversion ratio=$\overline{P}/\overline{R}$
$\delta$=Demand to resource conversion ratio=$\overline{R}/\overline{D}$

## 3.2.Assumptions and Mathematical Model

It is assumed that there is only one manufacturer's monopoly. Since, customer demand is price dependent, therefore, to analyse the behaviour of system with inter-related demand and price, Equation (1) is formulated.

$$D_t = D_{t-1} + \left(1 - (P_{t-1}/\overline{P})*\alpha\right) \tag{1}$$

Where To formulate pricing equation, two factors are considered; availability of resources and customer demand. Equation (2) shows that if resource availability increases, price decreases and if demand increases price increases to control the demand.

$$P_t = P_{t-1} - \left(\left(1 - \left(D_{t-1}/\overline{D}\right)\right)*\beta\right) + \left(\left(1 - \left(R_{t-1}/\overline{R}\right)\right)*\gamma\right) \tag{2}$$

Since, availability of resources depends on demand, therefore, Equation (3) is formulated with shows that with increase in demand, number of resources availability decreases.

$$R_t = R_{t-1} + \left(\left(1 - \left(D_{t-1}/\overline{D}\right)\right)*\delta\right) \tag{3}$$

Thus, a set of linear difference equations with integrated demand, resources and pricing of resource service is formulated.

To study the overall behaviour of system under these equations, simulation is performed.

## 3.3.Simulation Results

Simulation is performed based on real scenarios where each scenario is considered as a shop with different set of parameters at initial stage. In this way, each scenario is simulated under 3 cases based on the initial values of parameters.

### 3.3.1. Scenario I: Proposed Model

Proposed system is simulated using $\overline{D}$=15; $\overline{P}$=10; $\overline{R}$=7. The results are shown in figure 1.

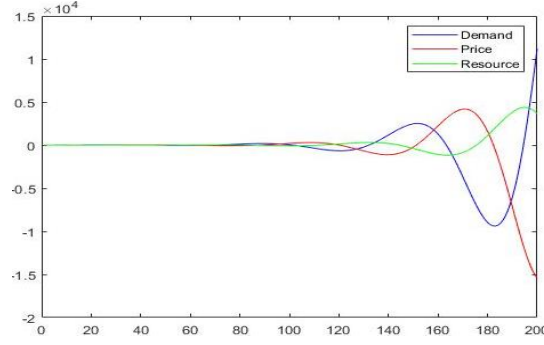Figure 1:Dynamics of systems under Scenario I (Divergent) with; $D_0=2$; $P_0=3$; $R_0=4$ where Red line represents price; Blue line represents demand; Green line represents resource availability

Figure 1 shows that system undergoes divergent behaviour due to interaction of demand, price and resource availability. Thus, as number of available resources increases, it effects demand due to lowering of price. Divergence of the system can be justified with the dynamics of CMfg in which there can be a wide range of virtually distributed resources adding to and leaving the system and as the system evolves, manufacturers add more resources to the cloud to generate revenue. However, this is not a real scenario in which quantities diverge in negative, therefore, to make system stable or converge, the memory of system is minimized using eigenvalue concept.

Based on the above situation, the corresponding eigenvalue for each equation is calculated and boundary of system between divergent and convergent behaviour is calculated by numerically adjusting the eigenvalues. The results are shown in figure 2
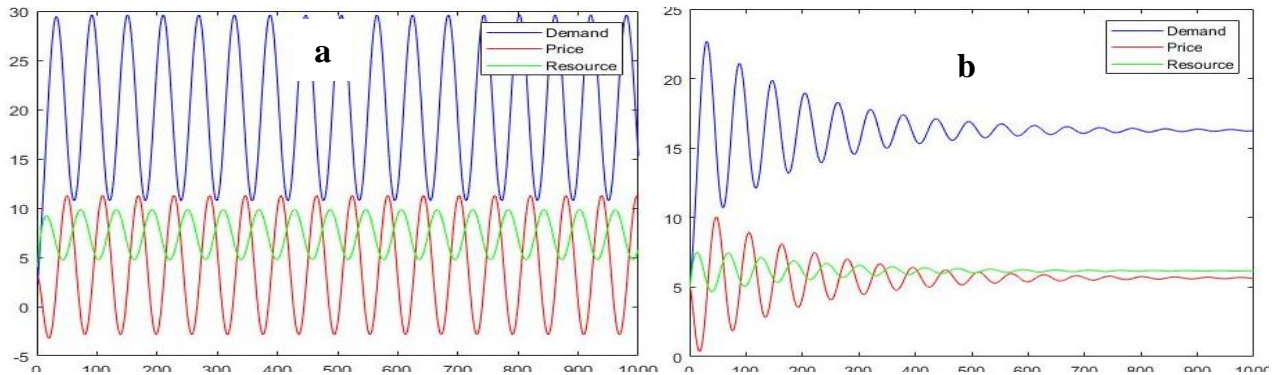


Figure 2: Eigenvalue adjustment (a) stable system on critical eigenvalues (b) converging system below critical values

Otherwise system can be made positive by adding minimum positive value constraint on the parameters. The result of this technique is shown in Figure 3
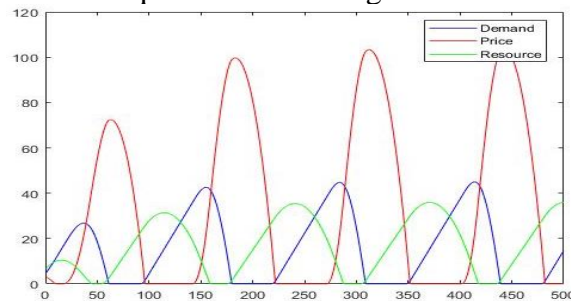


Figure 3: Diverging Behaviour of system with minimum positive value constraint

150

### 3.3.2. Scenario II: Natural Growth of Demand:

The demand is also affected by some external factor with a certain growth rate, thus placing a limit on maximum demand. To do so, Equation 1 is transformed to Equation 4 with remaining set of Equations 2 and 3.

$$D_t = D_{t-1} + \underbrace{\left(D_{t-1}*\left(1-D_{t-1}/\overline{D}\right)\right)}_{\text{Growth Factor}} + \left(\alpha*\left(1-P_{t-1}/\overline{P}\right)\right) \tag{4}$$

When system is simulated using real world scenario where $\overline{D}$=30; $\overline{P}$=20; $\overline{R}$=10, the results are presented in Figure 4 with different initial values.
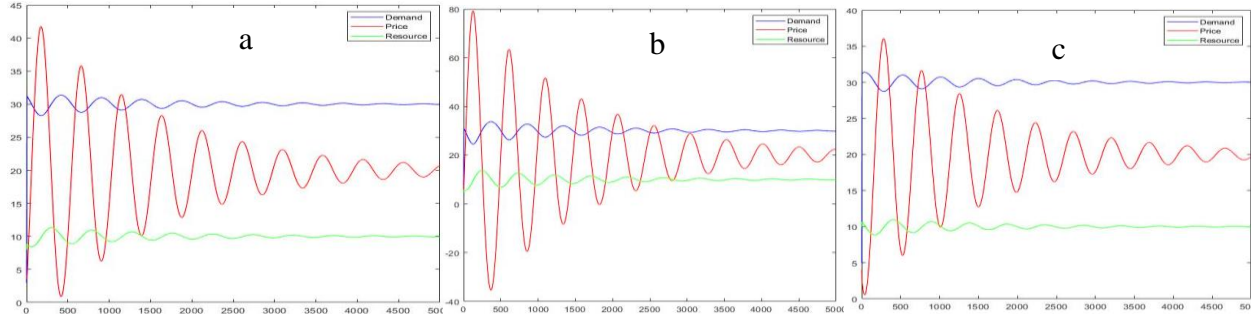


Figure 4: Dynamics of system under Scenario II (Convergence) (a) $D_0$=3; $P_0$=4; $R_0$=8.
(b) $D_0$=5; $P_0$=5; $R_0$=5. (c) $D_0$=5; $P_0$=4; $R_0$=10

Figure 4 shows that placing cap on demand factor, the system converges to a mean value as in predator-prey model. This convergence can be explained by the fact that increase in demand led to the shortage of available resources and made the resource curve steeper, as compared to scenario I described in Figure 1, due to which whole system shifted towards equilibrium [25]. In CMfg, external factors can be the ease of access to manufacturing resources which attract more customers and thus overall demand increases continuously. Thus, with more customer-centric system, the dynamics of prices greatly affect the working of a system. Therefore, manufacturers can generate more revenue by dynamically pricing their services based on customer demand

### 3.3.3. Scenario III: Restricted Available Resources

To analyse the significance of resource factor, the resource availability is restricted in this case. Since, with limited number of resources, it is not possible to increase the demand without the effect of resource limit, therefore, another factor is added to the demand equation which serves as *house-full* sign for the customers. Equation 1 is transformed to Equation 5 and Equation 3 is transformed to Equation 6.

$$D_t = D_{t-1} + \left(1-(P_{t-1}/\overline{P})*\alpha\right) - \underbrace{\left(R_{t-1}*\left(1-R_{t-1}/\overline{R}\right)\right)}_{\text{Restricted Resource Factor}} \tag{5}$$

$$R_t = R_{t-1} + \left(\left(1-\left(D_{t-1}/\overline{D}\right)\right)*\delta\right) + \underbrace{\left(R_{t-1}*\left(1-R_{t-1}/\overline{R}\right)\right)}_{\text{Restricted Resource Factor}} \tag{6}$$

When system is simulated using real world scenario where $\overline{D}$=60; $\overline{P}$=30; $\overline{R}$=20, the results are presented in Figure 5 with different initial values.

Figure 5 shows the natural convergent behaviour of system, since by limiting the resources also limits demand and thus system shift towards equilibrium, as in scenario II.
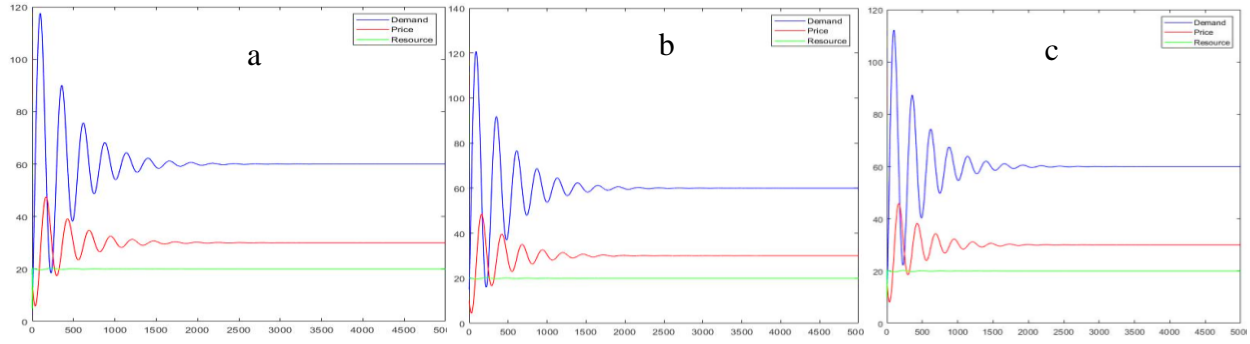


Figure 5: Dynamics of system under Scenario III (Convergence) (a) $D_0=15$; $P_0=10$; $R_0=20$. (b) $D_0=20$; $P_0=12$; $R_0=5$. (c) $D_0=20$; $P_0=15$; $R_0=10$

## 4. Findings

The results of simulation study show that the inter-connected variables lead to emergent properties of system in every iteration. Valuable findings from this study related to CMfg system and manufacturer's pricing includes;

- Variation in price has an obvious effect on the profit of manufacturer, however, resource availability can be controlled by manufacturer based on his profit margin and resource expenses. For instance, the resources that charge high maintenance cost can be avoided in low profit demand.
- Excess demand can be controlled with varying prices of same resources.
- Number of resources made available on cloud can be adjusted and cloud transaction cost can be avoided which, otherwise, is controlled by the cloud providers.

## 5. Conclusion

Complexity science enables to study the behaviour of complex systems by applying its tools and techniques. One such application has been presented in this work in which a complex CMfg system is considered. The emergent and dynamic behaviour customer demand in CMfg system facilitates manufacturers to generate revenue by pricing their resources according to customer demand. However, this aspect has not been given due attention in literature and so far, pricing in CMfg is based on CSP while keeping fixed prices of the manufacturers. To support this argument, a mathematical model is developed and simulated to study the effect of dynamic demand on prices and resource availability of manufacturers. Premise of our proposed model is based on known facts about manufacturers profit generation but these facts have not been attested in CMfg. Simulation results are compared with different scenarios where effect of demand restricted, and resource restricted environment on behaviour of CMfg is also studied. These results depict that with continuation of demand, CMfg system undergoes divergent behaviour. Therefore, considering fixed prices impedes profit maximization of manufacturer under such system. However, the restriction on demand and resources leads to convergence of the system in which system attains an equilibrium state. In such conditions, the selection of appropriate prices to bring the system in equilibrium can also affect overall profit of manufacturers.

## 5.1. Future Recommendations

Proposed work has highlighted a critical challenge of pricing in CMfg. However, there are some limitations in this work. To continue the current research direction, some suggestions for future work are provided:

- Cost factor can also be considered in estimating the price equation, since there is always some fixed cost and variable production cost which depends on resource available and price cannot be minimum than the fixed cost.
- Each resource price can be varied based on its costs and maintenance and then the effect of pricing on demand and resources can be analysed.
- The demand of each customer can be considered individually, henceforth, providing bulk order benefit to the customers with high demand rate.
- Market competition can be considered and effect of competitive prices on demand can be analysed.
- Since, cost of manufacturing decreases as the technology gets mature, therefore, this factor can also be added to the system along with price increase due to inflation.

## References

1.      Brailsford, S., et al. *Complex systems modeling for supply and demand in health and social care*. in *Simulation Conference (WSC), Proceedings of the 2011 Winter*. 2011. IEEE.
2.      Lewin, A.Y., *Application of complexity theory to organization science.* Organization Science, 1999. **10**(3): p. 215-215.
3.      Wu, D., et al., *Cloud manufacturing: Strategic vision and state-of-the-art.* Journal of Manufacturing Systems, 2013. **32**(4): p. 564-579.
4.      Brailsford, S.C. *System dynamics: What's in it for healthcare simulation modelers*. in *Simulation Conference, 2008. WSC 2008. Winter*. 2008. IEEE.
5.      Sayama, H., *Introduction to the modeling and analysis of complex systems*2015: Open SUNY Textbooks.
6.      Mönch, L., *Simulation-based benchmarking of production control schemes for complex manufacturing systems.* Control Engineering Practice, 2007. **15**(11): p. 1381-1393.
7.      Feistel, R. and W. Ebeling, *Evolution of complex systems: selforganisation, entropy and development*. Vol. 30. 1989: Springer.
8.      Lu, Y. and X. Xu. *Cloud manufacturing for a service-oriented paradigm shift*. in *Industrial Engineering and Engineering Management (IEEM), 2014 IEEE International Conference on*. 2014. IEEE.
9.      Wu, D., *Cloud-based design and manufacturing: a network perspective.* 2014.
10.     Ren, L., et al., *Cloud manufacturing: key characteristics and applications.* International Journal of Computer Integrated Manufacturing, 2014: p. 1-15.
11.     McAfee, R.P. and V. Te Velde, *Dynamic pricing in the airline industry.* forthcoming in Handbook on Economics and Information Systems, Ed: TJ Hendershott, Elsevier, 2006.
12.     Al-Roomi, M., et al., *Cloud computing pricing models: a survey.* International Journal of Grid & Distributed Computing, 2013. **6**(5): p. 93-106.
13.     Hsu, P.-F., S. Ray, and Y.-Y. Li-Hsieh, *Examining cloud computing adoption intention, pricing mechanism, and deployment model.* International Journal of Information Management, 2014. **34**(4): p. 474-488.

14. Li, H., J. Liu, and G. Tang. *A pricing algorithm for cloud computing resources*. in *Network Computing and Information Security (NCIS), 2011 International Conference on*. 2011. IEEE.

15. Sharma, B., et al. *Pricing cloud compute commodities: a novel financial economic model*. in *Proceedings of the 2012 12th IEEE/ACM international symposium on cluster, cloud and grid computing (ccgrid 2012)*. 2012. IEEE Computer Society.

16. Xu, H. and B. Li. *Maximizing revenue with dynamic cloud pricing: The infinite horizon case*. in *Communications (ICC), 2012 IEEE International Conference on*. 2012. IEEE.

17. Ranaldo, N. and E. Zimeo, *Capacity-driven utility model for service level agreement negotiation of cloud services*. Future Generation Computer Systems, 2016. **55**: p. 186-199.

18. GUO, W., W. PENG, and L. WANG, *The equilibrium of cloud manufacturing system under cost uncertainty*. Journal of Advanced Mechanical Design, Systems, and Manufacturing, 2017. **11**(2): p. JAMDSM0020-JAMDSM0020.

19. Li, Q. and J. Ma, *Research on price Stackelberg game model with probabilistic selling based on complex system theory*. Communications in Nonlinear Science and Numerical Simulation, 2016. **30**(1): p. 387-400.

20. Cheng, Y., et al. *Study on the utility model and utility equilibrium of resource service transaction in cloud manufacturing*. in *Industrial Engineering and Engineering Management (IEEM), 2010 IEEE International Conference on*. 2010. IEEE.

21. Lin, T.Y., et al., *Multi-centric management and optimized allocation of manufacturing resource and capability in cloud manufacturing system*. Proceedings of the Institution of Mechanical Engineers, Part B: Journal of Engineering Manufacture, 2016: p. 0954405415624364.

22. Argoneto, P. and P. Renna, *Supporting capacity sharing in the cloud manufacturing environment based on game theory and fuzzy logic*. Enterprise Information Systems, 2016. **10**(2): p. 193-210.

23. Liu, Y., et al., *Resource service sharing in cloud manufacturing based on the Gale–Shapley algorithm: advantages and challenge*. International Journal of Computer Integrated Manufacturing, 2017. **30**(4-5): p. 420-432.

24. Mankiw, N.G., *Teaching the principles of economics*. Eastern Economic Journal, 1998. **24**(4): p. 519-524.

25. Afonso, O. and P.B. Vasconcelos, *Computational economics: a concise introduction*. Vol. 25. 2015: Routledge.

# Physical Predictions from Dynamical Systems Modeling of an Electrodynamic Dissipative Structure

Ben De Bari, James Dixon, Bruce Kay
*University of Connecticut, USA*
Dilip Kondepudi
*Wake Forest University, USA*

**Abstract**

Oscillatory phenomena are ubiquitous to biological systems, and have been understood as fundamental to what biology does. Biological systems have been recognized as part of a class of non-equilibrium self-organizing systems. We study an electrodynamic non-equilibrium system, and utilize a dynamical systems model to account for its unanticipated oscillatory behavior. The Charge-Depletion Model (CDM) is built around assumptions of how charge accumulates in and is dissipated by the system, and accurately captures the oscillations observed in the physical system. The model is also used to make predictions of the behavior of the physical system, further corroborating the underlying assumptions of the model. Preferred states of increased dissipation support that the system is a dissipative structure, as well as a further hypothesis of Maximizing the rate of Entropy Production (MEP) in non-equilibrium systems. The phenomena of this and other non-equilibrium systems inform us on the nature of biological systems as physical systems.

## 1. Introduction

Oscillatory phenomena are integral to the behavior of biological systems, manifest across biological scales, from the neural central pattern generators [1], to inter-limb coordination in locomotion [2], to group-wide predator-prey population dynamics [3]. Rhythmicity and oscillations are also oft-encountered in dynamical systems theory, especially in non-conservative systems [4,5]. Non-conservative systems, like biological systems, are subject to energy flows; both involve intake and dissipation of energy or energy and matter. The rhythmicity of, for example, a forced mass-spring system comes about in part as a balancing of the intake and dissipation of energy [4,5]. Within this broad class of non-equilibrium systems are dissipative structures, self-organized flows of energy and matter which arise from, and in turn sustain, flows of energy [6]. These dissipative structures are functional in that their morphology comes about in order to increase the rate of entropy production [6]. It has been recognized that biological systems are dissipative structures [7-10]. Given that these dissipative structures depend on the flow of energy, we expect that they will tend towards states which increase that flow. The rate of entropy production is a measure of this flow of energy through the system. In each of these complex systems, emergent oscillatory phenomena can be understood as regulatory processes, coming about as a balancing of intake and dissipation of energy, resulting in a dynamic stability. We hold that thermodynamic metrics, in particular the rate of entropy production, can be used to predict and explain self-organization of complex systems and their consequent behavior.

Some electrodynamic dissipative systems have been studied by other groups, such as Jun et al., whose system consisted of conductive beads in oil subject to electric fields [11]. The author's work demonstrated that the system self-organizes into states of minimal electrical resistance; states which would increase the rate of energy dissipation and entropy production. Another system studied by Belkin et al., similarly evolved towards states of higher entropy production, and additionally exhibited quasi-periodic mechanical oscillations [12].

One such dissipative structure which our group works with is the Electrical Self-Organized Foraging Implementation (E-SOFI) [13]. The E-SOFI consists of conductive beads immersed in oil and subject to electric charge from a source electrode separated by an air-gap (Figure 1). The beads are surrounded by a metal grounded ring. The imposed charge incites dipole moments on the beads, and the forces drive the beads towards the grounding ring. Through all these mutual interactions, the beads self-organize into branching tree-like structures, extending out from the ring towards the source, and exhibiting swaying motions and traversing along the ring.

In previous work by our group, it was established that the system, in its morphological development and its behavior once formed, tends towards states of higher entropy production [13,14]. The rate of entropy production, $\Sigma$, is the product of the voltage and the current, divided by the temperature (Equation 1). Given the relative constancy of both temperature and voltage within trials, the rate of entropy production will be discussed with reference to the current through the system. The experiments herein serve to explore how its oscillations relate to the system's capacity for entropy production. The first experiment relates the current through the system to the trajectory of the terminal bead, and the following two experiments are tests of predictions from the model.
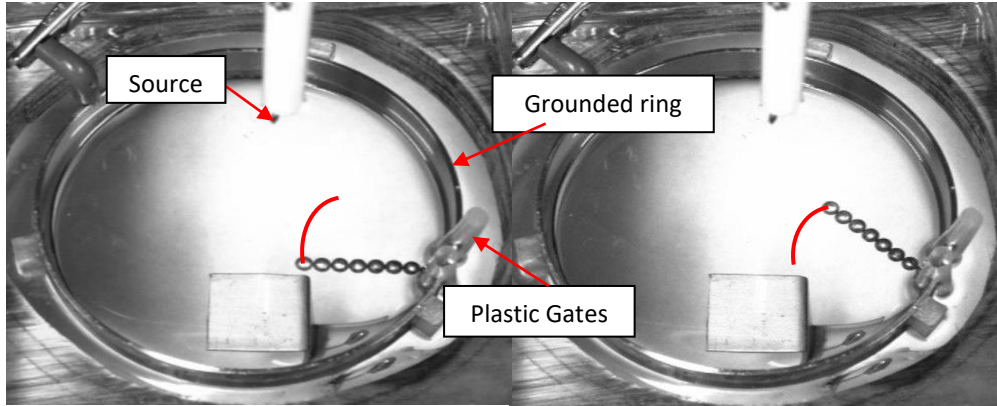
$$\Sigma = \frac{V * I}{T} \quad (1)$$



*Figure 1: Pictures of E-SOFI. Tree is at approximately each end-point of its oscillatory trajectory. The square plate beneath the dish is the stand for the magnet used in Experiment 3.*

## 2. Experiment 1: Oscillations in a Minimal Case
**Methods**

To isolate the oscillatory behavior for study, we modified the E-SOFI by imposing constraints on a singular tree (Figure 1). Plastic gates were placed on the ring electrode so that the bead at the base is unable to move along the ring, and the terminal bead's motion is limited to an arc. Under these constraints the system settles into relatively steady oscillations, sustained for trials as long as eight hours. The terminal bead nearest the source is the primary focus for the trajectory data.

The current through the system is recorded through analog to digital converters and Matlab programs written by Bruce Kay. The trajectory of the terminal bead is taken from video data, processed in ImageJ, and position and velocity data are determined by a Matlab tracking script written by Nicolas Oullett.

**Results**

The current and trajectory data from a subset of the trial is displayed in Figure 2. The trajectory data is presented as distance from the source, thus the troughs are the midpoint of the cycle (where the distance to the source is minimized) and the peaks correspond to either extrema of the bead's arc. Relative phase of the current peaks in cycles of bead displacement was determined to be -0.0368 radians on average; nearly coincident.
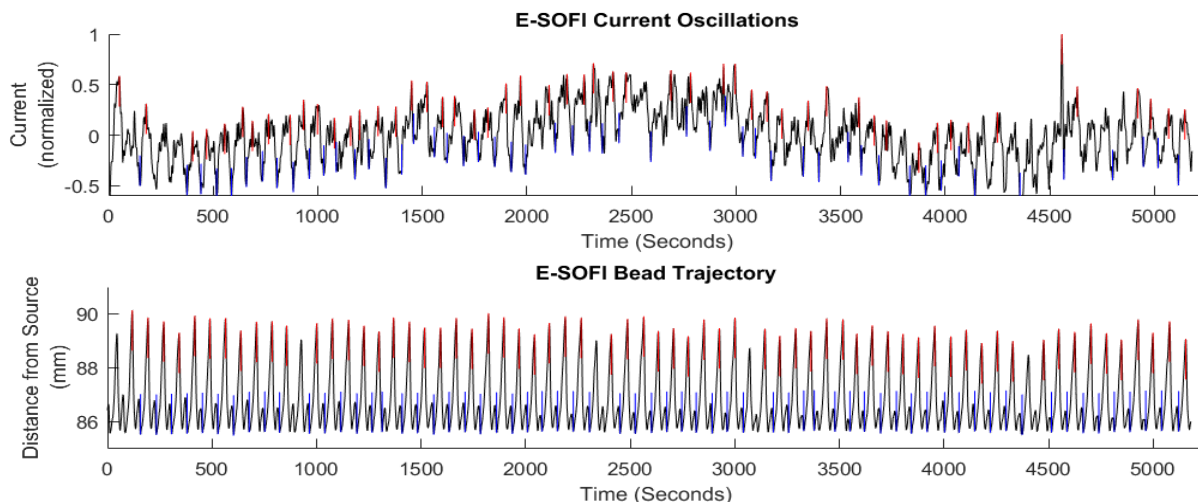


*Figure 2: E-Sofi current and position data. Peaks are picked out in red, troughs in blue. Current data is normalized for the accuracy of the peak-picking program.*

**Discussion**

The results of this experiment were unanticipated. Given an understanding that the system developed to states that maximized the flow of charge, it was expected that the current should be maximal when the tree is minimally displaced from the source, where it was expected there would be more charge. The opposite occurs in fact, with the current peaking nearly coincidentally (and slightly before) the bead is maximally displaced from the source. This quirk is what the Charge-Depletion Model (CDM) is aimed at explaining.

