

Analysis and Exploitation of Musician Social Networks for Recommendation and Discovery

Ben Fields, *Student Member, IEEE*, Kurt Jacobson, *Student member, IEEE*,
Christophe Rhodes, Mark Sandler, *Senior Member, IEEE* and Michael Casey, *Senior Member, IEEE*

Abstract—This paper presents an extensive analysis of a sample of a social network of musicians. The network sample is first analyzed using standard complex network techniques to verify that it has similar properties to other web-derived complex networks. Content-based pairwise dissimilarity values between the musical data associated with the network sample are computed, and the relationship between those content-based distances and distances from network theory explored. Following this exploration, hybrid graphs and distance measures are constructed, and used to examine the community structure of the artist network. Finally, results of these investigations are presented and considered in the light of recommendation and discovery applications with these hybrid measures as their basis.

I. INTRODUCTION

As more freely-available audio content continues to become accessible, listeners require more sophisticated tools to aid them in the discovery and organization of new music that they will find enjoyable. This need, along with the advent of Web-based social networks and the steady progress of signal-based music information retrieval have created an opportunity to exploit both social relationships and acoustic similarity in recommender systems. Combining these relations can provide a means to improving the understanding of the complex relationship between songs, which itself is required to improve song recommendation. One way to do that is to base recommendations on more information than is provided by a single distance measure between

songs, allowing the production of systems capable of mediating content-based recommendations with given social connections, and hence the construction of socially structured playlists.

Social networks present a way for nearly anyone to distribute their own