**3. The Charge-Depletion Model**

The intuition guiding the CDM is that the tree, by its conductivity, creates regions of low potential by conducting local charge to ground; hence charge-depletion. While the tree is depleting charge in one region of the dish, the constant supply from the source allows charge to build up on the oil surface at other regions of the dish. Given the tree's inclination towards states that increase dissipation, it will tend to move from regions of low potential to regions of high potential. The depletion of charge in one region by the tree and the saturation of charge in other regions interact to create the regions of alternating low and high potential, driving the E-SOFI's oscillations.

The CDM consists of two coupled equations, one for a one-dimensional charge distribution (representative of the portion of the electric field through which the terminal bead travels) and the other accounting for the forces on the bead. The charge distribution equation (Equation 2) has two terms, the negative term contributing to the depletion of charge and the positive term to its saturation, representing the current drawn by the tree and the current supplied by the source respectively. This equation determines the amount of charge $y$ at each location in a discretized

one-dimensional space $\{x_1, \dots, x_n\}$. $y_i$ is the amount of charge at the i$^{th}$ point, and $x_b$ is the location of the bead. $Cmax_i$ is the maximum possible charge at the i$^{th}$ point, a function of the capacitance of the oil and air. $\sigma$ determines the rate at which charge accumulates, representative of an applied voltage. The depletion term is most importantly a function of the location of the bead and varies with the inverse square of the distance as in the electric force. $c_1$ takes values between 0 and 1 and expresses the conductivity of the bead, i.e. the proportion of the charge at location $y_i$ which is conducted away. $c_2$ is a term to prevent the denominator from going to zero.

$$\dot{y} = -\frac{c_1 * y_i}{(x_i - x_b)^2 + c_2} + \sigma(Cmax_i - y_i) \quad (2)$$

The forces on the bead (Equation 3) are a damping coefficient $\beta$, representing the viscous contribution of the oil, and the force of the electric field on the bead. $v_b$ is the velocity of the bead. The force on a conducting sphere in an electric field is proportional to the product of the dipole moment and the gradient of the field. $\mu$ is representative of the dipole moment, and $\frac{dy}{dx}$ is a local derivative of the charge distribution, representing the gradient of the field.

$$\ddot{x}_b = -\beta v_b + \mu\frac{dy}{dx} \quad (3)$$

These two coupled equations adequately reproduce the oscillatory phenomena found in the E-SOFI, including the relative phase relationship of the current (calculated as the change in the total sum of charge $y$) and the bead's displacement from the source (Figure 3). Relative phase of current in terms of bead cycles under these parameters is found to be -0.569 radians on average; the two processes peak nearly coincidentally.
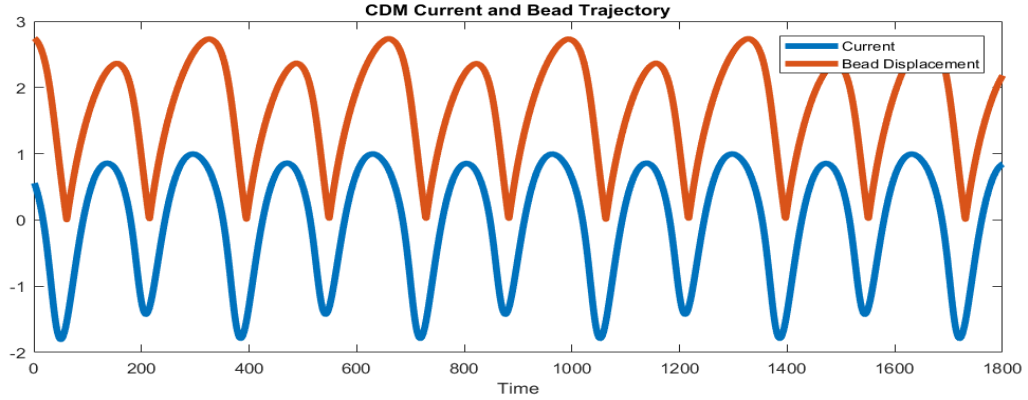


*Figure 3: Current and bead trajectory from CDM. Displacement, current, and time, are unitless with arbitrary values, and thus appear on the same axes. The shape of the charge distribution is parabolic with its maximum at x = 0 (here on the vertical axis). The displacement is distance from the source.*

The main phenomena of interest in the E-SOFI has been replicated successfully by the CDM. We find that for both systems the current is maximized when the bead is maximally displaced from the source. The guiding intuition regarding the nature of saturation and dissipation of charge in the dish is supported by the model. Additionally, the relative phase results from both the virtual and real system support the hypothesis of MEP. It was found that on average, the current peaks slightly before the bead reaches its maximal displacement. We suggest that this is evidence that the bead is continually moving towards states where there is more charge available; following a current peak, the charge available at that location begins to decrease and the bead moves up-gradient back towards regions of higher potential. Given that the flow of charge provides the forces

supporting the integrity of the tree, that the system tends towards states of increased current can be understood as inherently regulatory, reinforcing the conditions for its persistence. The system appears to be driven towards states that maximize access to and dissipation of charge, such that regulatory oscillations emerges to satisfy energetic constraints on the system.

## 4. Experiment Two: Voltage and Oscillation Frequency

If the behavior of the tree is a function of the charge available, then variations in the applied voltage are expected to have effects on the behavior of the tree. It was found that by varying parameter σ in the CDM, the term representative of voltage, the frequency of the bead's oscillations covaried linearly (Figure 4A), $R^2 = 0.9897$. Increased values of σ corresponded with increases in oscillation frequency. Experiment two sought this relationship in the E-SOFI.

### Methods

The E-SOFI is setup as in experiment one, with a singular tree gated by plastic constraints. Trials were approximately five minutes of constant voltage, from 14 to 26 kV. Video data of the tree was collected, and manually analyzed for oscillation frequency.

### Results

It was found that the oscillation frequency of the E-SOFI tree varied nearly linearly with increased voltage (Figure 4B), $R^2 = 0.9184$.
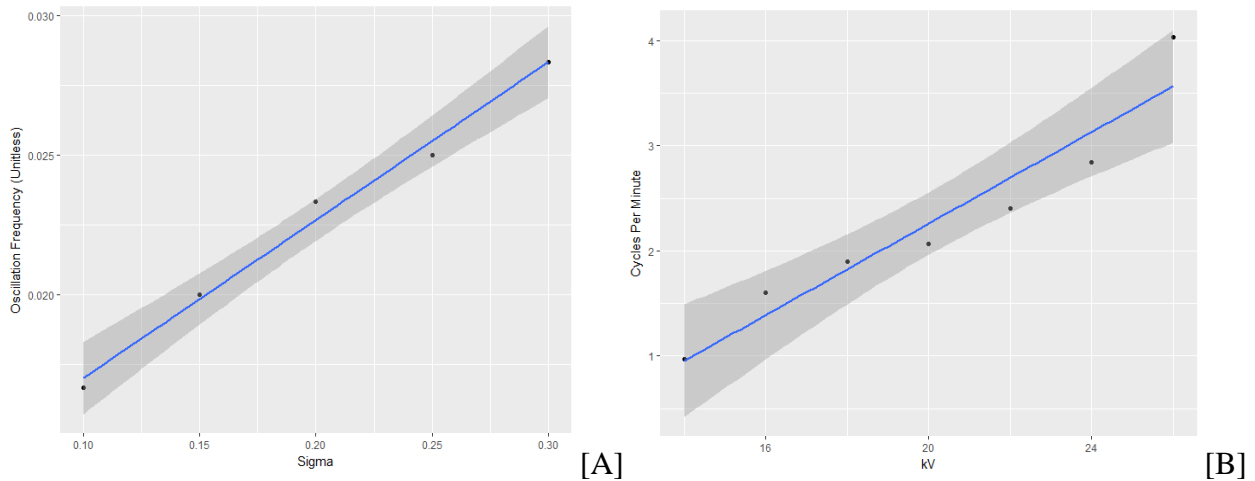


[A]       [B]

*Figure 4A: Oscillation frequency for varying values of sigma in the CDM. B: Oscillation frequency for varying voltage values in the E-SOFI.*

### Discussion

The hypothesis that oscillation frequency would vary with applied voltage was supported by both the CDM and E-SOFI data. These results can be interpreted in the context of MEP hypotheses; the behavior of the system is a function of its capacity to dissipate energy, and when more energy is accessible by the system, it will tend towards behaviors that acquire and dissipate that energy more quickly. These results are particularly interesting because they suggest that, in addition to morphological developments, behaviors of dissipative systems too are in accord with

159

increasing the rate of entropy production. Experiment three investigates other behaviors of the system and their relation to MEP principles.

## 5. Experiment Three: The Rebound Effect

From the understanding of how charge accumulates and is dissipated in the dish, we expect that charge will accumulate more in regions which the tree does not traverse. In experiment three, a new constraint is imposed on the bead (both virtual and real), locking it to one extreme of its cycle. While the bead is locked, charge in distal regions should accumulate to greater levels than if the tree had been free to access those regions. Once the constraint is lifted, we expect that the bead should move towards the other end of its cycle with greater velocity than if it had not been locked down. This anticipated increase in velocity is the rebound effect.

### Methods

In the CDM, this constraint is accomplished by adding another force to the bead's acceleration equation (Equation 4). The force is applied until movement stops, and then is removed, after which the bead resumes its oscillations.

$$\ddot{x}_b = -bv_b + \mu\frac{dy}{dx} + m(x_m) \quad (4)$$

In the E-SOFI, this is accomplished by replacing the terminal bead of the tree with a magnetic chrome bead of the same size, and introducing a magnet below the dish which locks the tree into place. In each trial, the tree oscillates for approximately five minutes, is locked with the magnet, released, and allowed to oscillate freely. The magnet is controlled by a stepper motor, and is programmed to move automatically to the same locking and release heights. Video data is collected, processed in ImageJ, and the trajectory and velocity of the terminal bead is calculated via a Matlab tracking program written by Nicolas Oullett.

| Release | Unconstrained | Trial |
|---------|---------------|-------|
| 51.28954 | 24.3683 | 1 |
| 53.62853 | 30.97155 | 2 |
| 58.94246 | 32.22161 | 3 |
| 53.81402 | 35.46877 | 4 |
| 60.21681 | 42.87612 | 5 |
| 68.28362 | 36.15041 | 6 |

*Table 1: Rebound and mean unconstrained velocities for each trial, mm/s*
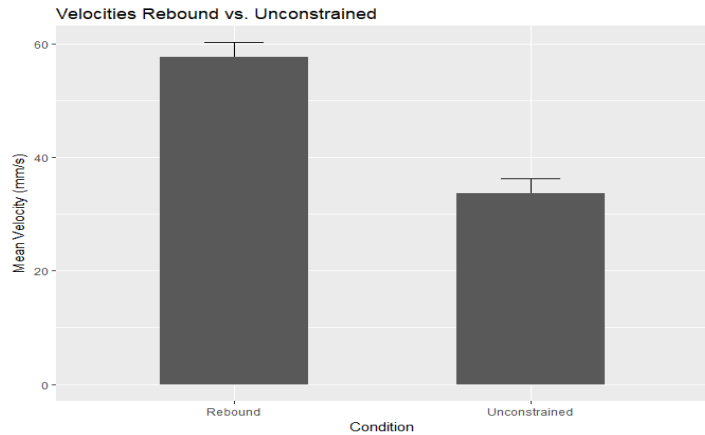


*Figure 5: Mean rebound and unconstrained velocities from E-SOFI*

### Results

The CDM exhibits the anticipated rebound effect (Figure 6) with a peak velocity immediately following the removal of the virtual magnetic constraint. Trajectory data from six trials was taken with the E-SOFI, and the peak velocities in each cycle were recorded. Mean

160

unconstrained and rebound peak velocities are shown in table 1. An independent samples T-test was conducted to determine the difference between rebound and unconstrained conditions. There is a significant difference between rebound ($M = 57.7$, $SD = 6.21$) and unconstrained ($M = 33.7$, $SD = 6.17$) conditions ($T(6.73)$, $p = 5.201 \times 10^{-5}$). Results are plotted in Figure 5.
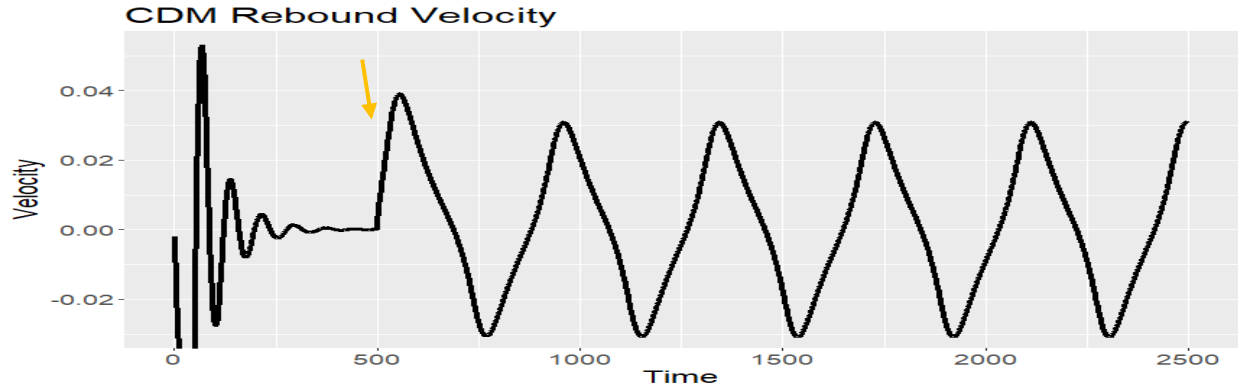


*Figure 6: Velocity profile from a run of the CDM. The magnetic force damps the bead until 500 time-steps. Upon removal of magnetic constraint (orange arrow) the velocity reaches a significantly greater value than under typical oscillatory conditions.*

## Discussion

The anticipated rebound effect was observed in both data from the CDM and the E-SOFI. This rebound effect was anticipated given the competition between the saturation and depletion of charge in the dish. That the effect is found in both the simulated and real system suggests that the CDM captures the essential features of the real behaviors, and that the charge-depletion hypothesis is a well-founded candidate explanation for the observed oscillatory phenomena. The rebound effect is consistent with MEP hypotheses as well, providing another situation in which the system is compelled towards states with higher energy available, and its behavior changes accordingly to reach those states.

## 6. Conclusions

That the E-SOFI would maximize the current while maximally displaced from the source was initially inexplicable. The experiments recounted here, and the work with the CDM, have served to further our understanding of the nature of the system, and to support an explanatory hypothesis. The strength of this work rests in the success of making physical predictions about the behavior of the E-SOFI, from the dynamical systems model of its behavior.

This work also has potential to further our understanding of behavior in the broader context of non-equilibrium systems. These results suggest that the behavior of self-organized systems is intrinsically related to the energetics; the system tends towards states with more charge available, which in turn allows for the persistence of the organization. This energy seeking behavior and consequent oscillations can be seen as inherently regulatory behaviors directed at maintaining a dynamic stability. While it would be difficult to support a claim that all behaviors of biological systems are directed only at increasing entropy production, it is exciting to observe a physical system which demonstrates end-directedness in its development and its behavior, and which exhibits such biological traits as energy-seeking, structural maintenance, and oscillatory regulatory phenomena. Study of these non-equilibrium systems is a promising endeavor which sheds light on the nature of biological systems and their behavior.

Works Cited

[1] F. Delcomyn, *Neural Basis of Rhythmic Behavior in Animals* Science, **210**, pp. 492-498, (1980)

[2] C. R. Gallistel, *The Organization of Action: A New Synthesis*, Ch. 4, 5. (1980).

[3] Lotka AJ, *Analytical Note on Certain Rhythmic Relations in Organic Systems,* Proceedings of the National Academy of Sciences of the United States of America, **6,** 7, pp. 410-415, (1920)

[4] B. Kay, *The Dimensionality of Movement Trajectories and the Degrees of Freedom Problem: A Tutorial*, Human Movement Science, **7**, pp. 343-364, (1988).

[5] R. Abraham, C. Shaw, *Dynamics: The Geometry of Behavior*, Ch. 2-3, (1982).

[6] D. Kondepudi, *Self-Organization, Entropy Production, and Physical Intelligence* Ecological Psychology, **24**, 33, (2012).

[7] Nicolis, G. *Physics of Far-From-Equilibrium Systems and Self Organization* New Physics, (1989).

[8] L. Von Bertalanffy, *The Model of the Open System*, General System Theory, pp. 139-154, (1973).

[9] S. Kaufmann, *Origins of Order* (Oxford University Press, Oxford, 1993).

[10] R. Rosen, *Essays on Life Itself* (Columbia University Press, New York, 2000).

[11] J. Jun *Dynamics of Self-Organization of Ramified Patterns in an Electrochemical System.* University of Illinois, Urbana-Champaign (2006).

[12] A. Belkin, *Self-Assembled Wiggling Structures and the Principle of Maximum Entropy Production*, Scientific Reports (Illinois), **5**, 1, (2015).

[13] J. Dixon, *End-Directedness and Context in Non-Living Dissipative Systems* World Scientific Review, **20**, 185, (2015).

[14] D. Kondepudi, *End-Directed Evolution and the Emergence of Energy-Seeking Behavior in a Complex System* Physical Review, **91**, 1, (2015).

# Long Range Synchrony can be achieved through Reaction-Diffusion processes in *Physarum polycephalum*

Abid Haque[a], Subash Ray[a], Gregory Weber[a], Simon Garnier[a]

[a] Federated Department of Biological Sciences, Rutgers Newark – NJIT, Newark, NJ 07102

Long range synchrony occurs when spatially distant parts of a system exhibit highly correlated behavior, even though they are not all directly exchanging information with each other. Synchronized behavior enables organisms to coordinate parts of their body to perform complex tasks such as walking, foraging for food, or communicating with each other. The high selective advantage of synchronized behavior has led to the evolution of complex centralized neural architectures, such as the brain. However, a majority of living organisms on earth do not possess such centralized control over their behavior, despite facing the same fundamental evolutionary pressures. Therefore, it becomes important to understand the mechanisms responsible for the emergence of synchronized behavior in the absence of centralized control. The protist *Physarum polycephalum,* commonly known as the acellular slime mold, has emerged as a model system for studying decentralized control systems that exhibit synchronized behavior over large distances. The slime mold is a unicellular organism which exists as a network of interconnected tubules, which transport nutrients and cell signals between various parts of the cell body. The tubules exhibit highly synchronized diametric contractions across distances spanning the length of the cell body. This enables the cytoplasm within the tubules to efficiently shuttle back and forth due to hydrostatic pressure differences. The exact mechanism through which the oscillations maintain synchrony across large distances is unknown, but experimental studies agree that an underlying $Ca^{2+}$ oscillator is locally responsible for the mechanical contractions. This study attempts a bottom-up approach to explain the emergence of long-range synchrony from local interactions between neighboring sections of the tubule. Using a Reaction-Diffusion model to study the biochemical oscillator system, we aim to estimate the parameters and nature of the dynamical system that leads to the emergence of synchronized oscillations. In the described model, a tubule is modeled as a linearly connected system of oscillators. We have encoded each oscillator using a two-dimensional system of ODEs, through which the $Ca^{2+}$ "signal" periodically changes its concentration. The oscillators are coupled, such that the chemical signal follows a concentration gradient, and diffuses into its neighboring oscillators. The key outcome of the simulation is the emergence of a traveling wave of signal concentrations in a tubule (Fig. 1). This traveling wave can periodically produce membrane contractions, which would enable shuttle streaming through periodically reversible pressure gradients. Slime molds have recently gained attention for their decision-making abilities, and synchronized membrane oscillation patterns are thought to be responsible for this cognitive-like behavior. Our study adds support to the growing evidence that cognitive behavior can emerge from localized interactions. Similar complex systems approaches can help us broaden the definition of cognition, and study its evolutionary origins.
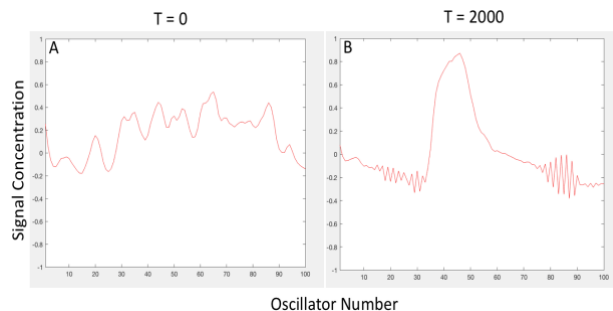


**Fig.1**: Traveling waves of signal concentration emerge in a linearly connected oscillator system. (A) Oscillators randomly initialized at timestep T = 0. (B) Emergence of a traveling wave (moving from left to right), due to diffusive coupling of oscillators. Snapshot shown here at timestep T = 2000.

# How Rankings Go Wrong:
## Structural Bias in Common Ranking Systems Viewed as Complex Systems
## work in progress

Patrick Grim, Jared Stolove, Natalia Jenuwine, Adrian Apaza, Hannah vanWingen, Jaikishan Prasad, Paulina Knoblock, Callum Hutchinson, Chengxi Li, Kyle Fitzpatrick, Chang Xu & Catherine Ming

Center for the Study of Complex Systems, University of Michigan

## Abstract

We introduce agent-based techniques to analyze inherent structural bias in abstract models of common ranking systems such as PageRank, HITs, and Reddit. In the complex dynamics of reputational loops, an element's ranking itself influences factors in terms of which rank is calculated, resulting in the amplification of divergence and the exaggeration of small random and path-dependent differences. Agent-based models of basic algorithms employed in PageRank, HIT, and Reddit are constructed. These models allow comparisons of bias dynamics and effects across a number of parameters.

## Introduction

In many real world examples, such as college ranking and online search, objects with very similar quality can end up with significantly different rank. Sometimes the data that a ranking system relies on may be dubious. But even in the best of conditions and with the cleanest data input, we would argue, the very structure of some familiar ranking systems can result in distorted informational output.

The basic idea of all familiar ranking systems is an attempt to read objective quality— what sites, papers, or posts are genuinely worth reading—from social measures of what is read and responded to by whom. In any such system there will be a looping factor: a site, paper, or post will be widely read because it is ranked highly, but will be ranked highly precisely because it is widely read. How severe the loop—and how far rank will come unglued from quality—will depend on the degree to which users base their choice of a site, a paper, or a post on rank as opposed to some independent judgment of quality.

The extent to which looping constitutes a distorting factor, however, will differ with different ranking algorithms. What we offer here is an advanced sketch of approaches with an incomplete sample of results regarding three familiar algorithms—PageRank, HITs and Reddit. A more complete paper with a fuller development of techniques and results will appear elsewhere.

## PageRank

PageRank analyzes the traffic between sites in order to determine the ranking of the websites. Using a network structure, PageRank treats each website as a node and hyperlinks between websites as directed links between the nodes. Each site is initially assigned a node value of 1 divided by the total number of nodes. At each time interval, PageRank divides the node value for each node by the number of outgoing links from that node, and this value is sent as an incoming value to the node at the other end of each out link. Each node's value is then replaced by the sum of its incoming values. PageRank also redistributes an extremely small amount of value equally between the nodes regardless of links, representing the possibility of an individual typing a URL

directly into the search bar rather than clicking a node. PageRank then ranks the sites in order of greatest to least node value.

We construct a series of simulations in which pages are assigned an inherent 'quality' between 1 and 100. At each step in the evolution of the model a small number of pages are added to the network—much as pages are progressively added on the internet, and roughly as nodes are added in a preferential attachment network (Fig. 1)
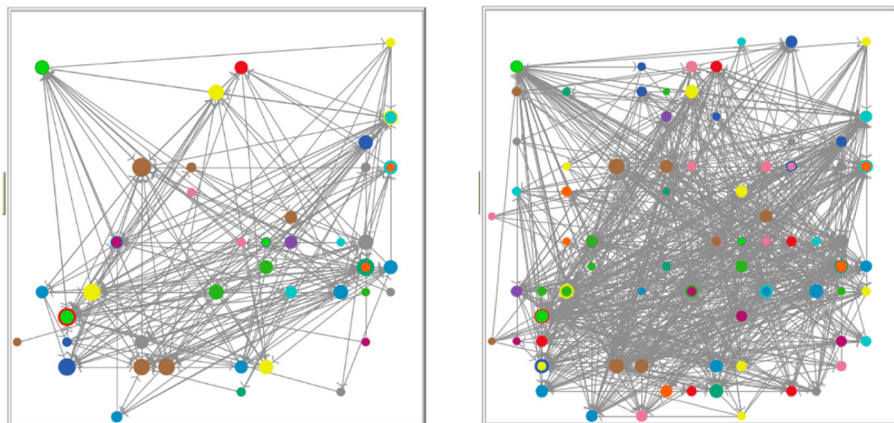


Figure 1. An example of the general progression over time of a model in which links from a node x to y are established probabilistically in terms of the quality of y alone. Stages 50 and 100 in a typical model evolution shown.

The relationship between rank and quality changes significantly as the proportion of link formation determined by rank increases. When links are determined by quality alone, the most highly-ranked pages all have relatively high quality. While many high-quality pages end up with low rank, low-quality pages cannot receive rank above a certain quantity. Also, the spread in rank between pages is moderate, with the most successful pages receiving around 3.5 times the average amount of rank. Lower quality pages are able to obtain a higher rank; and the difference between the highest and lowest ranked pages becomes much larger (Fig. 2).
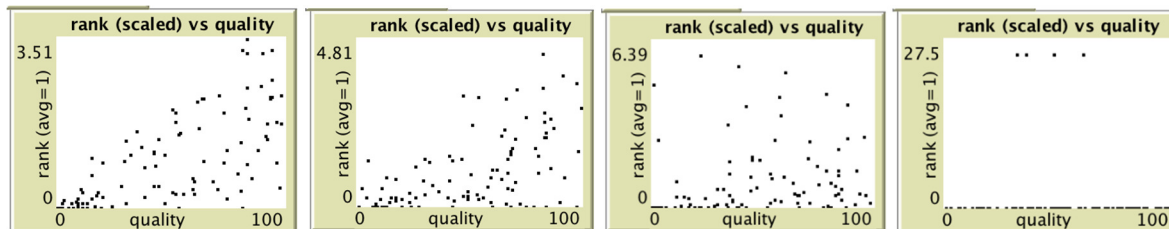


Figure 2. Rank vs. quality as the probability of link formation is calculated in terms of quality alone (left), rank .3, .7 and entirely in terms of rank (right).

We introduce 'discrepancy' as a measure for divergence of rank from quality. Over a wider range of cases, increasing discrepancy with increasing weight given to rank is shown in Figure 3.
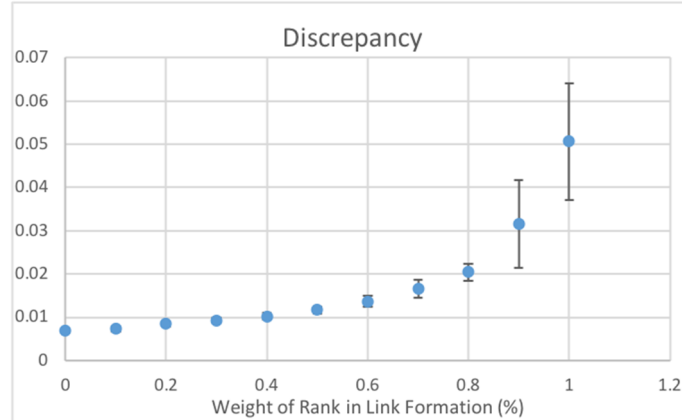
Figure 3. Weight of attention to rank as opposed to quality in link formation and the resultant discrepancy between modelled rank and quality of PageRank sites.

As agents respond more to rank—what we have targeted as the more realistic case—the discrepancy between rank and quality is higher, the average quality of a page linked to is lower, and fewer pages can dominate. If links are formed entirely in terms of rank, not surprisingly, the quality of pages linked to is essentially random. As rank is considered more, pages which are established early have a decisive advantage.

**HITs**

The HITS Algorithm is used by Academia.edu, a social networking website that shares papers and monitors their impact. . In general, for each web page, the scheme assigns two values recursively. One is Hub Score, the sum of authority scores of all the nodes that the site points to. The other is Authority Score, the sum of the hub scores of all nodes pointing to this specific web page. A page that suggests more widely recognized authorities would obtain a higher Hub Score, whereas a page that is recommended by more quality hubs would achieve higher Authority Score. Hence, the results produced by HITS Algorithm has an interaction between these two parameters.

In our model for the HITS algorithm both authors and papers are assigned a constant built-in 'quality' value. During setup, authors decide which papers to recommend solely based on the paper quality. Subsequently, papers are given authority scores based on the authoritativeness and number of authors that have recommended them, while authors are given hub scores calculated by the popularity and 'quality' of their papers. One consequence—both in our model and in a real instantiation of HITS—is that the ranking of an author whose papers are popular in the academic community but of average quality may carry more weight than that of another author with less popular but higher quality articles.

When paper judgment is set at 95% in terms of quality, we observe an expected correlation between quality and rank (Figs. 4, 5). At such a setting, authors are recommending papers on the basis of quality, for the most part. A small divergence remains due to the stochastic nature of the model.
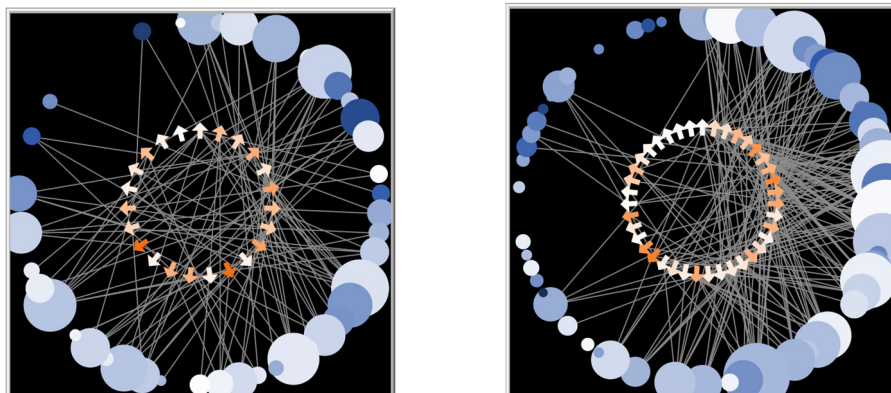
166

Fig. 4. A typical run of the model based solely on quality, steps 12 and 14 shown. Blue circles represent papers, with size representing their scores relative to other papers. Color saturation (blueness) represnts the degree of difference between HITS rank and 'quality.' Orange arrows represent authors, with saturation of orange representing the difference between their ranking and the average quality of their papers.
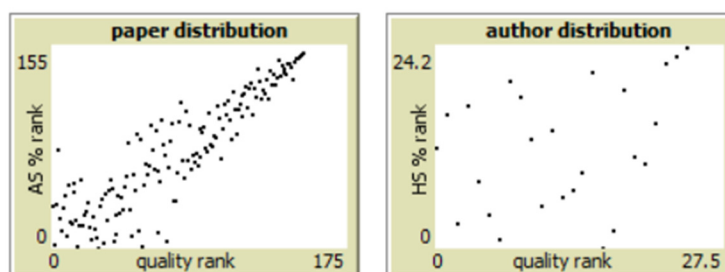


Fig. 5 Rank plotted against quality for papers (left) and authors (right) with links established with a setting of 95% on the basis of quality.

At a 'percentage-by-paper-quality' of 5%, on the other hand, authors make recommendations almost entirely in terms of established rank rather than the quality of each paper. In this case see no positive correlation between quality and rank (Figs. 6, 7).
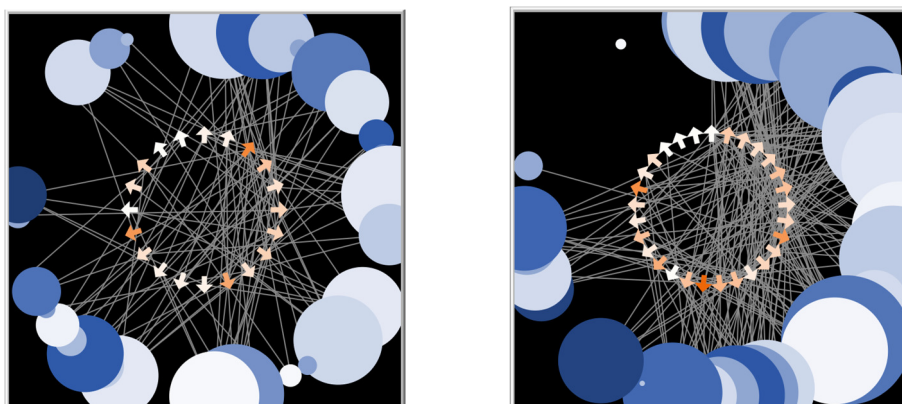


Fig. 6. A typical run of the model with ranking independent of quality, steps 12 and 14 shown. Blue circles represent papers, with size representing their scores relative to other papers. Color saturation (blueness) represents the degree of difference between HITS rank and 'quality.'

Orange arrows represent authors, with saturation of orange representing the difference between their ranking and the average quality of their papers.
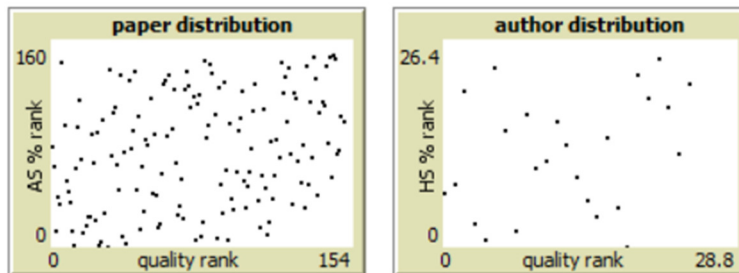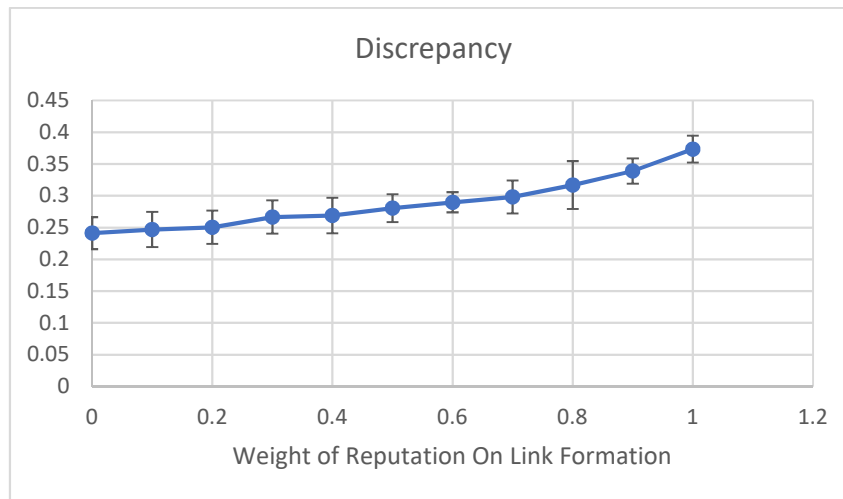


Fig. 7  Rank plotted against quality for papers (left) and authors (right) with links established with a setting of 95% on the basis of rank, only 5% on the basis of quality.

The HITS algorithm has a discrepancy that is an order of magnitude larger than that of PageRank.  As a recursive function with two reinforcing variables a higher degree of error propagation is anticipated.



**Reddit**

Reddit.com is a social news aggregation site where users can discuss and rate posted content. Content is distributed by topic among various 'subreddits' where posts with the highest net rankings (upvotes minus downvotes) rise to the top of the page. By ranking based on net score, Reddit is intended to create a preference toward non-controversial content, since a post which receives 50 upvotes and zero downvotes will be ranked the same as a post with 500 upvotes and 450 downvotes. Comments on posts are also voted upon and change order correspondingly. Reddit's main Front Page shows posts with the most upvotes across all subreddits, with the order of links changing constantly based both on the time of submission and user votes. A post's score will not decrease as time passes but newer posts will get a higher score.

Our model of Reddit simulates website users reading posts and voting on them. We model users as deciding to read a post and to give it an upvote or downvote based on two factors: its objective 'quality' (a number between 0 and 100) and its rank as determined by the Reddit algorithm.  Each modeled user is assigned a threshold that determines how likely they are to read,

upvote, or downvote a given post on average. During the simulation, users periodically "leave" the site, replaced by users with different thresholds. New posts with assigned 'qualities' are periodically created as well, then ranked according to Reddit's algorithm.

In order to get a baseline measure to help us evaluate the algorithm, we start off with two highly unrealistic situations: a "best case scenario" for both reading and voting and a "worst case scenario." The best case scenario considers a situation in which users are able to tell the objective quality of a post before even reading it, and make their decisions to read and vote based solely on quality. The worst case scenario describes a situation in which users read and vote on posts based solely off of rank, and therefore rank is totally unrelated to quality.
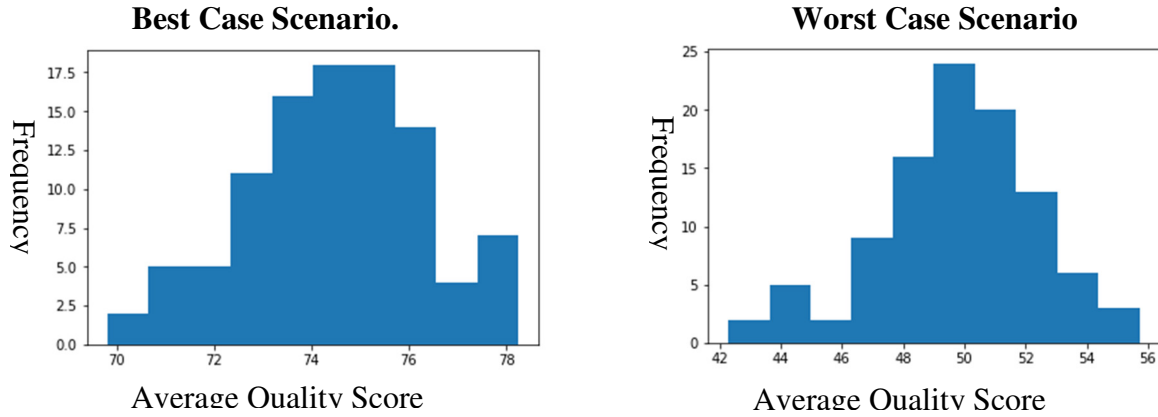


Figure 8. The distribution of average quality of a read post, or "average quality score", in both the best and worst case scenarios over 100 runs of the model.

The goal of any ranking algorithm is to offer users posts of high quality. We evaluate how well Reddit performs given users who base their reading and voting more on quality or more on rank. Figure 9 shows an evaluation of the Reddit algorithm's robustness to rank-bias based on how close the average score is to the best case scenario when we introduce varying levels at which reading and voting are based on rank.

|  |  | Voting Bias | | | | | |
|---|---|---|---|---|---|---|---|
|  |  | 0 | 0.2 | 0.4 | 0.6 | 0.8 | 1 |
| Reading Bias | 0 | 74.9 | 74.399 | 74.774 | 74.26 | 72.844 | 72.33 |
|  | 0.2 | 73.51 | 72.84 | 73.477 | 74.62 | 73.94 | 71.78 |
|  | 0.4 | 73.73 | 72.21 | 75.16 | 71.22 | 72.17 | 70.82 |
|  | 0.6 | 70.97 | 71.057 | 72.105 | 71.033 | 71.466 | 68.788 |
|  | 0.8 | 69.561 | 68.34 | 68.865 | 67.068 | 70.044 | 60.78 |
|  | 1 | 68.3827 | 66.83 | 65.77 | 65.98 | 63.01 | 49.025 |

Figure 2

Figure 9. The average quality of a post read in versions of the model with varying levels of Reading Bias and Voting Bias. Cells that are relatively closer to the best case scenario are colored green, those relatively closer to the worst case scenario are colored red.

The figure reveals that the algorithm performs quite well even at relatively high levels of both kinds of bias. For example, even with both bias parameters set to .6 rank, the resulting average quality score is 71.033, much closer to the ideal scenario than the worst case scenario.

Reddit's robustness in the face of rank-related bias appears to be due to the built-in penalization of older posts in the ranking algorithm. Reddit's algorithm provides newer posts with an advantage in two ways. First, newer posts are given a one-time increase in their initial raw score, called the "time damp." Second, additional upvotes only impact a posts' score logarithmically, so highly ranked posts must acquire exponentially more votes to counteract the advantage given to newer posts. PageRank, we've seen, is dominated by a small number of pages. Reddit's linear distribution of votes and quality avoids that consequence.

**The Looping Effect and Bias in Familiar Ranking Systems**

The basic idea of all familiar ranking systems is an attempt to read objective quality— what sites, papers, or posts are genuinely worth reading—from social measures of what is read and responded to by whom. In any such system there will be a looping factor: a site, paper, or post is widely read because it is ranked highly, but it is ranked highly precisely because it is widely read. The existence of loops in ranking systems appears to be inevitable. The degree to which particular algorithms are vulnerable to such loops, however, may differ. Our attempt here, using simple abstract models of three current systems, has been to make some first steps in evaluating that relative vulnerability.

# Application of Multi-Depot Multi-TSP Model to Humanitarian Logistics in Times of Disaster.

**Hayford K. D. Adjavor and Yong Wang**
**State University of New York - Binghamton.**

## Abstract

Humanitarian logistics have continued to enjoy attention in academia, and researchers from fields including engineering, natural, and social sciences have continued to work in this area. Evacuation of people and the distribution of relief items to affected communities in times of disaster require enormous work from all stakeholders in the humanitarian relief chain & logistics. There is a great sense of urgency on the part of governments, relief workers, public and private humanitarian organizations, as well as researchers. This urgent need to evacuate people and to deliver relief items, means that transportation (e.g. air, road, rail, marine), becomes a crucial part of the entire relief management effort. Consequently, the need for transportation during the pre-disaster planning (preparedness & response phase), and post-disaster (recovery & clean-up phase) cannot be overemphasized. This paper, therefore, investigates the evacuation efforts during Hurricane Katrina in the city of New Orleans, Louisiana. We first formulate the problem as a multi-depot multi-traveling salesperson problem *(MmTSP)*; well-known NP-hard problem. Network science technique is then used to corroborate and/or rationalize the best number of shelters needed. Geographic Information System *(GIS)* - *ArchGIS* - was then employed to provide a network map formed by the number of affected areas and shelters. To find feasible solutions to the evacuation problem, genetic algorithm *(GA)* technique was subsequently employed. The main finding of this work is the potential of the described strategy for future ground-breaking works in this growing field of humanitarian relief chain & logistics.

## Keywords

humanitarian logistics, network science, genetic algorithm

# Dynamics of Characteristics Associated with Dental Opioid Prescription Patterns in Emergency Departments

Shabnam Seyedzdaeh Sabounchi[1], Nasim S.Sabounchi[2], A. Serdar Atav[3], Leon Cosler[4]

[1] College of Community and Public Affairs,
Binghamton University, State University of New York
shabnam@binghamton.edu

[2] System Science and Industrial Engineering Department,
Binghamton University, State University of New York
Sabounchi@binghamton.edu

[3] Decker School of Nursing, Binghamton University, State University of New York
atav@binghamton.edu

[4] School of Pharmacy and Pharmaceutical Sciences,
Binghamton University, State University of New York
lcosler@binghamton.edu

## Abstract

The legitimate prescription of opioids to treat chronic non-cancer pain as well as the patients' opioid misuse has been increasing since late 1990's. Among different health care settings, emergency departments (EDs) have the highest number of patient visits with pain. Number of non-traumatic dental problem visits to EDs has also increased largely from 2000 to 2010. Our analysis originates from the systems theory approach to explain the dynamic complexity that characterizes the opioid epidemic in ED dental patient visits. The causal model shows the interactions between different variables and describes feedback structures that govern the dynamics of opioid prescriptions. By incorporating system dynamics modeling we are able to recommend modifications to current strategies for fighting the opioid epidemic. In this study, we analyze data from the National Hospital Ambulatory Medical Care Survey (NHAMCS) from 2002 to 2015. Patient visits with non-traumatic dental problems to EDs within the 50 states of United States and District of Columbia are investigated. Our findings demonstrate that the rate of opioid prescriptions has remained steady and constitutes more than half of all pain medication prescriptions for patients with dental pain. We use the data to validate the causal structure that underlies prescribing behaviors for opioid.

# Numerical Simulation of the Conduction and Propagation of Spatiotemporal Electrodynamics on the Heart Surface

Bing Yao

Hui Yang, Ph.D.

*Complex Systems Monitoring, Modeling and Analysis Laboratory,*

*The Pennsylvania State University, University Park, 16802, USA*

# ABSTRACT

Heart disease is one of the leading causes of morbidity and mortality in the United States. Effective medical treatments for heart diseases are of paramount importance to the patients' health and the whole family's wellbeing, thereby broadly impacting the society. To improve the cardiac care services, it is important to develop a better understanding of disease-altered cardiac electrodynamics. Numerical models and computer simulations facilitate the quantitative elucidation and understanding of heart functions. However, the electrical activities of the heart are varying in both space and time, and the human heart is with highly complex geometry. As such, simulation modeling of the conduction and propagation of cardiac electrodynamics involves a great level of complexity. In this work, we propose a novel efficient method to numerically simulate the conduction and propagation of spatiotemporal electrodynamics on the complex heart surface. Specifically, we first project the complex 3D heart surface into a 2D graph using the method of t-distributed stochastic neighbor embedding. Second, we propose a moving-least-square mesh-free method to simulate the nonlinear space-time electrodynamics of the heart on the 2D graph. Simulation results demonstrate that the proposed model not only efficiently simulates the cardiac electrodynamics, but also will be an effective tool to assist medical scientists in the heart-disease investigation.

# Fractal Pattern Recognition of Image Profiles for Manufacturing Process Monitoring and Control

Farhad Imani[1], Bing Yao[1], Ruimin Chen[1], Prahalada Rao[2] and Hui Yang[1]

[1] Department of Industrial and Manufacturing Engineering,
The Pennsylvania State University, University Park, PA
[2] Department of Mechanical and Materials Engineering,
University of Nebraska, Lincoln, NE

## Abstract

Nowadays manufacturing industry faces increasing demands to customize products according to personal needs. This trend leads to a proliferation of complex product designs. To cope with this complexity, manufacturing systems are equipped with advanced sensing capabilities. However, traditional statistical process control methods are not concerned with the stream of in-process imaging data. Also, very little has been done to investigate nonlinearity, irregularity, and inhomogeneity in image stream collected from manufacturing processes. This paper presents the multifractal spectrum and lacunarity measures to characterize irregular and inhomogeneous patterns of image profiles, as well as detect the hidden dynamics of the underlying manufacturing process. Experimental studies show that the proposed method not only effectively characterizes the surface finishes for quality control of ultra-precision machining but also provides an effective model to link process parameters with fractal characteristics of in-process images acquired from additive manufacturing. This, in turn, will allow a swift response to processes changes and consequently reduce the number of defective products. The proposed fractal method has strong potentials to be applied for process monitoring and control in a variety of domains such as ultra-precision machining, additive manufacturing, and biomanufacturing.

## A Causal Modeling Approach for the Study of Association Among Food Intake, Exercise, And Mental Distress

There is limited evidence available for the relationship between diet and mental wellbeing. Our goal in this study is to evaluate the effect of healthy diet, exercise, and healthy practices on mental wellbeing and compare the significance between younger (18–29 years) and older adults (30 years and older). We used 563 records of data collected internationally through different social media to identify these causal relationships. We applied the backward regression methods and found that in the young age groups, higher consumption of fast food (more than three times per week), negatively affects mental well-being. However, older adults achieve higher mental wellbeing with lower consumption of carbohydrates (such as rice and bread) and more fruit. This is in part due to the functionality of brain which is related to age, food and mental wellbeing. We developed a casual loop diagram shown in the figure below to demonstrate our hypothesis of these causal relationships that describe more healthy diet, healthy practice and exercise can increase brains level of dopamine which leads to higher mental well-being and consequently higher motivation to improve healthiness.
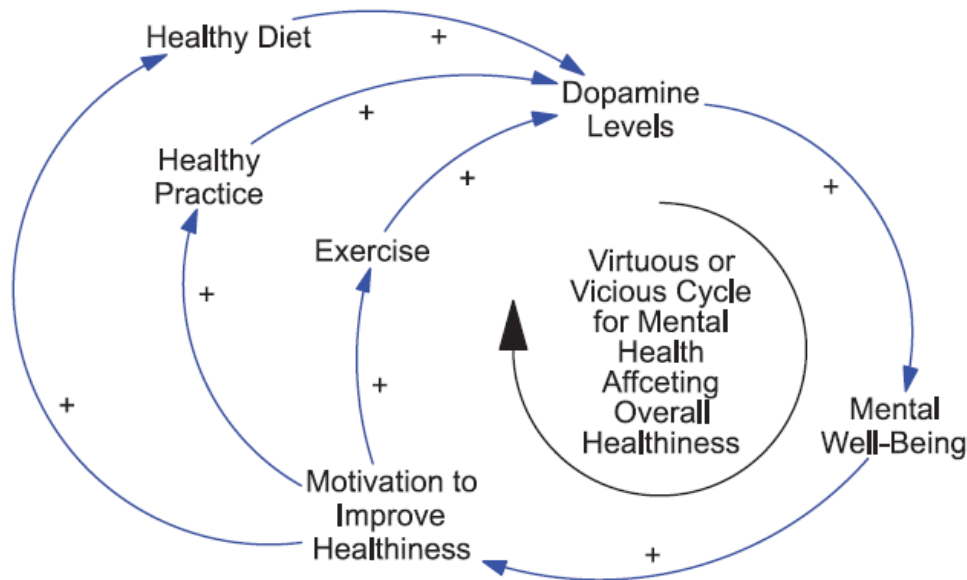


**Figure:** Casual loop diagram
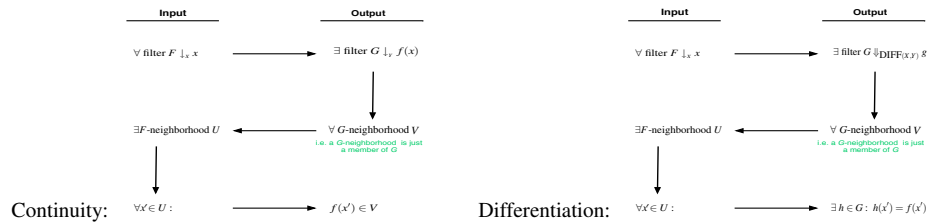
# Heterotic Complex Systems

Howard A. Blair

EECS, Syracuse University

Heterotic dynamical systems seamlessly combine dynamical system models such as classical/quantum, discrete/continuous, etc. To rigorously specify these dynamical systems as systems of simultaneous first-order differential equations over diverse types, differentiation on convergence spaces, can be used. Convergence spaces, built on the work of Henri Cartan who began with filters of sets as the basic notion to treat convergence, as distinct from open sets, is elucidated in full in Bourbaki [1]. Differentiation is performed natively, avoiding artificial numerical encodings of the elements of these domains, and allows for specifying the evolution of continuous and discrete states.

Category-theoretic characterizations of heterotic models of computation were introduced by Stepney et al. in [2]. Heterotic models of computation combine computational models such as classical/quantum (with restrictions on coupling of classical and quantum variables), digital/analog, synchronous/asynchronous, imperative/functional/relational, etc. to obtain increased computational power, both practically and theoretically. The term *heterotic* refers to extreme hybridization in the form of the seamless combining of diverse system paradigms. A relatively simple example can serve to illustrate a heterotic dynamical system.

Consider the set of real numbers with the usual basic Euclidean metric space topology where $d(x_1, x_2) = |x_1 - x_2|$. This metric space is a *convergence structure*. We can add the additional structure of a directed graph: the vertices of the digraph are the integers as a subset of the real numbers, and the set of edges is $\{(n, n-1) \mid n \text{ an integer}\} \cup \{(n, n+1) \mid n \text{ an integer}\}$. The result is a nontopological convergence space. Patten's work, [3] on discrete structure convergence spaces applies to formulating heterotic dynamical systems involving discrete structure data-types.

A *convergence space* [4] is a nonempty set $X$ of elements where each $x \in X$ has an associated family of filters, each of which is said to *converge* to $x$. The family of filters is closed under reverse inclusion and contains the *point*-filter consisting of all subsets of $X$ containing $x$. The class of all convergence spaces forms a cartesian closed category CONV, that in turn makes function spaces uniformly obtainable as convergence spaces, and makes function application and composition continuous. Continuity and homomorphism both reduce to preservation of filter convergence and enable differentiation operators to be defined on subspaces of the function spaces in agreement with differentiation operators on normed spaces and hence on Euclidean and Hilbert spaces. Rigorous definitions for continuity and differentiation are below:

[1] Bourbaki, N. *Topologie, generale*, Actualites Sci. Ind. **858** (1940), **916** (1942), **1029** (1947), **1045** (1948), **1084** (1949), Paris.

[2] Stepney, S. et. al. "A Framework for Heterotic Computing", 8th Int. Workshop on Quantum Physics and Logic, 2012.

[3] Patten, D. *Problems in the Theory of Convergence Spaces*, Ph.D. Dissertation, Syracuse University, 2014.

[4] Binz, E. "Recent results in the functional analytic investigations of convergence spaces" in Proc. 3rd Prague TopologicalSymp., 1971.

# An Agent Based Model on Schelling's Segregation by Languages

*Dieudonne Ouedraogo,Azadeh Ahkamiraad*

**Abstract**

The purpose of this study is to show how people are distributed and settled down in cities based on their location in neighborhoods of their preferred language. Our Project could be defined based on the conventional Schelling's Segregation (1969). Briefly, an agent in a small region may change its location when his/her majority neighbors are not speaking the same language as the agent's language. We have developed a new dynamic model for our study which bear semblance to Schelling's model. Our model uses dynamic classes containing iterators and constructors. Also, we modeled and analyzed the interaction between agents and the macroscopic effect of the interaction; which leads to segregation based on agent's attribute (language). We make various assumptions and simulate scenarios to analyze the dynamic behavior and the evolution of emergence. To depict the level of clustering, we are using the mean distance of every node in the cluster

## I. INTRODUCTION

The Schelling's Segregation model was a breakthrough and a significant milestone in expressing social dynamics and how it is possible for the society to emerge (or show emergent behavior). More work on the Schelling's Segregation may lead us into potential inclusion of all forms of visible segregation (Caste, Gender, Color)[1]. However, we must be aware of the fact that Visible Segregation is not the only factor interpreting human behavior. The choice of company is determined by multiple factors, not necessarily discriminative, but emotional and cognitive as well. This may be regraded as the vision for an all inclusive model of human behavior and ABM is a perfect tool for the classical model.

## II. EXPLANATION OF SYSTEM MODELED

Assumptions for Classical Schelling's Model[2][3], are as follows:

1) The agents in Schelling's problem are based two at- tributes(race) are spread over 2D space.

2) The randomly chosen agent looks around in his/her neighborhood

(a) If the agent finds other agents of the same type, he/she stays

(b) Else, the agent moves to a random location in 2-D space

If we keep the threshold of neighbors speaking the same Language low enough, we would find localized segregation taking place. But a High Threshold (Indicating a low level of tolerance) would cause the agents to not settle down in forms of colonies and keep moving around.

## III. EXPLANATION OF EXPERIMENTS AND MATHEMATICAL ANALYSIS

Thomas Schelling's model, does provide a apt explanation about how society decides to segregate and organize themselves. It also explains how people are inherently logical (but not necessarily rational at every instant), about the decision to aggregate with the same kind, or a community consisting of similar agents. However, as we see from social dynamics evolving over time, we find that time has caused societies to integrate much more than they used to. For instance, looking at the population distribution of the united states and finding the proportion of every ethnicity in the states using the histogram based on immigrant population alone, over the past century, (Figure 1) , we may conclude that the population is growing. The rise in population is
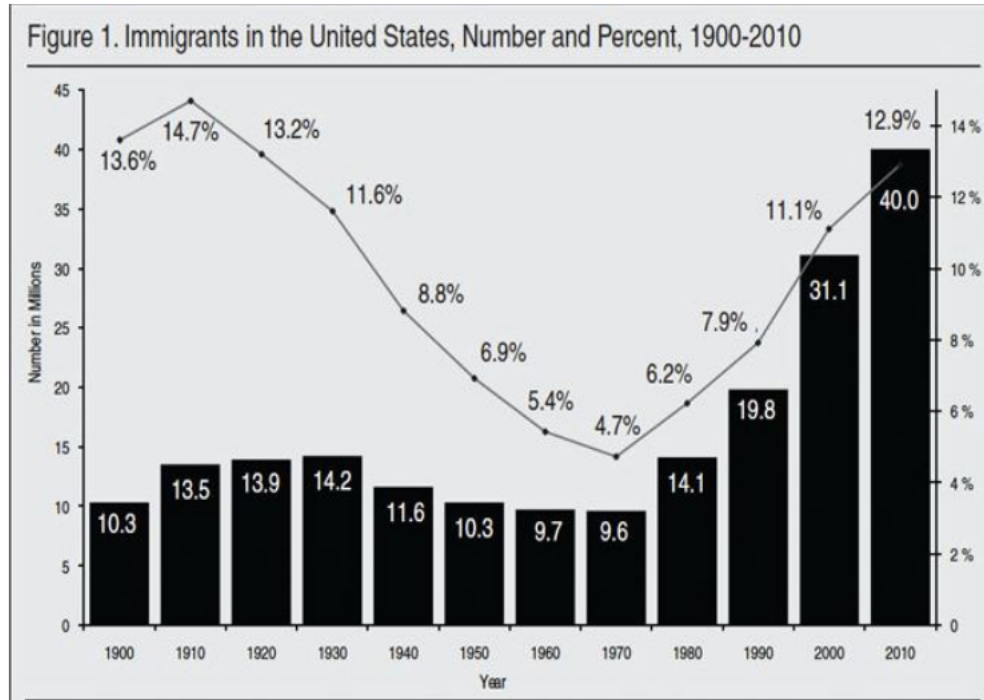
Figure 1: Figure1:Source:Decenial Censuses, 1990-2000, American Community Survery, 2010

40 million in 2010 as opposed to 9.6 million in 1970 would indicate that individuals don't just segregate, but selectively integrate as well. This behavior may be termed as localized segregation (or selective aggregation).

Based on this observation, we have introduced two key behaviors in our model. These assumptions are listed as follows:

1) Adaptation: People speaking the same language mostly come from the same culture. They tend to segregate, but similarity is not the sole metric of segregation. We have also included the 'majority in numbers' as a factor influencing segregation & adaptation, in our model. Con- sidering the population distribution of the United States as our modeling environment, we have the following four group that co-exist. We have simultaneously tried to observe their behavior:

(a) The Majority, i.e., English speaking popula- tion has the highest adaptation (or influence rate and population). Thus, over the years, the agents for other languages will try to adapt to the behavior exhibited by the English speaking population and hence convert to the same. They are represented by the color 'Blue'

(b) The second most common language observed is 'Spanish', represented by the color 'Red' and has the second highest adaptation rate and population

(c) We then define one of the minorities as 'Chinese', which is represented by the color 'Green' and have the third highest Adaptation/influence

(d) Lastly, we consider the Arabic population having the least adaptation/influence/population and are represented by the color 'Yellow'. All adaptation rates are kept constant.

2) Reproduction based on Adaptation: We also assume that the majority in numbers is a direct contributor to the overall increase in the number of agents, over the years. The number of agents belonging to a particular language community increase based on their adaptation rate and current population. The behavior of these two factors may be expressed with the following equation:

179

**T = a · english+b · spanish+c · chinese+d · arabic**

Where, a = 1.7, b = 1.1, c = 1.1, d = 1.05 are the fixed adaptation rates for selected languages, and english,spanish,chinese,arabic are the populations of the languages over the years. The spawn of agents and their locations over the map is implemented using organized randomness, so as to obtain different results while maintaining the order of majority. The addition of new agents each year is the expressed by the following equation: New Agents = Adaptation · Current Agents This would reflect the exponential increase of the majority language over time. Mean Distance: In order to analyze the level of segregation and sparsity, we have aggregated the positions and calculated the mean distance between individuals of one language. Individuals that are evenly spread out throughout the environment and have a centroid coinciding the approximate center of the environment, would not produce sufficient variation in the quantity. Whereas, multiple asymmetric clusters skewed across directions would cause the quantity to variate drastically.

3) Mathematical Analysis: We have simulated using two and four agents.

(a) Position of each agent is defined by x and y all in range of [0,1].

(b) Neighborhood is defined as the nearest element to the agent in term of Euclidean distance.

(c) We have also increased the neighborhood size to 40,60

(d) At the beginning of the process all agents are uniformly distributed in the square region [0,1]*[0,1]

(e) Agents move according to the following:

i) Go through all agents possible

ii) Pick a random position in in the square

iii) If the agent is pleased with the position

(tolerance level is met), agent relocates. Else, goto (i)


## IV. RESULTS

The results shown below emulate our projected outcomes appropriately. Change in Threshold and total number of agents, organized randomness imply that the mean position metric may be used to quantify the level of clustering. We observe the following:


### A. Moderate Threshold

The first case included taking agents in a similar manner to the population distribution of United States. The pop- ulation is as orthodox as it is accepting. The details for first two years is shown in figures 2a 2b

We further observed, the minorities are forced to aggre- gate. This maybe due to the inherent and moderately re- pulsive behavior. It shows induced rationality and process of intermingling. Figure 3


### B. Low Threshold

We adjusted a threshold of 0.2 and obtained the following results for the first two years ref. Figures 4a and 4b below.


### C. High Threshold

A threshold of 0.9 would produce high segregation, it also represents the case of a very orthodox population of individuals. The first two years have been shown as follows

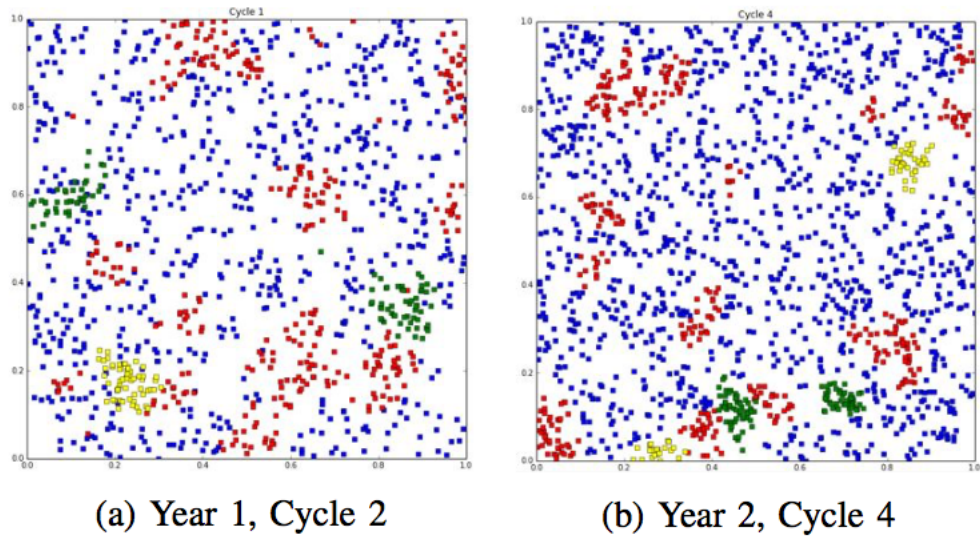(a) Year 1, Cycle 2          (b) Year 2, Cycle 4

Figure 2: Figure2:Trend of localized segregation/selective aggregation of the population for a moderate threshold
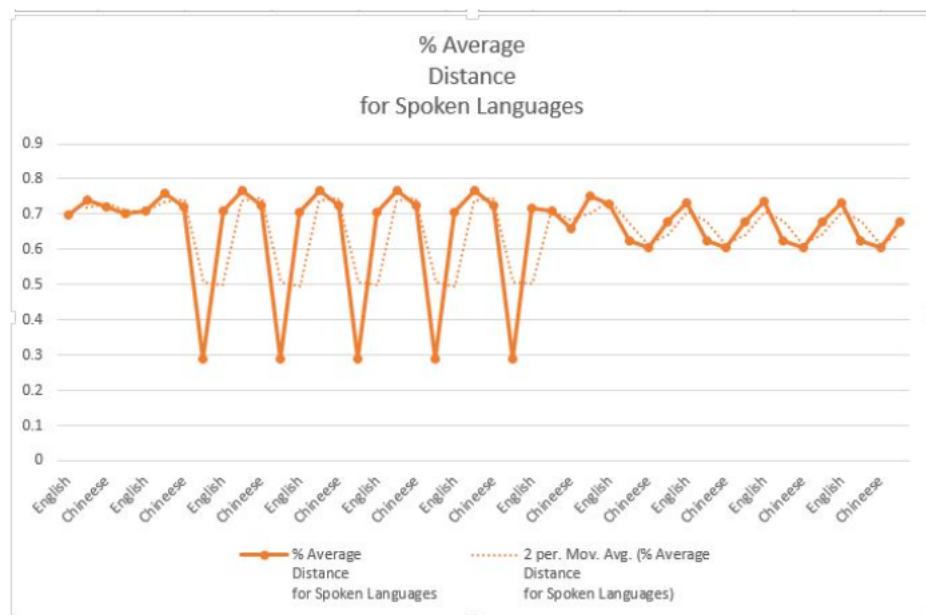


Figure 3: It can be observed from multiple simulations that the mean distance for each language is affected by, the population size, population increase and the tolerance level

(a) Year 1, Cycle 0

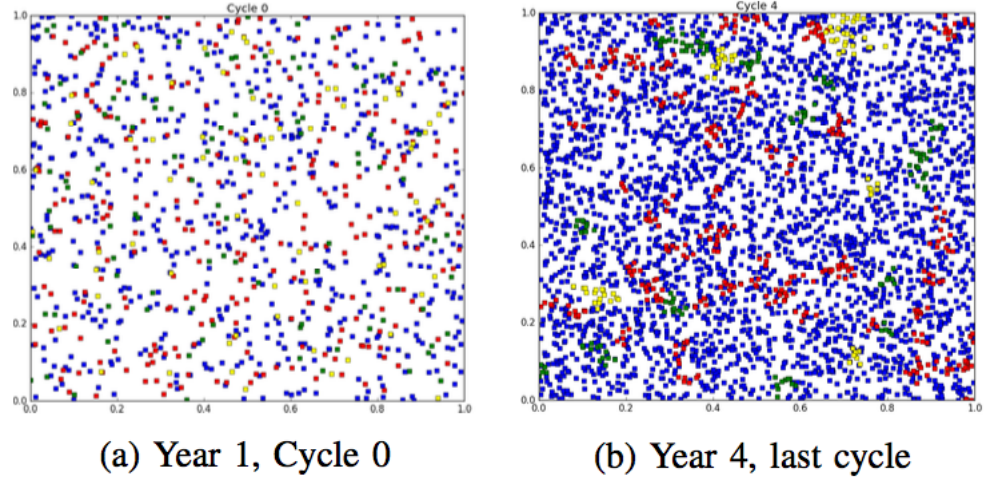(b) Year 4, last cycle

Figure 4: Figure 4:Mild Segregation is observed - coincides with the classical model(3)



(a) Year 1, all cycles



(b) Year 2, all cycles

Figure 5: Representation of a very orthodox population distribu- tion - Ideal Case
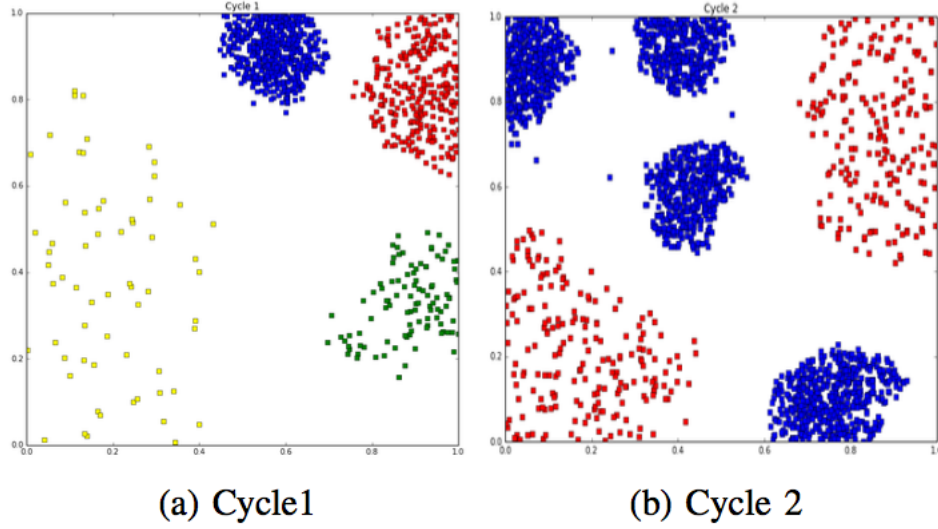
(a) Cycle1        (b) Cycle 2

Figure 6: Increasing the segregation neighborhood also leads to implementation of organized randomness. By implementing 20 neighbors and tolerance level of 0.5

### D. Variable Neighborhood

For every movement in the Moore neighborhood, we performed prior computations with 8 agents, just as Schelling had taken for the matrix based simulation, where every agent has eight neighbors. However, in a real society, this is not the case. Individuals look beyond just their own neighborhood, i.e. a broader community in order to make a choice to stay or leave, hence we assumed a larger neighborhood size (40 neighbors) and tried simulating the classic Schelling's model with this revised assumption. ref. Figure 6a, 6b, 7

We can clearly see from the cases above that even at a low tolerance level, if the neighborhood is expanded, high segregation is produced. As we raise the size of the neighborhood, the system partitions clearly into clusters of same language, segregation is very obvious. As the size of the neighborhood diminishes, segregation falls apart, the system tends to integrate, all other values remain constant.

If one holds values of the neighborhood constant or in small variation in certain interval, the perception is that even small tolerance level (and or population increase) leads to amplified segregation, but our set of simulations shows that a main contributor to segregation at the macroscopic level is the choice of the size of the neighbors, Schelling's assumption was based on values of neighborhoods confined in a region where segregation moves alongside tolerance. Therefore we have proved that Schelling's model was wrong.

There is a relation between size of neighborhood tolerance and segregation and our future goal is to quantify it.

### V. DISCUSSION OF RESULTS & SCOPE OF FUTURE STUDY

Our results show the different case scenarios. We can also implement language learning process by integrating a list of languages that gets appended every time an agent moves to a neighborhood where there is a majority of people speaking a different language. We would also like to express the mathematical relationship between neighborhood size and the tolerance level. We would like to use our model to explain migration between regions by including other variables such as employment, crime rate, climate etc

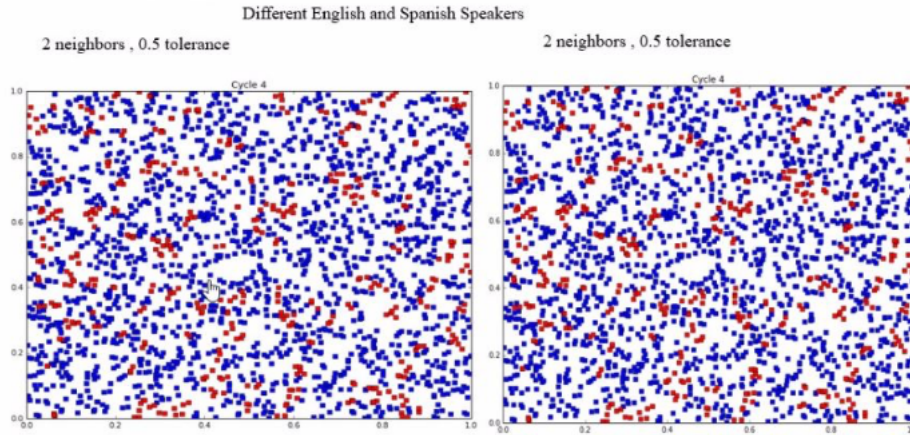Figure 7: Figure 7 :Simulation for tolerance = 0.5, neighborhood = 2

**REFERENCES**

[1] C. T. Schelling, "Dynamic models of segregation," Journal of Mathematical Sociology, vol. 1, pp. 143–186, 1971.

[2] H. Sayama, Modeling Complex systems, 9. 2015, vol. 53, pp. 1689–1699, ISBN: 9788578110796.

[3] C. T. Schelling, "Models of segregation," The American Economic Review, vol. 59, no. 2, pp. 488–493, 1969.

# Educating Tomorrow's Leaders in the Science of Connections: The USMA Network Science Minor

Jonathan W. Roginski

[1] Network Science Center
United States Military Academy, West Point, NY
jonathan.roginski@usma.edu

## Abstract

Educating the next generation of leaders using network thinking prepares them for life and leadership in a complex world which the relationships between entities are often as importanteven more importantthan the entities themselves. The whole is often more than the sum of its parts. Weve embraced preparation of our leaders for this complexity at the United States Military Academy. We exist to educate, train, and inspire our graduates in preparation lead an Army with the mission to fight and win our Nations wars. So, understanding the environment in which we operate is imperative because what works in one place will not in another and what worked yesterday may not work today. This challenge is not unique to the military. Indeed, the imperative of operational understanding and awareness extends well beyond the military operational art to leadership across any endeavor from understanding the sales and service climate in business to the intricacies of interpersonal dynamics in team science to the connections between elements of knowledge that define education. Our interdisciplinary Network Science Minor students grow classes in fundamentals and theory to experiences in modeling and application, culminating in a sponsored integrative capstone experience.

Network Science is the science of our times. The United States Military Academy's exploratory experience that is MA394: The Fundamentals of Network Science and moreover, the entirety of the Network Science Minor progression, will better enable you to thrive in a world characterized by its connections. Though we often hear that our world is complex, through these courses you will understand what that actually *means* and how to deal with that complexity, effectively modeling situations that resist traditional, reductionist approaches. More than that, you will refine your ability to communicate your insights both orally and in writing, critical skills so many lack. Network Science is inherently interdisciplinary; there is something in our minor program for everyone. We are excited about the opportunity to facilitate your network-enabled vision of the world.

*You'll walk away from our minor with tools and intuition into our connected world that makes you invaluable when leading and advising in a complex environment.* You'll be proficient in the use of network science technology and algorithms so your analysis will scale to large and interesting problems. Not only will you communicate confidently, but your competence will inspire confidence with whom you speak.

*Our expectations are simple* First, enjoy yourself. You are embarking upon this network science journey voluntarily and we know it. Your Professor wants nothing more than for you to succeed. Second, come prepared. Generally, these are math courses and math concepts. the topics build on each other. Finally, collaborate with each other and with your Professors. Network Science is a team sport and the whole of our program is greater than the sum of its parts.

To demonstrate the breadth of our program's research, the following are poster titles from recent network science students:

Bipartite Analysis of Modern Military Alliances (Bhoo)
Actions Speak Louder Than Words (Bullying, Conger)
Modeling Diffusion of Health Policy (Dula)
A Game of Math (Gleason)
Trading Behavior in MLB, ′14&′17 (Helton)
Spread of VE Through Social Network (Isham)
Importance of Amino Acids in Protein Composition (Milanesa)
Who's Who of MA103 and MA104 Citations (Murdock)
Pandora's Choice (Newton)
Army Soccer Passing Network (Nolasco)
Harry Potter Network Analysis (Palumbo)
Bad Guy in MCU (Quillen)
Interaction Between Genes and Diseases (Chloe Smith)
Network Science & Fake News (Sodgerel)
HITS and Political Science (Suba)
Predicting NFL Wins with PageRank (Sullenger)
The Structure of Influential People (Walas) International Nuclear Trade (Warnock)

We hope you enjoy this program and get everything you can out of your Professors, each other, and this experience!

## Acknowledgments

**Network Science Center**
West Point

# Promoting Obesity Prevention Policy: how counties can diffuse innovative preventative health policies

Kathryn Dula

Networks Science Center
United States Military Academy
kathryn.dula@usma.edu

## Abstract

This research generated 17 multiplex network experiments. Each multiplex network was generated by aggregating 5 layers from the absolute difference between county attributes. Each layer was created by adjoining counties with similar enough attribute values. The attributes included: percent of obese adults, percent poor, percent uninsured, number of physicians per 1,000 people in the population, and Popin Food Desert Scores. The layers were combined into a single multiplex network by weighting each layer's adjacency matrix using a pairwise comparison technique. The matrices were then summed up into a holistic adjacency matrix that represents the relative influence each county has over another. Two counties with and intra-influence score above .5 are adjoined in the multiplex network. Of the resulting 17 experiments, nine were too dense to warrant a diffusion analysis. Of the eight remaining, 5 diffusion trials were run, varying the number and identity of the initial active nodes as well as the proportion of active neighbors required to be activated in turn. The diffusion trials revealed that cliques in each experiment acted as blocking communities, requiring a large number of initially active nodes for complete diffusion to occur. In practical terms this means that lobby groups will not be able to concentrate efforts on a small subset of counties in West Virginia and achieve state-wide preventative health policy adoption. Lobby campaigns will have to span the state in order to be effective.

While research has improved how doctors handle weight loss plans with their patients on an individual level, we have yet to address environmental factors on a larger scale. The only entities that have the ability to change the resources and infrastructure in our shared environment is our government. Thus the United States is in need of effective preventative health policy that will shape our environments to foster healthier lifestyles to reverse this trend. This study will outline the initial conditions required to make the adoption of obesity prevention policies popular among county level governments in West Virginia. The influential relationship between counties can be modeled in a network where the nodes are counties and the weighted edges represent the extent of the influence they have on other counties' policy decisions. How strongly we value another government's decisions is a sum of many different county attributes. Each county attribute that impacts policy diffusion is represented in a separate layer of the network. Once each layer is generated by connecting counties whose attribute values are within a specific threshold of one another, we collapsed the layers into a single aggregate layer. The edges in this single layer represent the comprehensive influence that counties have over each other based upon the extent of attributes they share in common. Then using a linear threshold diffusion model, updated with a bounded confidence opinion dynamics model we can simulate how many counties will adopt an innovative health policy given an initial set of counties that have adopted a similar policy. This study conducts a sensitivity analysis using a nearly orthogonal latin hypercube experimental design to efficiently sample the solution space with 17 different experiments. Furthermore it analyzes the impact initial sets have on the complete diffusion of an innovative health policy across the network of counties in West Virginia.

# Agent-Based Modeling as a Method for Testing Counter-Terrorism Strategy

Cadet Jordan Isham
Department of Mathematical Sciences
United States Military Academy
West Point, New York

---

Advisor: LTC Jonathan Roginski
Department of Mathematical Sciences
United States Military Academy
West Point, New York

CONTACT:    CDT Jordan Isham, Department of Mathematical Sciences, USMA,
West Point, NY, 10997 Tel: (802)-238-3007; e-mail: Jordan.Isham@usma.edu

LTC Jonathan Roginski, Department of Mathematical Sciences, USMA, West Point, NY,
10997 Tel: (845-938-0267); e-mail: Jonathan.Roginski@usma.edu

## Abstract

The implementation of any single counter-terrorism strategy is both costly and time intensive. Thus, an appropriate application of network modeling and simulation, specifically through agent-based modeling (ABM), serves a useful low-cost tool to test various strategic approaches prior to their implementation. In this regard, the nature of agent-based simulation allows for a robust examination of the dynamic and complex system of terrorism. ABM's are structured to represent a collection of agents that interact and change over time, where the interactions between agents give rise to macro-systems and highlight targeted characteristics of a given network.[1] ABM can be used to analyze social networks, even dark networks, using a stochastic approach to adjust micro-entities and demonstrate significant macro-changes in the network. This study explores the potential for network modeling through network analysis and simulation analysis of a targeted-leader approach. Namely, results come from a topographical network analysis and simulation analysis based on the removal of the top two leaders at various points in the network's development.[2] As depicted through simulation analysis, the slope of the log-log distribution of nodes versus degrees helps to explain the structure of the terrorist's organization regarding its fragmentation and robustness. Under a preferential attachment model, targeting key leaders early in the stages of development fragments the network more than when the network has grown in size. However, the fragmented components prove much more robust in the later stages of fragmentation than in the early stages.

---

[1] Nianogo, Roch A., and Onyebuchi A. Arah. "Agent-Based Modeling of Noncommunicable Diseases: A Systematic Review." *American Journal of Public Health* 105, no. 3 (2015). doi:10.2105/ajph.2014.302426.

[2] J.P Keller, "Dismantling Terrorist Networks: Evaluating Strategic Options Using Agent-Based Modeling."

## 1. Introduction

The limitations of most social network analysis' stem from 1) confounding variables, 2) an unclear direction or degree of relationships, and/or 3) the difficulty in modeling findings to large populations as a result of lack of data. The potential confounding variables of homophily and the shared environments both suggest that agent-based modeling may serve as a viable alternative to social network analysis. This study addresses those limitations using agent-based modeling. Stochastic computer simulations of agents assess macro-level patterns that arise from micro-level interactions. Agent-based approaches are proven to be useful when individual agent behavior is complex, thus allowing for learning and adaptation, feedback loops, and/or reciprocity. It is also useful in heterogeneous environments which can influence agent behavior and interaction. Lastly, ABM is useful when agents are not fixed in space or time; and when inter-agent interactions are complex, non-linear, and influence agent behavior. The limitations to agent-based models cause for a rather parsimonious model with overarching assumptions, a qualitative rather than a quantitative interpretation, and difficulty of validating the final model with very limited data.[3]

The specific ABM software used is NetLogo, which is the most widely used open source ABM. NetLogo was designed to be easily readable and allows for many variations of ABM. Essentially, NetLogo is a coding program that serves as a computational model for users to take certain inputs, manipulate them in algorithmic ways, and generate outputs. ABM relies on simple rules to observe agent interactions. For example, when an ant leaves their nest, walks in a trail to look for food, and then returns to their nest, NetLogo models the process. In order to tell the ant to keep looking for food if it has not found food yet, the code would say: if not carrying-food? [look-for-food]. Thus, NetLogo is both a modeling language and an integrated environment designed to make ABM easy to build. To build on the above ant model, the code could include move-towards-nest and wander-for-food functions. With the ability to incorporate randomness into the mode, NetLogo proves useful in modeling real-world interactions.[4]

One particular study that serves as the basis for this research, *Dismantling terrorist networks: Evaluating strategic options using agent-based modeling by* Jared Keller, Kevin Desouza, and Yuan Lin, describes the application of ABM for testing four counter-terrorism strategies to include 1) leader-focused, grassroots, geographic, and random intervention strategies. Current debate surrounds the effectiveness of leader-targeted killings. While some scholars argue that "targeting the group's leadership reduces its operational capability by eliminating most highly skilled members and diverting time towards the protection of its leaders, others posit that "leadership decapitation is likely to increase the number of willing recruits for terrorist groups to exploit."[5] The Noordin terrorist plot, which attempted to kill the Indonesian president, suggested that a targeted killing of Noordin may temporarily slow the network's capabilities but that newfound networks may split from the original network with the same or stronger ideal. In this regard, many variables complicate the strategy of leader decapitation to include the timing of an attack in respect to the group's life cycle, the characteristics of the organization (violent, clandestine, and/ or values based), the process of leadership turnover, the

[3] Abdulrahman M. El-Sayed,, Peter Scarborough, Lars Seemann, and Sandro Galea. "Social Network Analysis and Agent-Based Modeling in Social Epidemiology".

[4] U. Wilensky (2005). NetLogo Preferential Attachment model. http://ccl.northwestern.edu/netlogo/models/PreferentialAttachment. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL.

[5] Bryan C. Price "Targeting Top Terrorists: How Leadership Decapitation Contributes to Counterterrorism." *International Security*, vol. 36, no. 4, 2012, pp. 9–46., doi:10.1162/isec_a_00075.

size and number of connections per node, etc. Unlike the leader-focused strategy, the grassroots-focused strategy eliminates operatives in network "who are responsible for conducting activities and carrying out tasks".[6]

## 2. Assumptions
Although the interactions of this model do not necessarily reflect real-world data, social science qualitative analysis drives the agent-types and interactions. Additionally, the foundation of the model stems from the random-graph model, extending to the small-world model, and then to the Barabasi-Albert model containing a preferential attachment mechanism which means that highly-connected nodes continue to develop relationships faster than those with fewer relationships (the rich get richer). Through simulation analysis, we are able to create a log-log distribution over time with the degree on the $x$ axis and the number of nodes on the $y$ axis. Using the equation suggested by Sean Gourley, the probability of a terrorist attack is modeled by $P(x) = Cx^{-a}$ where $a$ represents the slope of the log-log distribution of $\frac{log(number\ of\ nodes)}{log(degree\ of\ network)}$.[7] Similar to Gourley, our findings suggest that the organizational structure can be interpreted through the slope of the log-log degree *vs.* node distribution. A high $a$ indicates a fragmented network with more components of lower robustness while a lower $a$ indicates more a robust network with fewer components. With the creation of a mode for intervention to kill the top two leaders based on their degree, results surround a comparison of $a$ before any intervention and after a targeted intervention. In order to accept the results from any ABM (due to a shortage of data), we must assume that an edge implies a close social bond between two people, and thus a local cluster acts as a close-knit micro-community of family and friends.

## 3. Model Formulation
One ABM under the Threat Anticipation Program focuses on the effect of social grievances in the Middle East through the simulation of a complex socioeconomic system with thousands of agents each having several variables and interactions within multiple types of social networks.[8] As a continuation of previous models, this study explores how effectively ABM can reflect the spread of radicalization in a social network. After removing the nodes with the most connections in the network, much like leader-targeted killings, we re-analyze the dynamics and summary statistics. As the model progresses, agents will continue to create connections with other agents in the system based on the number of connections and the lifespan of the node in the system.

The use of a known data set helps to validate the analysis of the network because the number nodes, edges, and the nature of the network is largely understood. The nature of dark networks implies that connections between individuals are often unclear. In the Noordin network, 561 edges exist between 145 nodes, where clear relationships are indeed visualized but the network only shown known edges between actors involved in the network. In terrorist networks dark networks have a much higher quantity of uncertainty in edges, clusters, closeness, and edge weight. A more accurate depiction will test probability among network connections, adjusting the transitivity to look at instances with triadic closure. Although running code does not pinpoint triadic closure, using probability helps explore the potentiality of relationships within the network. The dark network can be explored by simulating the various types of edges

---

[6] J.P Keller, "Dismantling Terrorist Networks: Evaluating Strategic Options Using Agent-Based Modeling."
[7] Sean Gourley. "The Mathematics of War." *TED Talk*.
[8] JP Keller, "Dismantling Terrorist Networks: Evaluating Strategic Options Using Agent-Based Modeling," Page 6.

in regard to strong and weak ties by adding/ removing edges and vertices based on random probabilities of their relationships depicted through a series of 50 runs for each probability and centrality measure. In his report on dark networks, Mark Granovetter points out that "if it were not for weak ties, clusters would not be connected at all."[9] He goes on to say that strong ties serve as a much greater motivation to be a source of support in times of uncertainty, such as in terrorist networks. In the Noordin network, a false identification of connected clusters would complicate the entire counter-terrorism operation. Since a dark network may have connections unseen by topology and connections that may not truly contain substance, simulation runs with probabilities help identify and aggregate the most influential actors. A topology network also fails to depict various measures of centrality—they only cluster individuals and highlight the degree of their connections (whether true connections or not). While betweenness assumes that an actor has power over any two other actors when in the shortest path, closeness measures the ability of information to get to an actor, thus demonstrating centrality. In the Noordin network, the degree of an actor is not necessarily as important as their closeness or betweenness when it comes to passing along secret information or influencing a particular terrorist attack. Particularly important to a terrorist network; eigenvector centrality assumes that connections to central actors are weighted heavier than connections to peripheral actors. The identification of the top two leaders accounts for uncertainty by running three probabilities to explore various network dynamics based on adding edges, removing edges, and removing vertices for four centrality measures including degree, closeness, betweenness, and eigenvector centrality.

Thus, in order to collect summary statistics, we run a network analysis in Gephi both before the removal of the top two nodes (where there are 145 nodes) and after the removal of the top two nodes. Then, transitioning to the simulation in NetLogo, we collect statistics in the form of degree distribution and with a log-log plot of the degree vs. number of nodes when there are 145 nodes and calculate the organization structure modeled by $a$ over time. It's helpful to compare the type of findings possible with network analysis of a known network and the findings of simulation analysis.

### 4. Results

Through a topographical network analysis of the Noordin network, we identified the top two leaders of the network by running 36 aggregated, averaged, and graphed networks (Appendix A). The top two strongest tied actors in each data set identified Noordin and Azhari Husin, who appear most often at the top of the network, with many standard deviations above the other actors (Appendix B). When using probabilities to add and remove edges, these two individuals continuously have the highest betweenness, closeness, degree, and eigenvector centrality. At the point of the network where there are 145 known nodes and 561 known edges and before the removal of these two nodes, the summary statistics include an average degree of 15.994. After the removal of these two nodes, the summary statistics include an average degree of 12.95. However, the topographical network analysis does not allow us to analyze the network as it changes over *time*. Thus, simulation results in NetLogo also surround a comparison of $a$ before and after a targeted intervention. A targeted killing of the top two nodes resulted in the $a$ of -2.66 when the network was at 143 nodes, and -2.65 after the network was at 1,000 nodes. As a network continues to grow, the change in $a$ decreases over time after the removal of the top two nodes. Thus, the removal of key leaders early in the stages of a network's development

---

[9]Sean Everton, TRACKING, DESTABILIZING AND DISRUPTING DARK NETWORKS WITH SOCIAL NETWORKS ANALYSIS, Naval Postgraduate School.

fragments the network more than when the network has grown in size. Ultimately, this study shows how the organizational structure of a dark terrorist network can be understood through both a topographical network analysis and simulation analysis. Specifically, the slope of the log-log distribution represents the organizational structure of the network which can then be used to predict the probability of attack. Where Gourley's study analyzes organizational structure based on the size and frequency of attack, this study uses the number and degree of nodes in the Noordin network.

## 5. **Future Directions**

There are several limitations of the model that we must keep in mind, such as the simplistic representation of agent interactions. In order to inspire confidence and explanatory power in the model, verification (relevance to the model's concept), validation (relevance to real-world phenomena), and replication are all useful. Particularly because ABM are stochastic in nature, replication serves even more use, and multiple runs are needed to confirm the accuracy of the model.  Another limitation too, is the researcher's ability to code with NetLogo. Although many limitations exist, there are also many benefits to using ABM. Particularly, the randomness of the network structure allowed for a practical application to a social network where the social bonds and degree of nodes are not always the same. Thus, through replicating the simulation across trials, the mean results were collected. Additionally, the analysis through time-stamps effectively model the propensity for growth and allow for the concept of preferential attachment to serve relevant. Finally, the flexibility of ABM proves incredibly useful in allowing for sensitivity analysis, and further manipulation depending on the research.

References

1. El-Sayed, Abdulrahman M., Peter Scarborough, Lars Seemann, and Sandro Galea. "Social network analysis and agent-based modeling in social epidemiology." *Epidemiologic Perspectives & Innovations* 9, no. 1 (2012): 1.

2. Gourley, Sean. "The Mathematics of War." *TED Talk*.

3. Keller, Jared P., et al. "Dismantling Terrorist Networks: Evaluating Strategic Options Using Agent-Based Modeling." *Technological Forecasting and Social Change*, vol. 77, no. 7, 2010, pp. 1014–1036., doi:10.1016/j.techfore.2010.02.007.

4. Nianogo, Roch A., and Onyebuchi A. Arah. "Agent-Based Modeling of Noncommunicable Diseases: A Systematic Review." *American Journal of Public Health* 105, no. 3 (2015). doi:10.2105/ajph.2014.302426.

5. Price, Bryan C. "Targeting Top Terrorists: How Leadership Decapitation Contributes to Counterterrorism." *International Security*, vol. 36, no. 4, 2012, pp. 9–46., doi:10.1162/isec_a_00075.

6. Reinares , Fernando. "Differential Association Explaining Jihadi Radicalization in Spain." *Combating Terrorism Center at West Point* , June 2017 https://ctc.usma.edu/posts/differential-association-explaining-jihadi-radicalization-in-spain-a-quantitative-study

7. Wilensky, Uri, and William Rand. *Introduction to agent-based modeling: modeling natural, social, and engineered complex systems with NetLogo*. Cambridge, MA: The MIT Press, 2015.

8. Wilensky, U. (2005). NetLogo Preferential Attachment model. http://ccl.northwestern.edu/netlogo/models/PreferentialAttachment. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL.

9. Wilensky, U. (1999). NetLogo. http://ccl.northwestern.edu/netlogo/. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL.

Appendix A

| Probability | .25 | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Network Action | Add Edges | | | | Remove Edges | | | | Remove Vertices | | | |
| Centrality Measure | B/w | Close | Degree | Eigen | B/w | Close | Degree | Eigen | B/w | Close | Degree | Eigen |
| Node #1 | 51 | 51 | 51 | 51 | 51 | 51 | 51 | - | 34 | 21 | 2 | 2 |
| Node #2 | 76 | 76 | 76 | 76 | 76 | 76 | 76 | - | 36 | 20 | 34 | 21 |
| Node #3 | 91 | 91 | 91 | 91 | 91 | 91 | 91 | - | 35 | 2 | 21 | 23 |
| Probability | .50 | | | | | | | | | | | |
| Network Action | Add Edges | | | | Remove Edges | | | | Remove Vertices | | | |
| Centrality Measure | B/w | Close | Degree | Eigen | B/w | Close | Degree | Eigen | B/w | Close | Degree | Eigen |
| Node #1 | 51 | 51 | 51 | 51 | 51 | 51 | 51 | - | 23 | 14 | 16 | 2 |
| Node #2 | 76 | 76 | 76 | 76 | 76 | 76 | 76 | - | 24 | 16 | 14 | 1 |
| Node #3 | 91 | 91 | 91 | 91 | 11 | 91 | 91 | - | 27 | 12 | 2 | 15 |
| Probability | .75 | | | | | | | | | | | |
| Action | Add Edges | | | | Remove Edges | | | | Remove Vertices | | | |
| Centrality Measure | B/w | Close | Degree | Eigen | B/w | Close | Degree | Eigen | B/w | Close | Degree | Eigen |
| Node #1 | 51 | 51 | 51 | 51 | 51 | 51 | 51 | - | 18 | 1 | 18 | 1 |
| Node #2 | 76 | 76 | 76 | 76 | 76 | 76 | 76 | - | 22 | 8 | 11 | 11 |
| Node #3 | 91 | 91 | 91 | 91 | 91 | 91 | 91 | - | 20 | 12 | 1 | 9 |

**Table 1.** After running python code on the adjacency matrix with p=.25, p=.50, and p=.75 (desensitize our results) where n=50 for each centrality measure (betweenness, closeness, degree, and eigenvector) and also for each network change (add edges, remove edge, and remove node), the result is 36 aggregated, averaged, and graphed networks where the top 3 strongest probability ties are depicted. The highlighted sets are graphed in Appendix B below.
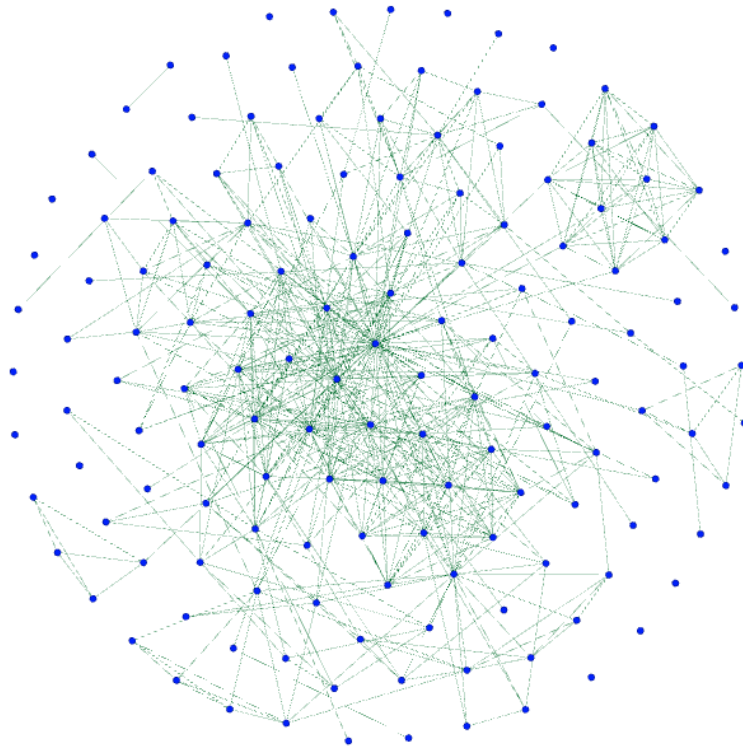
**Graphs 1-4.** Below includes the graphs for the highlighted data sets in Appendix A above. The x axis represents the node, labeled 0-144 (adjusted to match the 145 actors in Noordin network). The y-axis represents the probability of the occurrence for the centrality measure labeled in the graph's respective title. Thus, Noordin and Husin's standard deviations are visualized.
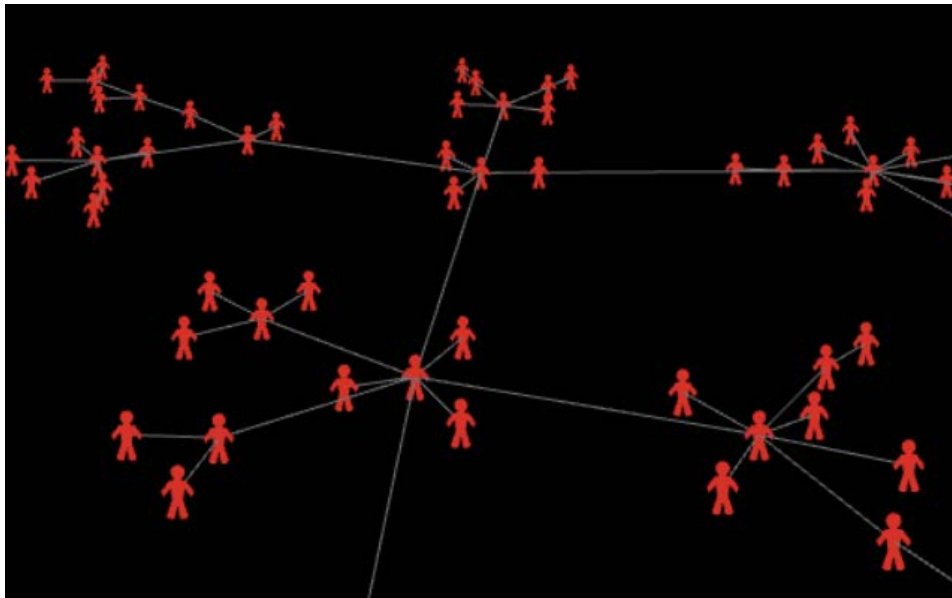


Add Edge if p=.5 for Betwenness



Remove Edge if p=.25 for Degree

**Figures 1 and 2.** Below includes the network visualizations for the regular Noordin network on the right and an adjusted network on the left. As the probabilities for adding edges increases, the modularity decreases. For example, the network on the left has been adjusted by removing edges who have less than a .75 probability, thus a higher modularity. The visualizations help understand various features and types of relationships within the network.

Original Noordin Top



NetLogo Snapshot

# Twitter Reveals: Using Twitter Analytics to Predict Public Protests

Mohsen Bahrami[1,3], Yasin Findik[2,3], Burcin Bozkaya[1,3], Selim Balcisoy[2,3]

[1]School of Management, Sabanci University, Istanbul, Turkey
[2]Faculty of Engineering and Natural Science, Sabanci University, Istanbul, Turkey
[3]Behavioral Analytics and Visualization Lab, Sabanci University, Istanbul, Turkey

**Abstract** Citizens participate in mass demonstrations to express themselves and exercise their democratic rights. However, because of the large number of participants, protests may lead to violence and destruction, and hence can be costly. Thus, it is important to predict such demonstrations in advance to safeguard against such damages. Twitter has been used as a tool by protestors for planning, organizing, and announcing many of the recent protests worldwide, such as those that led to the Arab Spring, Britain riots, and those against Mr. Trump after the presidential election in the U.S. In this paper, we aim to predict protests by means of machine learning algorithms. In particular, we consider the case of protests against the then-president-elect Mr. Trump after the results of the presidential election were announced in November 2016. Our findings confirm that Twitter can be used as a powerful tool for predicting future protests with an average prediction accuracy of over 75% (up to 100%). We further validate our model by predicting the protests after President Trump's travel ban executive order. An important contribution of our study is the inclusion of event-specific features for prediction purposes, which helps to achieve high levels of accuracy.

## 1. Introduction

By means of protests, people express their interests, needs, approval or disapproval of a particular situation, and try to bring a better future to their society. Even though majorities of protests have been reported to be peaceful [1], because of the large number of participants in demonstrations, protests may lead to violence and destruction. Thus, it is important to predict such demonstrations in advance to protect against such damages and reduce their expected costs. During the last decade, social media and especially Twitter have been widely used as organization, information and mobilization tools for protests. Twitter is arguably the most prominent and popular microblogging and social networking website worldwide. It is shown that 85% of tweets are related to news [2] and citizen journalists use Twitter to disseminate information and news in real-time, faster than the famous news agencies, and even correct if there is misinformation. Previous studies have shown Twitter's key role in recent protests such as those that led to the Arab Spring [3], and London riots [4]. Other studies in the literature include those that perform tweet content and sentiment analysis during events and protests [2], investigate Twitter usage at the times of crises and natural disasters [5], describe, model, and interpret the user networks and relationships between social networks as well as social movements [6, 7]. Here our focus is on studies that propose models for predicting protests using social media and especially Twitter. We observe two main approaches to prediction: the first based on the properties of online user networks, interactions on social media, and activity cascades [8, 9], and the second approach based on the features driven from aggregated user posts and their contents [3]. In our study, we take the second approach for prediction, for which we also find a significant number of studies in the literature. In this paper, we aim to predict public protests by means of machine learning algorithms using the features extracted from the collected tweets calling for demonstration. In particular, we

consider the case of the U.S. post-election protests in 2016, and make predictions for each of the fifty states in the U.S. Furthermore, we test our model by predicting the protests held in the U.S. airports after president Trump's executive order #13769, and attempt to explain that both protest cases were likely motivated by the same reasons.

The main contribution of our work is two-fold: similar to previous studies, the results of our study confirm that Twitter analytics can help to predict major protests successfully, and we show that the inclusion of regional features driven by reasons and motivators for protests can help to make more accurate predictions.

## 2. Methodology

In order to predict a protest using Twitter data with an online and real-time system, we propose an algorithm which is comprised of five steps. Figure 1 shows the flowchart of the proposed algorithm. Since we collected the data after the protests happened, we already have information without performing two initial steps. As a result, in this study we investigate steps three to five and refer to previous research to explain the methods used for performing the first two steps.



Figure 1. Protest prediction system flowchart

The first step is to search for early signals of a protest in tweet contents [10]. Whenever the amount of signals in a fixed time interval pass a certain threshold, the system starts the second step. The second step consists of finding main trending hashtags and reasons of a protest from the collected tweets or user posts [11]. One approach is to identify the main trending hashtags and then analyze the content of tweets, which use those hashtags [12]. Users may create several hashtags during the early hours of tweeting about the protest, but users will coordinate on using a few or even one of them [13, 14]. For steps 3-5 of the algorithm, we collect and prepare the Twitter data for analysis, then we extract input features / predictors of the model based on our dataset, and finally we select and apply a classifier for prediction. In what follows, we describe the general approach and then implement our methodology at each step for the specific study case. We make predictions at the state level, which means for fifty U.S. states, we make fifty predictions per day using a rolling origin prediction approach.

### 2.1 Data Collection & Cleansing

We collected the tweets, which were tweeted, from November 9 to 15, 2016, containing #NotMyPresident and from January 27 to 31, 2017, containing #muslimban and #travelban. The data we collected amount to more than 4 million tweets and retweets. In order to match tweets with states, only those tweets with known geo-location tags are useful. Furthermore, while hashtags hint to the topic of each tweet, they do not guarantee that the collected tweets under a given hashtag all have relevant contents. Invariably a noticeable percentage of the collected tweets have irrelevant contents, and thus they should be removed from the dataset. Removing tweets with irrelevant content reduces the data size to about 0.47 million tweets and retweets, which we consider as the population in this study Since retweeting can be viewed as an agreement and recommendation mechanism [15], we consider each retweet as a single tweet, independent from the original tweet. For simplicity, we thus refer to "tweets and retweets" only as "tweets" in the sequel.

## 2.2 Feature Extraction

### 2.2.1 Extracting features from the collected tweets

Most of the recent research indicates and agrees on some generic and common tweeting behavior of users before and during most protest and crisis events. These tweeting signals are shown to be common in most of the protest cases studied, independent from where, when, and why they are happening. We use variations of these features as inputs in our prediction model. The first group of features we extract from the data set pertains to the count of tweets in the context of a protest. As time gets closer to the protest day, an increase in tweeting activity is noticeable [16]. In our work, we calculate average tweet count per hour (e.g. average hourly tweeting pace). Tweeting pace can be used instead of tweet count per day when the time intervals of the collected data sets for model training are different. The second group of features include time, date, and place mentions, which we extracted from the collected tweet. These tweets are very important since they explicitly announce the protest times and venues [10, 17].Table 1 shows some examples of different ways of time, date, and place mentions. By searching for the mentions in the collected tweets, we derive mention counts for each state (place) and each particular date (time).

| Different Time Date Mention Examples |
| --- |
| Protest in my city planned for tomorrow evening #NotMyPresident |
| #NotMyPresident Anti-Trump rally planned for Downtown Indianapolis on Saturday |
| #LosAngeles high schools will be walking out November 14th 9:15AM. All protests will lead to City Hall. #protest #notmypresident |

Afterwards, using the Bag-of-Words method, we calculate the average number of violent words per tweet in a daily basis at the state level. Then we perform a sentiment analysis on tweets to find the percentage, and polarity (i.e. strength of the negativity or positivity) of those tweets with negative sentiment, broken down by data and state. The negativity percentage and its polarity are indicators of users' sentiment about the event which is considered as a predictor of protest, and the average number of violent words per tweet is considered as an indicator of the probability of protest resulting in violence [18].

### 2.2.2. Event Specific Features

There are several reasons which motivate protests and social unrest such as absence of democracy and freedom, political corruption, social injustice, police violence, and unstable economic conditions resulting in a high unemployment rate, poverty, and rising food prices. We argue that twitting behavioral features are not able to fully represent regional motivators of the protest; moreover Twitter penetration is not geographically comparable for one to safely consider metrics like tweet counts. We propose that the inclusion of the features related to the reasons of each protest would help improve the prediction results. To the best of our knowledge, we are the first to include regional features as motivators for protests used as independent variables in a prediction model. To construct an improved model, hence we attempt to consider the features which motivate citizens to participate in a particular protest event. In order to find these kinds of features, the reasons and goals of protests should be investigated carefully. Furthermore, we propose that in case of similarity between protests, these features could potentially be used as common predictors for the same kinds of protests. In the case of post-election protests, which were against Mr. Trump's presidency, we contend that there may be a negative association between the probability of protests happening and Mr. Trump's percentage of votes in each particular state. That is, we argue that those who voted for him are unlikely to participate

in a protest against him. Hence, we include Mr. Trump's percentage of votes as a predictor variable in our model. Figure 3 shows the states where Mr. Trump got the majority of the state votes, as shown in red [19].
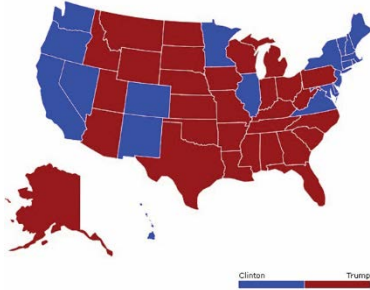


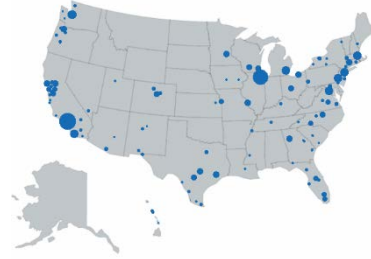Figure 3. State-level vote distribution      Figure 4. County-based lead vote distribution map

Total votes in each state is a collection of votes of all counties inside the state; thus, it is important to consider the population effect of large cities and counties as well. To elucidate, in some states such as Texas, even though the majority of the votes is in favor of Mr. Trump, in four counties, Mrs. Clinton has led with a very large margin of vote count compared to her competitors. Investigating the first two days of post-election protests, the corresponding protest maps, and the lead vote map (see Figure 4), we argue that there is also an association between protesting states and the large lead vote sizes in favor of Mrs. Clinton. In order to add large county lead votes to our model, we define a binary variable. The variable is equal to 1 for state $i$ if there is a county in that state with leading vote more than a given threshold in favor of Mrs. Clinton, and 0 otherwise. Figure 4 shows the significant lead votes of Mrs. Clinton in each state. The circle size is proportional to the number of votes received by Mrs. Clinton in each county.
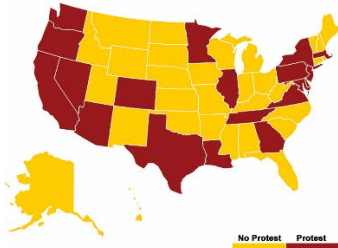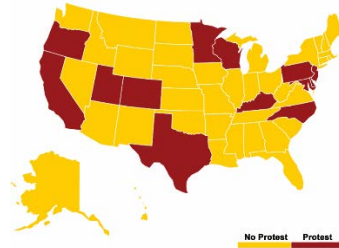


Figure 5. November 9, 2016 protest map      Figure 6. November 10, 2016 protest map

Figures 5 and 6 show the protest maps at the state level in the first two days after the election, on November, 9-10, 2016 [20]. One can visually detect an association between variables visualized in Figures 3 and 4, and the protest maps in Figures 5 and 6.

### 2.2.3 Feature Selection

To select the most significant features for our analysis among all extracted ones, before prediction for each next day, we use Wrapper feature selection method [21] with four classifiers namely: C4.5, Naïve Bayes, Logistic Regression (LR), and Support Vector Machines (SVM). We utilize a bi-directional Best First search on the set of extracted features using the training data with a 10-fold cross validation. We then chose the features which are significant in at least 40% of the folds on average of those four classifiers results. Afterwards, to avoid multi-collinearity, we calculate the pairwise correlation for all features and if there is a high correlation between two features, we remove the one with less significance.

Finally, after the feature selection procedure, we use the following seven variables as the predictors of our model:

- Daily average tweeting pace for each state
- Number of date, time, and place mentions
- Average polarity of daily tweets with negative sentiment
- Average number of violent words per tweet for each state
- Total daily tweet count with negative sentiment in data set as an indicator of overall Twitter users' activity related to the protest topic
- Percentage of President Trump's vote in each state
- A binary variable indicating if the county lead vote size is larger than a threshold in each state for Mrs. Clinton. After some parameter tuning, we used 150,000 votes for Mrs. Clinton as the threshold for the binary indicator variable as large lead vote states.

## 2.3 Classifier Selection

As we defined previously, the model's output is a binary variable that indicates whether a protest in each state on a particular day is predicted or not. We have tested various classifiers from different families of classifiers: C4.5, Naïve Bayes, Logistic Regression, and SVM. We obtained the overall best results among these using Logistic Regression with Logit; hence, in the next section we report only the results produced by this classifier. Logistic Regression's output is the probability of protest in each state on a particular date. We decide to convert the probability value into a binary value. To do this, we need a cutoff value to convert this probability into a binary classification. We allow the algorithm to determine the best cutoff value by utilizing the Single Rule decision tree (OneR method) based on the training data. If the predicted probability is above the cutoff, the model result is converted to 1 meaning that there will be a protest, otherwise it is converted to 0.

## 3. Analysis Results

The case we investigate is a chain of protests after the 2016 presidential elections, and unlike some previous studies, the time interval between the announcement of election results and the demonstrations is very short (less than 18 hours). Moreover, there are no tweets recorded prior to the first day of protests; to be able to predict such quickly formed protests, we train our model with each previous day's tweet counts and other predictors mentioned above, for predicting the next day's protests. In our predictions, we handle date/time mentions differently in that we need to consider more days of history simply because the date/time mention might be referring to *any* future date. Each day represents 50 data points, one for each state. After each day, 50 new points are added to the training data. For example, at the end of the 4th day, 250 data points will be available to train the model for 5th day's predictions. Our algorithm represents a progressive rolling prediction model, where we fit and tune the model parameters with all available data collected from previous days.

We present the performance results (prediction accuracies) of the two models we trained in Tables 2 and 3. In both tables, we report three predictor performance measures as True Positive Rate (TPR), True Negative Rate (TNR), and Overall accuracy. The first model (Model 1 presented in Table 2) excludes the event-specific predictor variables, whereas the second model (Model 2 presented in Table 3) includes all predictor variables.

Table 2.Post-election protest prediction results using features extracted from Twitter (Model 1)

| Predicted date | Nov 11 | Nov 12 | Nov 13 | Nov 14 | Nov 15 | Nov 16 |
|---|---|---|---|---|---|---|
| TPR | 50% | 46.66% | 72.72% | 71.43% | 100% | 100% |
| TNR | 97.5% | 97.14% | 82.05% | 88.37% | 87.5% | 88% |
| Overall accuracy | 88% | 82% | 80% | 86% | 88% | 88% |

We find that the event-specific variables contribute to the accuracy of predictions in the early days of the protests, which is quite useful because of the limited number of twitter records during that period. As the protests saturate, both models seem to perform similarly for predicting the protests.

Table 3. Post-election protests prediction results using features extracted from Twitter and event specific features (Model 2)

| Predicted date | Nov 11 | Nov 12 | Nov 13 | Nov 14 | Nov 15 | Nov 16 |
|---|---|---|---|---|---|---|
| TPR | 80% | 53.33% | 63.33% | 71.43% | 100% | 100% |
| TNR | 57.5% | 88.57% | 84.61% | 88.37% | 89.58% | 90% |
| Overall accuracy | 62% | 78% | 80% | 86% | 90% | 90% |

Figure 7 shows the ROC curve for models 1 and 2 based on the total prediction results of November 11 to 16. Area under ROC curves (AUC) show that both models are successful in prediction of protests; moreover, the results confirm that inclusion of event specific features can increase the performance of the predictive model. (AUC of model 2 is larger.)
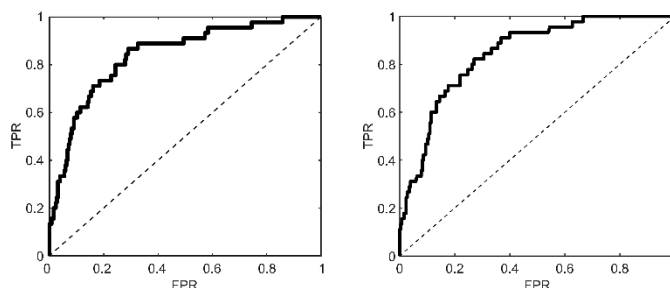


Figure 7. ROC curves of Model 1 (left) and Model 2 (right) based on overall prediction results for post-election protests, AUC1 = 84.07%, AUC2= 84.42%

Next, we use the same two models trained by the post-election protests data of November 10-16, 2016 to predict the protests of January 29, 2017 against the presidential executive order #13769. The results of the two models are shown in Tables 4 and 5. We were unable to collect enough tweet records for January 29 and 30, and hence we are unable to present prediction results for January 30 and 31.

| Table 4. Muslim-ban protest prediction accuracy indices using Model 1 | | Table 5. Muslim-ban protest prediction accuracy indices using Model 2 | |
|---|---|---|---|
| Predicted date | Jan 29, 2017 | Predicted date | Jan 29, 2017 |
| TPR | 73.33% | TPR | 66.66% |
| TNR | 65.71% | TNR | 71.42% |
| Overall accuracy | 68% | Overall accuracy | 78.66% |

As the results indicate, the models trained with the post-election protests could predict the protests against the executive order #13769 with a reasonably high accuracy. Again the second model including event specific features, generally speaking, produces better results and higher overall accuracy. Figure 8 shows the ROC curve for the models 1 and 2 classifiers confirming better performance of model 2 with larger AUC than model 1. Based on these results, we find that the post-election and post-travel-

ban activity on Twitter along with the most relevant hashtags provide a significantly high predictive power. Furthermore, we observe that using event specific features further help to increase the predictive power of our model and hence the resulting prediction accuracy values.
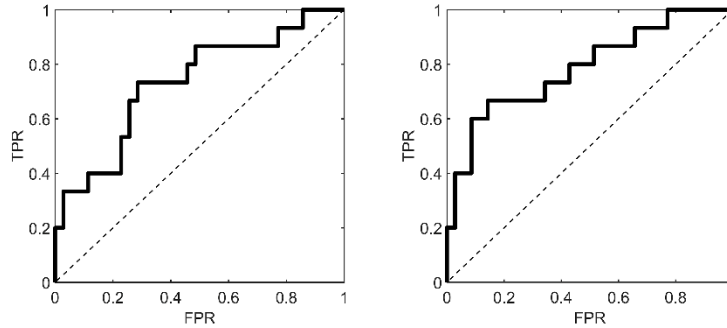


Figure 8. ROC curves of Model 1 (left) and Model 2 (right) used for prediction of protests against the executive order #13769, AUC1 = 73.33%, AUC2 = 78.66%

The results also suggest that there is a similarity between the two cases considered, because the models, trained with voting statistics and post-election tweets, can predict both the post-election protests and the protest against executive order #13769, with high accuracy. Our approach of using the Wrapper method also has further shown that the voting statistics and the post-election specific features always end up being in the set of most effective predictors. Consequently, one may speculate that both cases of protests were motivated by a set of common reasons, yet we believe there might be other features and reasons as well involved in amplifying the chain of protests under study.

## 4. Conclusion

In this study, we present a classification-based prediction model for predicting mass protests based on Twitter data and tweeting behavior of users. We additionally consider and include event-specific features (e.g. voting statistics) in our model to improve the prediction performance. Both models (with and without event-specific features) predict the post-election protests with high accuracy and true-positive, true-negative rates. In order to show the robustness of our proposed models and the similarity of two cases of mass protests that are potentially related, we have trained our models with one case of protests (i.e. post-election events) and tested them with a second case of protests (i.e. the Muslim-ban events due to presidential executive order #13769). We find satisfying levels of prediction accuracy with the latter case of protests. Our analysis and case study results suggest that the features extracted from Twitter and the models we have developed can potentially be used for protest prediction in the countries where the protests are considered as legal and Twitter is actively used. Furthermore, we note that adding event specific features improves the prediction performance for both cases of protests we studied, which suggests that such features may serve as common motivators or reasons for the related protest events. Finally, our study results emphasize the key role of social media, especially Twitter, in recent protests as an organization and information tool, and that the Twitter has the power of revealing answers for many research questions.

## References

1. McLeod, D. M., & Hertog, J. K. Social control, social change and the mass media's role in the regulation of protest groups. Mass media, social control, and social change: A macrosocial perspective, 1999; 305-330.
2. Kwak, H., Lee, C., Park, H., & Moon, S. What is Twitter, a social network or a news media? In Proceedings of the 19th ACM international conference on World Wide Web 2010; pp. 591-600.

3. Steinert-Threlkeld, Z. C., Mocanu, D., Vespignani, A., & Fowler, J. (2015). Online social networks and offline protest. EPJ Data Science, 4(1), 19.

4. Cheong, M., Ray, S., & Green, D. interpreting the 2011 London riots from twitter metadata. In Intelligent Systems Design and Applications (ISDA), 12th International Conference on IEEE 2012; pp. 915-920.

5. Brown, S. Twitter Usage in Times of Crisis. Open Access Journals for School Teachers in Indonesia, 2011; 29.

6. Gupta, A., Joshi, A., & Kumaraguru, P. Identifying and characterizing user communities on twitter during crisis events. In Proceedings of the 2012 workshop on Data-driven user behavioral modelling and mining from social media 2012; pp. 23-26. ACM.

7. Bajpai, K., & Jaiswal, A. A framework for analyzing collective action events on Twitter. In Proceedings of the 8th International ISCRAM Conference 2011.

8. Cadena, J., Korkmaz, G., Kuhlman, C. J., Marathe, A., Ramakrishnan, N., & Vullikanti, A. (2015). Forecasting social unrest using activity cascades. PloS one, 10(6), e0128879.

9. González-Bailón, S., Borge-Holthoefer, J., Rivero, A., & Moreno, Y. (2011). The dynamics of protest recruitment through an online network. Scientific reports, 1, 197.

10. Xu, J., Lu, T. C., Compton, R., & Allen, D. Civil unrest prediction: A tumblr-based exploration. In International Conference on Social Computing, Behavioral-Cultural Modeling, and Prediction. Springer International Publishing 2014; pp. 403-411.

11. Ramakrishnan, N., Butler, P., Muthiah, S., Self, N., Khandpur, R., Saraf, P., ... & Kuhlman, C. (2014, August). 'Beating the news' with EMBERS: forecasting civil unrest using open source indicators. In Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 1799-1808). ACM.

12. Romero, D. M., Meeder, B., & Kleinberg, J. (2011, March). 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 (pp. 695-704). ACM.

13. Korolov, R., Lu, D., Wang, J., Zhou, G., Bonial, C., Voss, C., ... & Ji, H. (2016, August). On predicting social unrest using social media. In Advances in Social Networks Analysis and Mining (ASONAM), 2016 IEEE/ACM International Conference on (pp. 89-95). IEEE.

14. Lehmann, J., Gonçalves, B., Ramasco, J. J., & Cattuto, C. (2012, April). Dynamical classes of collective attention in twitter. In Proceedings of the 21st international conference on World Wide Web (pp. 251-260). ACM.

15. Starbird, K., & Palen, L. Pass it on? : Retweeting in mass emergency. International Community on Information Systems for Crisis Response and Management 2010; pp. 1-10.

16. Lotan, G., Graeff, E., Ananny, M., Gaffney, D., & Pearce, I. (2011). The Arab Spring| the revolutions were tweeted: Information flows during the 2011 Tunisian and Egyptian revolutions. International journal of communication, 5, 31.

17. Compton, R., Lee, C., Lu, T. C., De Silva, L., & Macy, M. Detecting future social unrest in unprocessed twitter data: "emerging phenomena and big data". In Intelligence and Security Informatics (ISI), IEEE International Conference On 2013; pp. 56-60.

18. Kallus, N. Predicting crowd behavior with big public data. In Proceedings of the 23rd ACM International Conference on World Wide Web 2014; pp. 625-630.

19. https://www.nytimes.com/elections/results/president

20. https://www.nytimes.com/interactive/2016/11/12/us/elections/photographs-from-anti-trump-protests.html

21. Kohavi, R., & John, G. H. (1997). Wrappers for feature subset selection. Artificial intelligence, 97(1-2), 273-324.

# The Effect of Network Structural Properties on the Performance of Crawling Approaches

Katchaguy Areekijseree and Sucheta Soundarajan

Department of EECS, Syracuse University, NY. {*kareekij, susounda*}@syr.edu

## 1 Introduction

The literature on network crawling contains a large number of algorithms, but while one method may work very well on a certain type of network, it may perform poorly on other networks. The goal of our work is to understand why certain methods perform well on certain networks, and to give guidance to users in selecting an appropriate network crawling technique. We explore the relationship between the performance of various network crawling algorithms and network structure: in particular, we investigate those network structural properties that govern the ability of a crawler being able to move between regions of a graph, and show that these features have a strong effect on the performance of crawling methods.



Figure 1: Results of the controlled experiments.

## 2 Network Crawling vs. Structural Properties

We perform an analysis of seven of the most popular network crawling techniques, evaluated on controlled synthetic and real networks. Let $G = (V, E)$ be an unobserved, undirected network. A starting node and number of allowable queries $b$ are given. In each step, the crawler makes query on an observed-but-unqueried node, and we assume that all neighbors are returned in response to a query. The process stops after $b$ queries. The ouput is a sample graph $S = (V', E')$, which contains all nodes/edges observed. We consider *node coverage* as our evaluation criteria. Three structural properties of interest are *Modularity*, *average degree* and *average community size*. We categorize the seven crawling algorithms into three groups based on its perfromance, as shown in Table 1. **G1 - Node Importance-based methods**: maximum observed degree [1] and observed PageRank, **G2 - Random Walk** and **G3 - Graph Traversal methods**: BFS, DFS, Snowball and Random.

Table 1: Categorization and summary of the performances of crawling approaches.

| Property | G1 | G2 | G3 |
|---|---|---|---|
| Modularity | Excellent performance when modularity is low. | | Stable |
| Average community size | Strong performance when communities are large if modularity is high. Community size does not matter if modularity is low | Stable | |
| Average degree | Strong performance when average degree is extremely low ($\leq$10) even if modularity is high. Otherwise, stable | | Performance improvement when average degree increases. |
| **Best Method** | **MOD** | **RW** | **BFS** |

## 3 Results and Conclusion

Our experiments show several results. 1) the performance of G1 methods significantly improves as the size of the community increases, but decreases as modularity increases: these methods have difficulty transitioning to new communities. 2) The performance of G2 method is unaffected by these properties. These methods move freely between different communities. (3) The performance of G3 methods improves as average degree increases, and is unaffected by community-based properties. E.g., Figure 1 shows an example of results on synthetic networks with varying average degree and community sizes. These results give insight into different types of crawling algorithms, and help guide a user in selecting a method.

## References

[1] K. Avrachenkov, P. Basu, G. Neglia, B. Ribeiro, and D. Towsley. Pay few, influence most: Online myopic network covering. In *Computer Communications Workshops*, 2014.
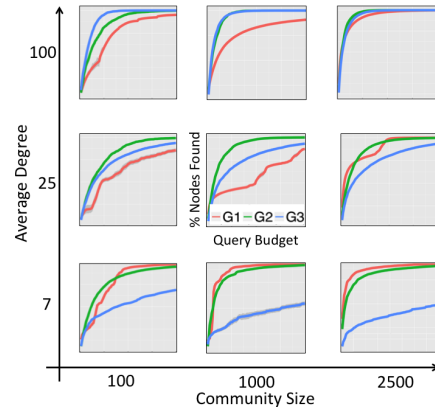
# LSTM BASED SOIL MOISTURE PREDICTION

Shivanand Venkanna Sheshappanavar[1], Chilukuri K. Mohan[1] and David G. Chandler[2]

[1] Department of Electrical Eng. and Computer Science
Syracuse University, New York
ssheshap/mohan@syr.edu
[2] Department of Civil and Environmental Eng.
Syracuse University, New York
dgchandl@syr.edu

## Abstract

Soil moisture content is an important variable that has a considerable impact on agricultural processes and practical weather-related concerns such as flooding and drought. We address the problem of predicting soil moisture by applying recurrent neural networks that use *Long Short-Term Memory (LSTM)* models. The success of our approach is evaluated using a dataset obtained from ground-based sensor infrastructure networks. Feature reduction using a mutual information approach is shown to be more effective than feature extraction using principal component analysis.

## 1   Introduction

Soil moisture has a colossal impact on several hydrological processes including infiltration, evapotranspiration and subsurface flow. Measurement and prediction of soil moisture help us obtain deeper insights into the localized dynamics of critical ecological processes. Accurate prediction of soil moisture allows the quantification of drought conditions and the prediction of flash floods caused by precipitation run-off. Economic consequences of accurate soil moisture prediction are significant, assisting improvements in crop productivity and agricultural management practices, and permitting precise control over the root zone environment, leads to healthier crops and higher yields. Weather forecasts can be improved, since high soil moisture results in high evaporation, increasing the likelihood of moisture convergence. In addition, monitoring soil moisture provides us with better understanding of how water, energy and carbon are exchanged between land and air.

Neural networks such as multilayer perceptrons have been useful in prediction of stream-flow based on snow accumulation, along with recommends application of Principal Component Analysis (PCA) [1]. More recently, deep networks have been used for prediction, providing greater flexibility in mapping diverse, complex functions. The risk of over-fitting the data (resulting in poor generalization) can be mitigated using regularization techniques or by reducing the feature space dimensionality. Researchers have recently used Deep Belief Networks (DBN) and other techniques for feature learning or extraction [2].

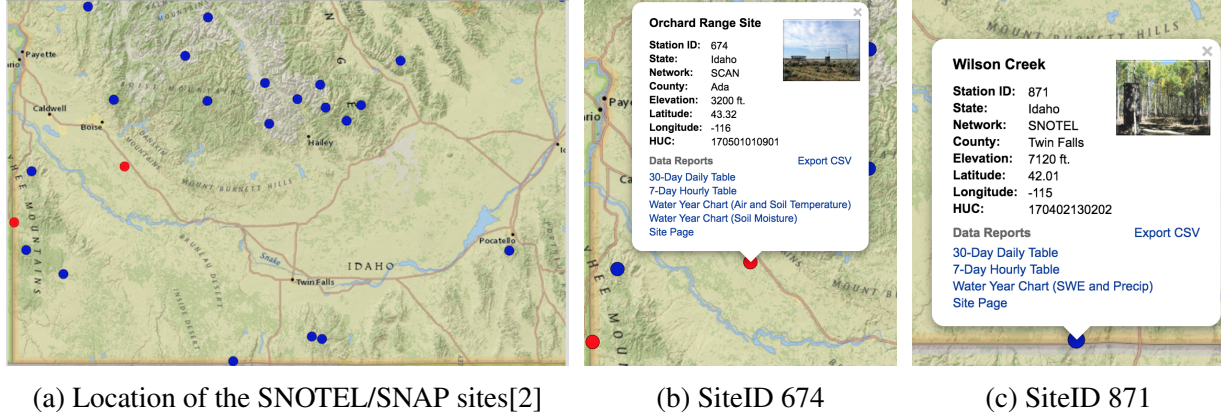(a) Location of the SNOTEL/SNAP sites[2]   (b) SiteID 674   (c) SiteID 871

Figure 1: Two sites selected for this experiment[2]

Recurrent Neural Networks with Long Short-Term Memory (LSTM) rank among the state-of-the-art networks for predicting future values of a time series, with potential application to hydrology [3]. Deep Feed-forward Neural Networks have also been used on datasets obtained using Visible Infrared Imaging Radiometer Suite (VIIRS) from cropland China [4].

This paper presents the application of LSTM model for soil moisture prediction, using datasets collected from ground using Soil Climate Analysis Network (SCAN) and Snow Telemetry (SNOTEL) networks, unlike prior works which focused on SMAP datasets. We apply this approach to data from different soil climate networks, and evaluate two feature extraction methods on this data. We were able to find the most significant features for soil moisture prediction, and shown the effectiveness of our approach compared to feature extraction methods previously used in hydrology.

## 2   Data Collection

Two important sources of data used in this paper are from networks SCAN and SNO-TEL, managed and operated by National Resources Conservation Services (NRCS) and Natural Water and Climate Center (NWCC) [5]. Figure 1 (a) shows a group of sites (SCAN-red, SNOTEL-blue) in Idaho. We selected 2 sites one from each network, Orchard Range Site (SCAN Station ID 674) and Wilson Creek (SNOTEL Station ID 871) located in Ada and Twin Falls county respectively of Idaho State as in figure 1 (b) and (c). Datasets with predefined features can be downloaded from [6] by specifying start and end dates, selecting hourly options for an annual period. A

Table 1: List of features with tag names.

| **All Feature Names (total - 21)** |
| --- |
| (PREC.I-1, PREC.I-2) Precipitation Accumulation |
| (TOBS.I-1) Air Temperature Observed |
| (STO.I-1:-2,8,20) Soil Temperature Observed |
| (SAL.I-1:-2,8,20) Salinity |
| (RDC.I-1:-2,8,20) Real Dielectric Constant |
| (BATT.I-1,BATT.I-2,BATT.I-3) Battery |
| (WDIRV.H-1) Wind Direction Average |
| (WSPDX.H-1) Wind Speed Maximum |
| (WSPDV.H-1) Wind Speed Average |
| (SRADV.H-1) Solar Radiation Average |
| (RHUMV.H-1) Relative Humidity Average |
| (RHUMX.H-1) Relative Humidity Maximum |

| **Soil Moisture Targets (total - 3)** |
| --- |
| (SMS.I-1:-2,8,20) Soil Moisture Percent |

dataset spanning for 5 years (01/01/2012 to 12/31/2016) has been collected, treated for missing values, divided into training (4 years, 2012-2015) and testing (3 months, Jan-March 2016) sets.

Distribution of NAs

PREC.i.1..in

Distribution of NAs

SMS.I.1..8..pct.slit

Distribution of NAs

STO.I.1..20..degC

Distribution of NAs

WXPDX.H.1..mph

Distribution of NAs

SRADV.H.1..watt.

Distribution of NAs

SAL.I.1..2..gram.

Figure 2: Plot of missing values (red vertical lines) and actual values (blue dots) of few features.

Figure 3: General statistics of missing values

Table 1 lists the features for station ID 674 (other sites may have less or more). For some of the features values are missing at a given site so we have analyzed the data for missingness by plotting the dataset for NA distributions as shown in figure 2 (due to space constraints and for proper visibility including only few original features in figure 2). The red vertical lines represent the missing values. General statistics of missing values shown in figure 3 where first row represents number of data points with no missing

| | SMS.I.1..8..pct...silt. | STO.I.1..8..degC. | SAL.I.1..8..gram. | RDC.I.1..8..unit. | SMS.I.1..2..pct...silt. | SAL.I.1..2..gram. | RDC.I.1..2..unit. | TOBS.I.1..degC. | BATT.I.1..volt. | BATT.I.2..volt. | STO.I.1..2..degC. | SMS.I.1..20..pct...silt. | STO.I.1..20..degC. | SAL.I.1..20..gram. | RDC.I.1..20..unit. | WDIRV.H.1..degr. | WSPDX.H.1..mph. | WSPDV.H.1..mph. | SRADV.H.1..watt. | RHUMV.H.1..pct. | RHUMX.H.1..pct. | PREC.I.1..in. | PREC.I.2..in. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 39742 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 |
| 432 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 |
| 2901 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 |
| 110 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 2 |
| 60 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 5 |
| 43 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 6 |
| 5 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 7 |
| 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 7 |
| 291 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 10 |
| 199 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 11 |
| 30 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 12 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 13 |
| 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 13 |
| 6 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 14 |
| 2 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 14 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 18 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 21 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 22 |
| | 43 | 43 | 43 | 43 | 48 | 48 | 48 | 65 | 65 | 65 | 246 | 540 | 540 | 540 | 540 | 558 | 558 | 558 | 558 | 558 | 558 | 609 | 3074 | 9948 |

values for any of the features. A '1' represents that a value exists for the feature in the given row and a '0' represents missing value. The corresponding number on left most column represents the number of such rows in the dataset. For example only one row is missing all values except for precipitation accumulation 'PREC.I.1..in.' as represented by penultimate row of the table. We removed features 10 percentage or more missing values for example dropped battery 'BATT.I-3' feature resulting in the feature count to 19 as shown in table 3. Mutual information between the

features and the visualization of features as in figure 2 suggested omitting soil moisture percent, two salinity attributes, one real dielectric constant and wind direction average features. We plotted density of the observed data and imputed data as shown in figure 4 to conclude that the observed values neither are identical nor similar to the imputed values (Missing Not at Random).
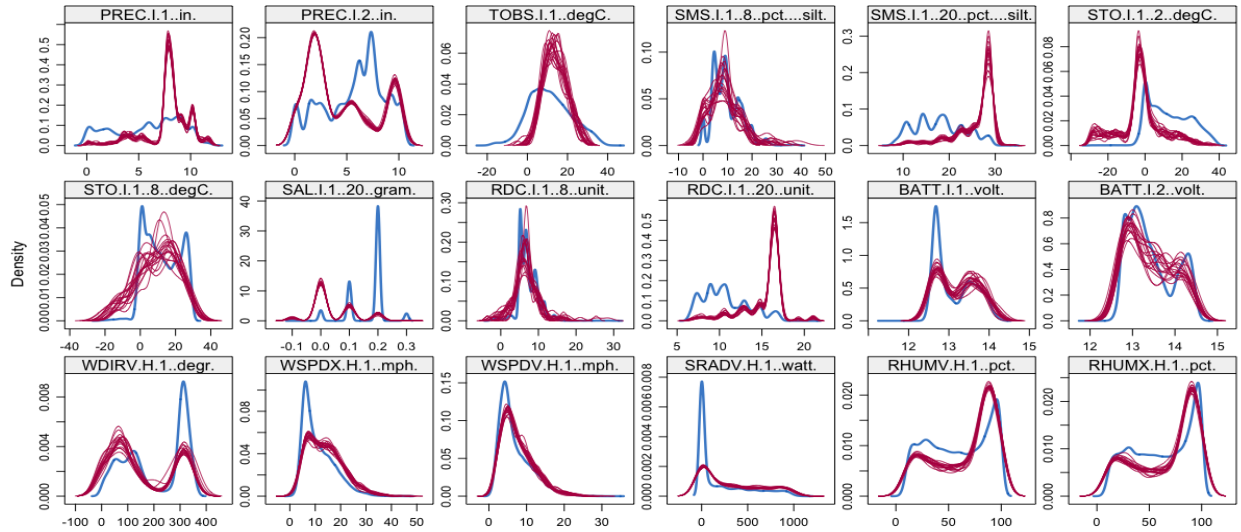


Figure 4: The observed values(blue) are neither identical nor similar to the imputed values(red) which is why the missing of values is concluded as Missing Not at Random (MNAR).

We used two methods for filling in the missing values and reduce feature dimensionality:

- Dataset I: Exponential moving average imputation using R package imputeTS [7] by taking the exponential moving average (with a window of k=24 indicating 24 hours or daily average as a replacement value) for each of the feature and filled the missing values.

- Dataset II: Multiple Imputation by Chained Equations 'MICE', is another R package [8] for filling missing values. MICE fills in data based on each feature and generates internally say a 'P' number of datasets where 'P' is equal to the number of columns (features(15)+targets(3)) in the dataset giving priority to each column. So we used 'with' function to explicitly get hold of the P number of datasets (in case of StationID 674 P=15+3=18) and 'pool' function to combine the datasets (in this case 18) to get a final MICE filled dataset.

- Dataset III: Using PCA on Dataset II we further reduced the features to just 8 from 15, since eight of the original features showed higher mutual information.

In each case, we have three subset datasets labeled A, B and C. All three subsets have only one target soil moisture content feature i.e. either of 2-inch or 8-inch or 20-inch.

- Subset A has all 15 features.

- Subset B has only those features which are relevant to that particular depth for example at depth 8-inch the features set does not include 2-inch or 20-inch soil temperature observed and salinity at 20-inch.

- Subset C is same as subset A but does not have any wind or air temperature related features.

In total 3x3 = 9 subsets are considered if all the features and targets are available (or less as in case of StationID 674 with 6 datasets).

Table 2: Types of Datasets generated for Site ID 674 with 8-inch as target soil moisture.

| Category | Dataset II - MICE filled | Soil Moisture | Datasets generated |
|---|---|---|---|
| 8-inch A | PREC.I-1,TOBS.I-1,STO.I-1:-2,STO.I-1:-8,STO.I-1:-20, SAL.I-1:-20,RDC.I-1:-8,RDC.I-1:-20,BATT.I-1,BATT.I-2, WSPDX.H-1,WSPDV.H-1,SRADV.H-1,RHUMV.H-1,RHUMX.H-1 | SMS.I-1:-8 | MICE-PCA (8-inch A) |
| 8-inch B | PREC.I-1,TOBS.I-1,STO.I-1:-8,RDC.I-1:-8, BATT.I-1,BATT.I-2,WSPDX.H-1,WSPDV.H-1, SRADV.H-1,RHUMV.H-1,RHUMX.H-1 | SMS.I-1:-8 | MICE-PCA (8-inch B) |
| 8-inch C | PREC.I-1,STO.I-1:-2,STO.I-1:-8,STO.I-1:-20, SAL.I-1:-20,RDC.I-1:-8,RDC.I-1:-20,BATT.I-1, BATT.I-2,SRADV.H-1,RHUMV.H-1,RHUMX.H-1 | SMS.I-1:-8 | MICE-PCA (8-inch C) |

Table 2 gives a glimpse of feature subsets, the target feature for each of A, B and C category belonging to dataset II generated by MICE filling of missing values. It also lists the datasets which can be generated from these subsets. Similarly keeping 2-inch, 20-inch data as the target value, 6 more subset datasets can be generated.

## 3 Method

Our method, illustrated in figure 5, includes data collection, data treatment, feature reduction or extraction, then the application of an LSTM network for soil moisture prediction. The setup and the layer wise workings of LSTM has been thoroughly discussed in [3]. We used the model prescribed in KERAS (a python library) for the core implementation of LSTM with epochs 50, batch size 72, 100 hidden units, linear activation as parameters with adam optimizer, root mean squared error as loss function using 48 months of data for training and 3 months of data in testing the model.
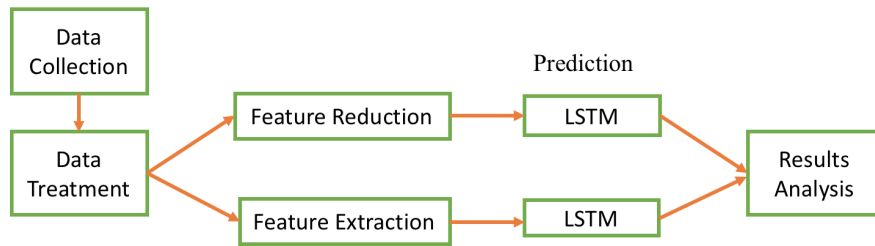


Figure 5: Overview of methodology

### 3.1 Feature Extraction using Principal Component Analysis (PCA)

Feature extraction is essential for effective summarization of all the features but on a reduced or minimum dimensions. We used a well-established statistical method PCA for features extraction which transforms data into a new space and still hold most of the original information. Useful features that are linear combinations of data attributes are derived by this approach. So our 15 features dataset was finally reduced to just 8 features reducing the dimension of dataset which showed more mutual information than compared to the original dataset.

Table 3: Feature extraction methods used and number of features after each method

| Feature extraction | Feature category | Original Features | Features based on percentage of missing values | Features based on visualization and mutual information | Features from PCA |
|---|---|---|---|---|---|
| Feature Names | Precipitation Accumulation | 1+1 | 1+0 | 1+0 | Eight significant principal components |
| | Air Temperature Observed | 1 | 1 | 1 | |
| | Soil Temperature Observed | 1+1+1 | 1+1+1 | 1+1+1 | |
| | Salinity | 1+1+1 | 1+1+1 | 0+0+1 | |
| | Real Dielectric Constant | 1+1+1 | 1+1+1 | 0+1+1 | |
| | Battery | 1+1+1 | 1+1+0 | 1+1+0 | |
| | Wind Direction Average | 1 | 1 | 0 | |
| | Wind Speed Maximum | 1 | 1 | 1 | |
| | Wind Speed Average | 1 | 1 | 1 | |
| | Solar Radiation Average | 1 | 1 | 1 | |
| | Relative Humidity Average | 1 | 1 | 1 | |
| | Relative Humidity Maximum | 1 | 1 | 1 | |
| total=N | | 21 | 19 | 15 | 8 |

## 4 Experimental Results

Applying LSTM to predict soil moisture content on 6 and 9 of the data subsets from Station ID 674 and Station ID 871 respectively we arrived at results presented in Table 4 and 5 with results of two datasets (II and III) including all three subsets (A,B and C). Due to space constraints we are only presenting the prediction and RMSE plots of station ID 674 in comparison with a well know machine learning method called Support Vector Machine (SVM) using radial basis function as kernel. All results are obtained with 50 epochs, using adam optimizer. For station 674 using dataset II (MICE filled) the RMSE was lowest among all combinations as seen in table 4.



Figure 6: Station ID 674 8-inch Dataset II Subset-Category A, B and C prediction results(a,b,c) with root mean square errors(d,e,f)
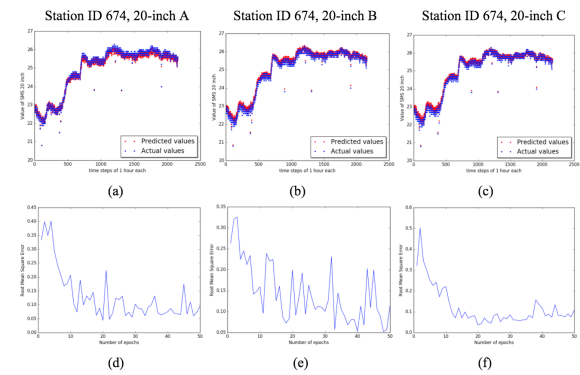
Figure 7: Station ID 674 20-inch Dataset II Subset-Category A, B and C prediction results(a,b,c) with root mean square errors(d,e,f)

Figures 7,8 and 10 show plots of predicted and actual values of soil moisture for station ID 674, plots of RMSE for each of the A, B and C subsets, for each of 8 and 20-inch, for each datasets II and III. Greater variations among values in 8-inch as shown in figure 7 when compared to values in Figure 8 can be attributed to the depth of the sensor device itself i.e. at 8-inch interference could be more compared to 20-inch depth. Figure 9 shows similar plot for station ID 871.
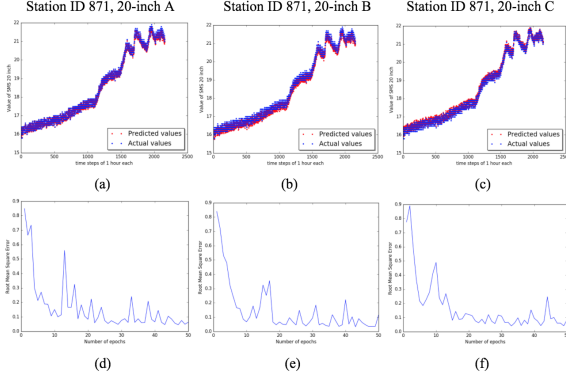
Figure 8: Station ID 871 20-inch Dataset II Subset-Category A, B and C prediction results(a,b,c) with root mean square errors(d,e,f)
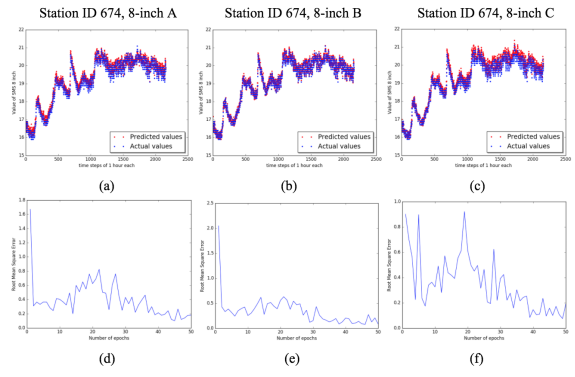
Figure 9: Station ID 674 8-inch Dataset III (MICE+PCA) Subset A, B and C prediction results (a,b,c) with root mean square errors(d,e,f)

Table 4: Station ID 674 RMSE values using LSTM and SVM: minimum, maximum and average values(LSTM) using Dataset II(missing values filled by MICE) and Dataset III(PCA on Dataset II)

| 674 | LSTM | | | | | | SVM | |
| | MICE | | | MICE+PCA | | | MICE | MICE+PCA |
| RMSE | min | max | avg | min | max | avg | | |
| 8-inch A | 0.103 | 1.209 | 0.373 | 0.101 | 1.671 | 0.398 | 6.217 | 8.201 |
| 8-inch B | 0.052 | **0.326** | 0.143 | 0.079 | 2.048 | **0.340** | 6.325 | 8.205 |
| 8-inch C | 0.088 | 1.340 | 0.345 | **0.075** | 0.922 | 0.351 | 6.222 | 7.499 |
| 20-inch A | 0.044 | 0.401 | 0.131 | 0.390 | 1.450 | 0.756 | 11.113 | 13.895 |
| 20-inch B | 0.049 | 0.702 | 0.301 | 0.451 | **0.837** | 0.546 | 12.058 | 17.131 |
| 20-inch C | **0.037** | 0.502 | **0.121** | 0.275 | 0.878 | 0.393 | 11.335 | 15.559 |

A slight skewness can be seen in figure 7 (c) and more skewness in figure 10 (c) where majority of the predicted values are slightly larger than the actual values which is why the red dots appear above the blue dots. This skewness can be attributed non inclusion of wind related features in the C subset of respective datasets, especially at 8-inch those features must be significant which is why we don't see any skewness for 20-inch plots on C subset datasets (figures 8(c), 9(c)). The greater proximity in actual and predicted values in the plots of figure 8 and 9 can be attributed to lesser disturbances at 20-inch than at 8-inch. Much worse SVM results of Station 871 are excluded from table 5. Also from tables 4 and 5 it is clear that RMSE was higher whenever PCA was used for feature selection.

Table 5: Station ID 871 RMSE values: minimum, maximum and average values using Dataset II(missing values filled by MICE) and Dataset III(PCA feature extraction on Dataset II)

| 871 | MICE | | | MICE+PCA | | |
| RMSE | min | max | avg | min | max | avg |
| 2-inch A | 0.288 | 1.441 | 0.625 | 3.905 | 8.36 | 6.83 |
| 2-inch B | 0.294 | 2.181 | 0.836 | 0.585 | 2.201 | 0.818 |
| 2-inch C | 0.312 | 1.443 | 0.627 | 0.525 | 2.209 | 0.850 |
| 8-inch A | 0.195 | 1.505 | 0.614 | 4.004 | 6.920 | 5.65 |
| 8-inch B | 0.385 | 1.352 | 0.714 | 0.159 | 1.542 | 0.667 |
| 8-inch C | 0.183 | 2.080 | 0.600 | 0.101 | 1.573 | 0.592 |
| 20-inch A | 0.046 | 0.850 | 0.165 | 1.843 | 3.506 | 2.715 |
| 20-inch B | **0.033** | **0.840** | **0.149** | 0.086 | **0.979** | **0.235** |
| 20-inch C | 0.041 | 0.891 | 0.170 | **0.052** | 1.020 | 0.287 |

# 5 Conclusions

This paper has explored the application of an LSTM based neural network learning model for soil moisture prediction, applied to real (SCAN and SNOTEL) data. In our work the data treatment and feature extraction steps played significant role in obtaining better results. Application of *MICE* for missing values imputation gave best results. We created a benchmark dataset of 6 and 9 subsets for sites 674 and 871 respectively based on sensor depth. PCA based feature extraction with 8 significant features performed well but with some degradation compared to the initial 15 features. Achieved RMSE minimization during both training and testing with an overall minimum RMSEs in range of 0.03 to 0.4 at different depths for majority of the datasets created for both the sites.

In future work, we plan to improve results by using RNNs for filling the missing data as proposed in [9], to make our novel approach an end to end neural network based soil moisture content prediction and forecasting system. Also fine tune LSTM parameters, to achieve negligible RMSE. As Neural networks are sensitive to outliers in the data, we plan to adopt methods in [10] for outlier detection and replacement, and investigate if statistical features can further improve the results.

## References

[1] Tarnpradab, Sansiri, et al. "Neural networks for prediction of streamflow based on snow accumulation." Computational Intelligence for Engineering Solutions (CIES), 2014 IEEE Symposium on. IEEE, 2014.

[2] Bai, Yun, et al. "Daily reservoir inflow forecasting using multiscale deep feature learning with hybrid models." Journal of Hydrology532 (2016): 193-206.

[3] Fang, Kuai, et al. "Prolongation of SMAP to Spatio-temporally Seamless Coverage of Continental US Using a Deep Learning Neural Network." arXiv preprint arXiv:1707.06611 (2017).

[4] Zhang, Dongying, et al. "Upscaling of Surface Soil Moisture Using a Deep Learning Model with VIIRS RDR." ISPRS International Journal of Geo-Information 6.5 (2017): 130.

[5] Schaefer, Garry L., and Ron F. Paetzold. "SNOTEL (SNOwpack TELemetry) and SCAN (soil climate analysis network)." Proc. Intl. Workshop on Automated Wea. Stations for Appl. in Agr. and Water Resour. Mgmt. 2001.

[6] NRCS website for SCAN and SNOTEL datasets. https://www.wcc.nrcs.usda.gov/scan/

[7] Moritz, Steffen, and Thomas Bartz-Beielstein. "imputeTS: Time Series Missing Value Imputation." R package version 0.4 (2015).

[8] Buuren, Stef, and Karin Groothuis-Oudshoorn. "mice: Multivariate imputation by chained equations in R." Journal of statistical software 45, no. 3 (2011).

[9] Che, Zhengping, et al. "Recurrent neural networks for multivariate time series with missing values." arXiv preprint arXiv:1606.01865 (2016).

[10] Zhao, Zhiruo, Chilukuri K. Mohan, and Kishan G. Mehrotra. "Adaptive Sampling and Learning for Unsupervised Outlier Detection." In FLAIRS Conference, pp. 460-466. 2016.

# Modified Neural Network Based Approach for Dynamic Collision-free Trajectory Generation

Hang Yin[1], Chilukuri K. Mohan[2] and Utpal Roy[3]

[1,3]Mechanical & Aerospace Engineering
College of Engineering and Computer Science
Syracuse University, NY
hayin@syr.edu, uroy@syr.edu

[2]Electrical Engineering and Computer Science
College of Engineering and Computer Science
Syracuse University, NY
ckmohan@syr.edu

## Abstract

In this paper, the author proposed an effective method to achieve both obstacle-avoidance and target-tracking for an autonomous agent in a complex environment. The proposed method is based on a topologically organized neural network, where the dynamics of each neuron is characterized by a set of governing equations. The activities of neurons create an artificial potential field that produces a feasible continuous path to guide the autonomous agent towards to the target while avoiding static and/or dynamic obstacles.

## 1 Introduction

An autonomous agent/vehicle should safely carry out the mission in hazardous and populated environments. To perform its task, an autonomous agent should reach a given target without collision with obstacles such as walls, trees, humans, and so on.

In the past, numerous researchers have worked on the path planning problem. Most of the works so far deal with static environments and used global methods, which could be viewed as a search process for a path in a graph. However, global methods limit the real-time capabilities of autonomous agents in a cluttered environment, especially when new information, about changes in the environment, is becoming available continuously, because of the time needed to perform the planning task. Later, the idea of using an artificial potential field around each obstacle was proposed combined with an attractive potential around the target. Unfortunately, the artificial potential field method suffers from a minimal local problem which may result in trajectory generation failure.

Several neural network models were proposed to generate real-time trajectories through learning. Muniz [1] proposed a neural network model for the navigation of an autonomous agent, which can generate dynamic trajectories with obstacle avoidance through unsupervised learning. However, this model is computationally complicated since it incorporates the vector associative map model and the direction-to-rotation effector control transform model [2] [3]. Hong et al. [4]proposed a modified pulse-coupled neural network (MPCNN) model for real-time collision-free path planning of agents in nonstationary environments. This model has to mesh the space into the grid which means some space which is relatively empty may have the same number of neurons with another part of the space where exists obstacles. Furthermore, the generated path is the sequence of neuron's parents from the agent to the target. To have a smooth trajectory, space has to mesh into a fine grid which required more computation power when the neural network should be updated.

Inspired by Fritzke's work [5] where the structure of the neural network is changing all the time based on the observed error. New neurons will be inserted into the network if the observed error is accumulated locally. If

we can allocate more neurons at the location where more obstacles are observed and allocate fewer neurons in the space, we may reduce the computation time to construct the artificial potential field.

## 2 Review of Neural Network Approach Proposed by Muniz

The neural network architecture model proposed by Muniz [1] is a discrete topologically organized map that has been used in many neural network models. This model is expressed in a finite-dimensional state space. Essentially, the objective is to represent a static or dynamic environment with an artificial potential field which results in the generation of a feasible path through the given environment. The dynamics of the $i$th neuron in the neural network is characterized by a governing equation show as follows:

$$\frac{dx_i}{dt} = -Ax_i + (B - x_i)\left([I_i]^- + \sum_{j=1}^{k} \omega_{ij}[x_{ij}]^-\right) - (D + x_i)[I_i]^+$$

where $k$ is the number of neighboring neurons of the $i$th neuron; $[I_i]^- + \sum_{j=1}^{k} \omega_{ij}[x_{ij}]^-$ is the excitatory inputs; $[I_i]^+$ is the inhibitory inputs; $A$ is a nonnegative constants representing the passive decay rate; $B$ is the upper bounds of the neural activity; $D$ is the lower bounds of the neural activity; and $I_i$ is the external input to $i$th neuron. The value of $I_i$ is defined as follows:

$$I_i = \begin{cases} -E, & if\ there\ is\ a\ target \\ E, & if\ there\ is\ an\ obstacle \\ 0, & otherwise \end{cases}$$

where $E$ is a large positive constant which defined as $E \gg B$. Function $[I_i]^+$ is a linear-above-threshold function defined as $[I_i]^+ = \max(I_i, 0)$ and similarly the function $[I_i]^-$ is defined as $[I_i]^- = \max(-I_i, 0)$. The connection weight $\omega_{ij}$ from the $j$th neuron to the $i$th neuron is a function of distance defined as:

$$\omega_{ij} = \frac{\mu}{d_{ij}}$$

where $d_{ij} = |q_i - q_j|$ is the Euclidean distance between $i$th neuron's position with $j$th neuron's position.

This neural network guarantees that the positive neural activity can propagate to the whole state space through lateral neural connections, while the negative activity stays local since there are no inhibitory connections among neurons. Therefore, the target globally influences the whole space to attract the agent, while the obstacles have the only local effect to avoid collisions. Also, the activity propagation from the target is blocked when it hits any obstacles.

For example, in order to guide the agent (marked with red triangular in Fig. 1) out of the U shape obstacle (marked with black), and to its target position (marked with a red cross), the original NN approach has to partition the space into grid shown in Fig. 2.
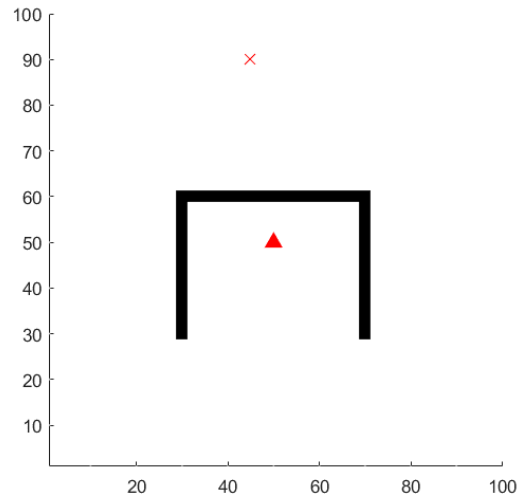


Fig. 1. Concave U shape obstacle with the agent

Each neuron, it is affected by its neighbors only. The neighbor of the *i*th neuron is defined as: *i*th neuron's neighbors are neurons which are connected with *i*th neuron directly. The structure of the neurons showed as Fig. 3. Then the artificial potential field is generated based on the activity of each neurons (Fig. 4).
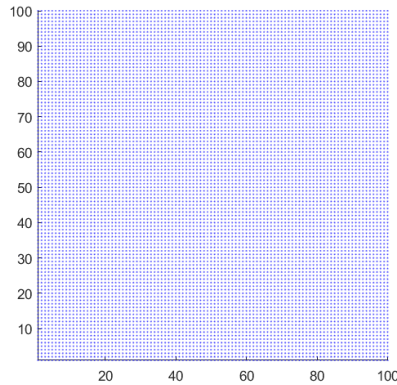


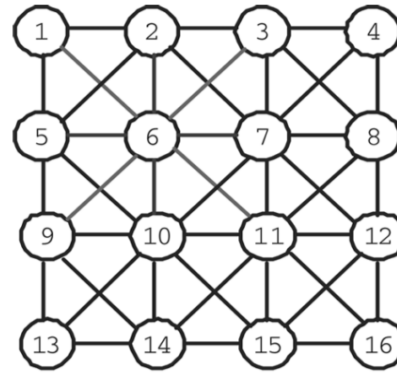Fig. 2. NN approach neuron allocation
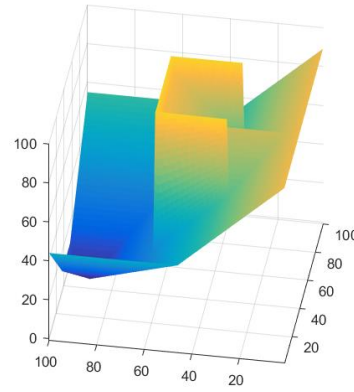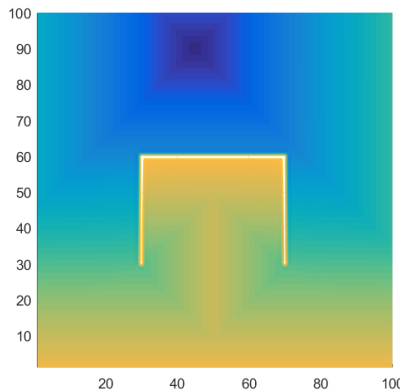


Fig. 3. NN approach neuron allocation



Fig. 4. Potential field generated by using original neural network approach

The original neural network approach evenly distributes neurons in the whole space, which means no matter how complex or simple the environment is, the density of neurons in the whole space will not change accordingly. In another word, if some space is relatively empty, this algorithm will allocate the same number of neurons than another more complex environment. This results in carrying out unnecessary computation in the neurons which don't affect the result much.

Therefore, if we can allocate neurons smartly by allocating more neurons in the space where the environment is complex and fewer neurons in free space and reducing the total number of neurons and increasing the efficiency of the algorithm.

## 3 Modified Neural Network Approach

Inspired by Fritzke's work, Growing Cell Structures (GCS) network [5], which the structure of the neural network is changing all the time based on the observed error. New neurons will be inserted into the network if the observed error is accumulated locally, and neurons with least signal counter value are chosen for deletion. The deletion mechanism allows GCS to approximate each region accurately either in static environment or dynamic environment. Moreover, the Voronoi diagram has been employed to partition the space.

Voronoi diagram partitioning a space into regions based on distance to points in a specific subset of the space. That set of points is specified beforehand, and for each point, there is a corresponding region consisting of all

points closer to that seed than to any other. These regions are called Voronoi cells. The Voronoi diagram of a set of points is dual to its Delaunay triangulation [6]. One important property of Voronoi diagram is its geometric stability of Voronoi diagrams, e.g., a change caused by some translation or distortion, yields a small change in the shape of the Voronoi cells.

Therefore, if we can employ Voronoi diagram guiding us allocate neurons in the space smartly, we may reduce the total number of neurons and increase the efficiency of the algorithm. Fig. 5 shows partitioning results of the given environment. Notice, the red triangle is the agent; the black blocks are obstacles, and the red cross is the target. This Voronoi diagram is created by given the locations of the agent, the target, and the obstacles. Based on the partitioning result, new neurons are inserted into the location of Voronoi cell vertex.

As it is evident from the Fig. 5, the neurons have been allocated according to the given environment. More neurons are located in the space where lots of obstacles have been observed, while fewer neurons have been allocated in the space where is relatively sparse. Overall, around 400 neurons have been used. In comparison, 10000 neurons will be used if we employ original neural network approach.



Fig. 5. Space partitioning by using Voronoi diagram

Unlike the original neural network where it is pretty easy to define one neuron's neighbor. For freely allocated neurons, there exists no specific geometry to describe one neuron's cell or its neighbors. Hence, how to define and how to find one neuron's neighbors becomes a critical issue.

The neighbors of $i$th neuron have been defined in the following way: $i$th neuron's neighbors are the neurons whose Voronoi cells connect with $i$th neuron's Voronoi cell directly. For instance (as shown Fig. 6), the neuron "A" locate in the center has four neighbors.

Then the neural network is updated by using the governing equations until the whole network becomes stable. To generate a smooth artificial potential field, interpolation techniques have been employed. Specifically, the scattered data interpolation [7] method has been used to estimate the potential value of a point sitting among neurons. Furthermore, the original governing equation has been modified as follows.



Fig. 6. Neuron structure and neighbors

$$\frac{dx_i^{att}}{dt} = -Ax_i + (B - x_i)\left([I_i]^- + \sum_{j=1}^{k} \omega_{ij}[x_{ij}]^-\right)$$

$$\frac{dx_i^{rep}}{dt} = -Ax_i - (D + x_i)[I_i]^+$$

where $x_i^{att}$ is the attractivity behavior and $x_i^{rep}$ is the repulsive behavior. Therefore, the global attractive potential and local repulsive potential field is defined as:

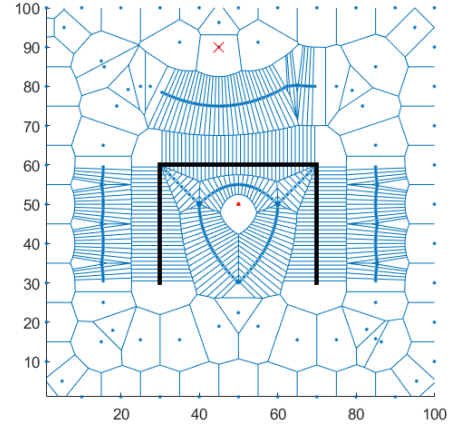$$U_{ga} = -\frac{1}{2}k^{ga}x_i^{att}, \qquad U_{lr} = \frac{k^{lr}}{2x_i^{rep}}$$

where $k^{ga}$ is the magnitude of the global attractive potential and $k^{lr}$ is the magnitude of the local repulsive potential. The generated potential field is shown in Fig. 7.

Since we use much fewer neurons to generate the artificial potential field, the shape of the repulsive potential field is not as good as the potential field generated by using the original neural network. To increase the computational speed, it is necessary to neglect some non-essential environment.
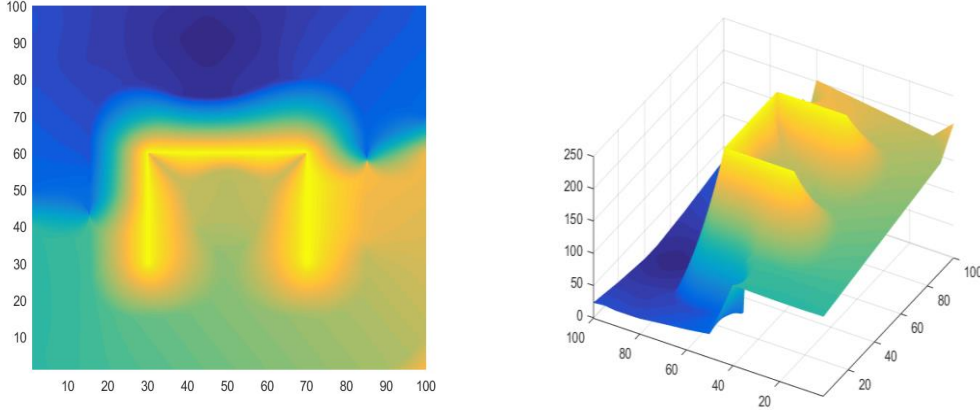


Fig. 7. APF generated by using Modified NN

After the artificial potential field has been generated, the obstacle-free trajectory can be established. The autonomous agent in this field will be attracted or repulsed away according to the gradients of the potential fields. The induced force $F$, which indicates the force applied on the agent to guide the agent, is:

$$F = F_{ga} + F_{lr} = -\nabla U_{ga} - \nabla U_{lr}$$

## 4 Simulation Results

We use a random generated map to test the real-time trajectory generation performance of our algorithm. 400 obstacles are randomly generated in a 100 by 100 unit length 2D environment. An autonomous agent located at $(1,1)$ and a moving target located at $(40,90)$. The velocity of the moving target is $\boldsymbol{v} = (0.3,0)$. The autonomous agent can detect obstacles within a limited range only. At the beginning of the simulation, the agent detects 11 obstacles' existence. As time goes by, more and more obstacles are detected, and new neurons been inserted to the neural network. At end of the simulation, the agent successfully arrived the target location without collide with obstacles.
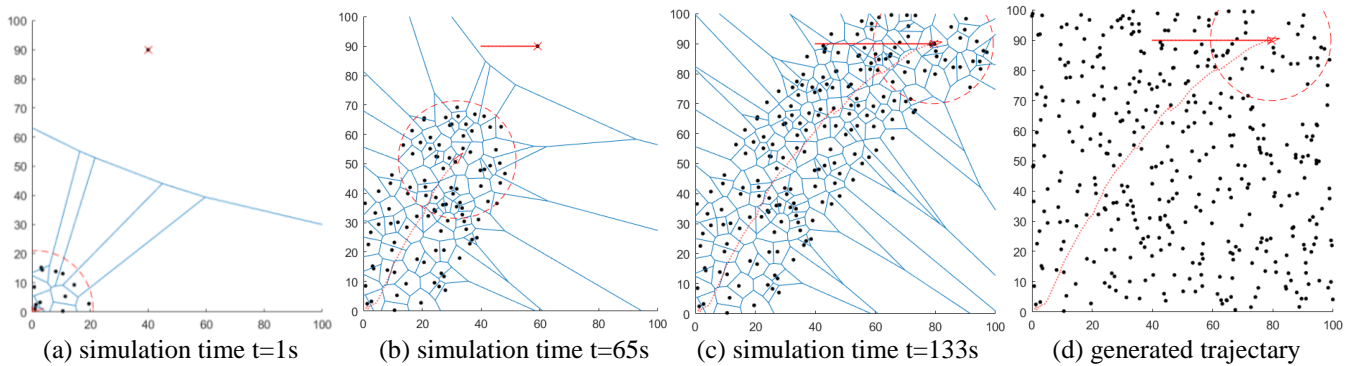


(a) simulation time t=1s     (b) simulation time t=65s     (c) simulation time t=133s     (d) generated trajectary

Fig. 8. Target tracking simulation

218

This paper introduces Voronoi diagram into the neural network approach which may require more computational time to finish. Furthermore, in order to generate the artificial potential field, each neuron's activities have been separated into two parts. So, it is difficult to compare the modified algorithm with the original one on the computational time basis.



Fig. 8. APF generated by using modified NN

To compare NN, MNN and A* method (A* is a widely used pathfinding and graph traversal algorithm [8]). Monte Carlo simulation has been employed. Overall, a hundred different scenarios have been tested. For each scenario, 40 obstacles have been randomly generated in a 100 by 100 unit length environment. For the experiment, an autonomous agent and its target are located at two separate location: (1,1) and (90,90) respectively. After the environment is established, the associated potential field is created, which helps in generating agent's trajectory based on the gradient descent of the developed potential field. Time consumption for generating a feasible trajectory has been recorded and the performance of the generated potential field has been evaluated by comparing trajectory for a different scenario. All statistics of 100 tests are collected and they have been illustrated in the following Fig. 10 and Fig. 11.

As we can see, A* algorithm (average runtime: 1.86s) spend 2 times more time than neural network method (average runtime: 0.81s), and 4 times more than modified neural network (average runtime: 0.42s). By applying Voronoi diagram, we can reduce the computational time dramatically, even though the Voronoi diagram related calculation is relatively complicated. Also, the generated trajectory has similar performance, though the modified neural network has a larger
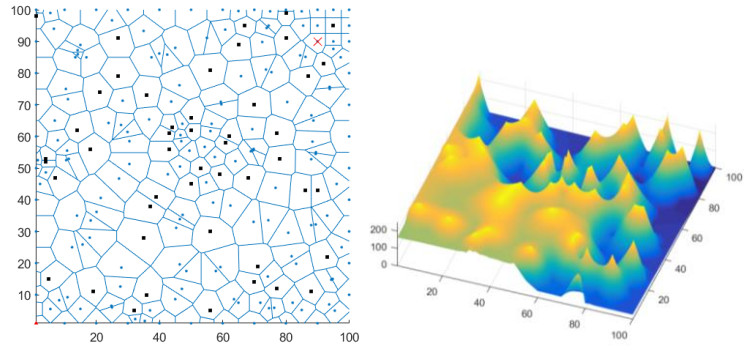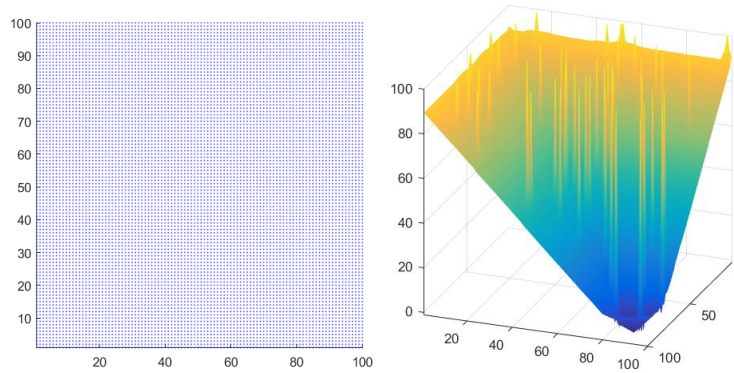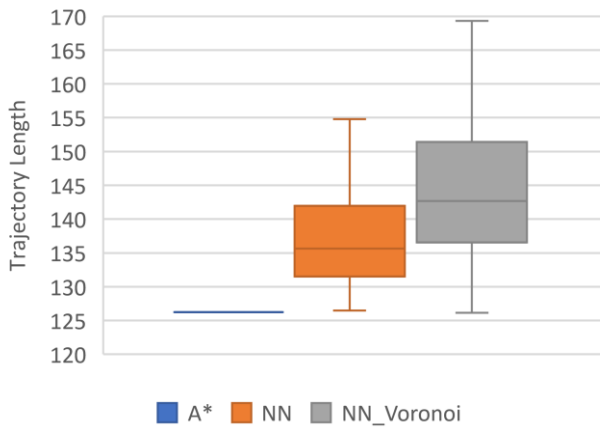


Fig. 9. APF generated by using original NN
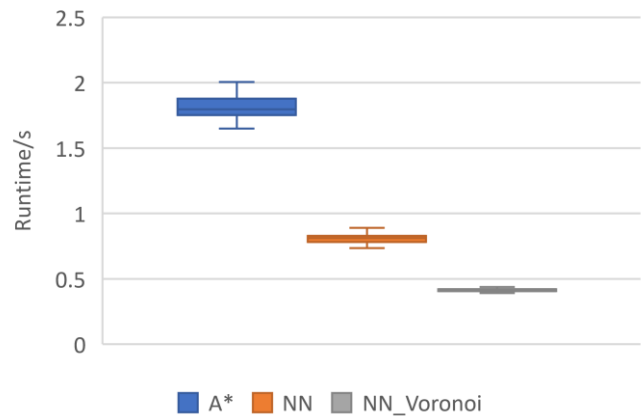


Fig. 10. Trajectory length



Fig. 11. Time consumption

variance. Therefore, we can conclude that, by employing Voronoi diagram, we can reduce the computational time while generating a feasible trajectory.

In addition, to check the robustness of the modified algorithm, a maze-like environment was built. The generated artificial potential field and generated trajectory are showed in Fig. 12 and Fig. 13.
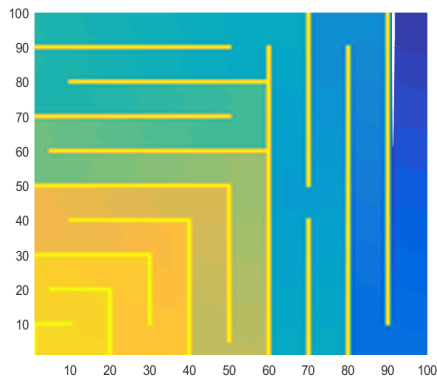


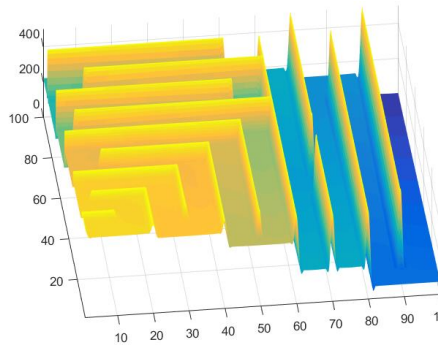Fig. 12. Generated artificial potetnial field by using MNN

Fig. 13. MNN trajectory

## 5 Conclusion

In this paper, a modified neural network approach is proposed for the dynamic collision-free trajectory generation in an arbitrarily dynamic environment. By employing the Voronoi diagram, neurons can be allocated based on the environmental information rather than evenly distributed in the space blindly. The potential field is then generated through the activities of the neurons that represent the varying environment. The real-time trajectory is then generated by applying gradient decent rules to the artificial potential field.

## Reference

[1] F. e. a. Muniz, "Neural controller for a mobile robot in a nonstationary environment," in *IFAC Conference Intelligent Autonomous Vehicle*, Helsinki, 1995.

[2] P. Gaudiano, E. Zalama and J. L. Coronado, "An unsupervised neural network for low-level control of a mobile robot: Noise resistance, stability, and hardware implementation," in *IEEE Transictions Systems, Man, Cybernetic*, 1996.

[3] E. Zalama, "A real-time, unsupervised neural network for the low-level control of a mobile robot in a nonstationary environment," *Neural Networks,* vol. 8, pp. 103-123, 1995.

[4] H. Qu, S. X. Yang, A. R. Willms and Z. Yi, "Real-time robot path planning based on a modified pulse-coupled neural network model," *IEEE Transactions on Neural Networks,* vol. 20, no. 11, pp. 1724-1739, 2009.

[5] B. Fritzke, "Growing cell structures—A self-organizing network for unsupervised and supervised learning," *Neural Network,* vol. 7, no. 9, pp. 1441-1460, 1994.

[6] "Voronoi diagram," Wikipedia, [Online]. Available: https://en.wikipedia.org/wiki/Voronoi_diagram.

[7] R. Franke, "Scattered Data Interpolation: Tests of Some Methods," *Mathematics of Computation,* vol. 38, 1982.

[8] P. E. Hart, N. J. Nilsson and B. Raphael, "A Formal Basis for the Heuristic Determination of Minimum Cost Paths," *IEEE Transactions on System Science and Cybernetics,* vol. 4, no. 2, pp. 100-107, 1968.

[9] L. Wyard-Scott and Q.-H. M. Meng, "A potential maze solving algorithm for a micromouse robot," in *Communications, Computers, and Signal Processing*, Victoria, 1995.

[10] A. Zelinsky, "Using path transforms to guide the search for find path in 2D," *International Journal Robotic Research,* vol. 13, no. 4, pp. 315-325, 1994.

[11] J. Ilari and C. Torras, "2D path planning: A configuration space heuristic approach," *International Journal Robotic Research,* vol. 9, no. 1, pp. 75-91, 1990.

[12] R. Glasius, "Population coding in a neural net for trajectory formation," *Network: Computation Neural System,* vol. 5, pp. 549-563, 1994.

[13] O. Khatib, "Real-time obstacle avoidance for manipulators and mobile robots," *International Journal Robotic Research,* vol. 5, pp. 90-98, 1986.

[14] C. Seshadri and A. Ghosh, "Optimum path planning for robot manipulators amid static and dynamic obstacles," *IEEE Transaction Systems, Man, and Cybernetics Society,* vol. 23, no. 2, pp. 576-584, 1993.

[15] Z. X. Li and T. D. Bui, "Robot path planning using fluid model," *Journal of Intellegence Robotic System,* vol. 21, pp. 29-50, 1998.

[16] G. Oriolo, G. Ulivi and M. Vendittelli, "Real-time map building and navigation for autonomous robots in unknown environments," *IEEE Transactions on Systems, Man, and Cybernetics,* vol. 23, no. 2, pp. 576-584, 1993.

[17] S. X. Yang and M. Meng, "Neural Network Approaches to Dynamica Collision-Free Trajectory Generation," *Systems, Man, and Cybernetics,* vol. 31, no. 3, 2001.

[18] Chen, Yi-Wen; Chiu, Wei-Yu, "Optimal robot path planning system by using a neural network-based approach," in *Automatic Control Conference (CACS), 2015 International*, 2015.

[19] Momi, Elena D.; Kranendonk, Laurens; Valenti, Marta; Enayati, Nima; Ferrigno, Giancarlo, "A Neural Network-based Apprach for Trajectory Planning in Robot-Human Handover Tasks," *Frontiers in Robotics and AI,* vol. 3, 2016.

[20] D. Simon, "The application of neural networks to optimal robot trajectory planning," *Robotics and Autonomous Systems,* vol. 11, no. 1, pp. 23-34, 1993.

[21] Y. Bassil, "Neural Network Model for Path-Planning of Robotic Rover Systems," *International Journal of Science and Technology,* vol. 2, no. 2, 2012.

[22] D. Tsankova, "Neural Networks Based Path Planning," *International Journal of Engineering,* vol. 3, 2011.

[23] H. Ouarda, "Neural Path Planning for Mobile Robots," *International Journal of Systems Appllications, Engineering & Development,* vol. 5, no. 3, 2011.

[24] V. Burzevski and C. K. Mohan, "Hierarchical Growing Cell Structures," 1996.

# Locating People of Interest in Social Networks

Pivithuru Wijegunawardana     Sucheta Soundarajan

{ppwijegu,susounda}@syr.edu, Syracuse University, Syracuse, NY

**Abstract**

The focus of the current research is to identify people of interest in dark networks, which represent illegal or covert activity. In such networks, people are unlikely to disclose accurate information when queried. We present `RedLearn`, an algorithm for sampling dark networks with the goal of identifying as many nodes of interest as possible . Results on real-world networks, show that `RedLearn` achieves up to a 340% improvement over the next best sampling algorithm.

## 1   Introduction

In this paper, we consider the problem of sampling a 'dark' network (i.e., a network representing illegal or covert activity) with the intention of observing as many "persons of interest" (POI) as possible given a limited query budget. Here, a POI is a node possessing a certain attribute (e.g., individuals involved in a criminal action). We assume that we begin with knowledge of one POI in the network, with the rest of the network unobserved (both in terms of topology as well as node attributes) and, we continue to locate POIs and uncover the network by querying nodes. Once we query a node ($v$), we assume that we get to know a) Whether $v$ is a POI or not, b) Who $v$'s neighbors are and, c) What color estimates $v$ would give to its neighbors.

A major complicating factor in this problem is that due to the covert nature of the networks being studied, one cannot expect the observed information to be reliable. Therefore, we consider that a) POTs can hide edges between POI nodes b) nodes can misreport whether their neighbors are POIs.

## 2   Methodology

We label a POI node as a red node and a non POI node as a blue node in the network. We present `RedLearn` [1], a learning-based algorithm for sampling networks where the learning algorithm (logistic regression classifier) predicts the probability of an observed node being red. `RedLearn` uses network structure-based features such as number of red/blue neighors, number of red traingles a node is a part of to learn the patterns of connections between red nodes (e.g., homophily vs. anti-homophily). Neighbor color estimates-based features such as how many red/blue neighbors have estimated the node color as red/blue, what is the probability a node color is red based on these estimations, to learn the relationship between what a node estimates about its neighbors' colors and the true colors of those neighbors. We formulate the probability of a node misreporting its neighbor color propotional to the hierarchy of both parties in the criminal organization and honesty of the person reporting.
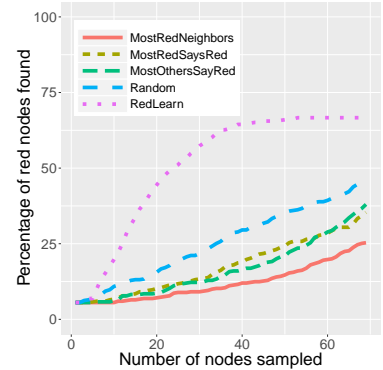


Figure 1: Comparison of POI sampling algorithms in NoordinTop terrorist network when only red nodes are aware of other red nodes and no homophily is present between red nodes.

We show that in cases where the POIs are likely to be connected to other POIs, a simple algorithm of choosing the node with the most red neighbors works well. Figure 1, shows a comparison of POI sampling algorithms, in the more realistic scenario where POIs hide their connections with other POIs in NoordinTop terrorsit network. In this network we have assigned node colors based on whther they use some communication medium or not. In this scenario, `RedLearn` shows outstanding performance, beating the next best strategy by up to 340%.

## References

[1] Wijegunawardana P, Ojha V, Gera R, Soundarajan S. Seeing Red: Locating People of Interest in Networks. In Workshop on Complex Networks CompleNet 2017 Mar 21 (pp. 141-150).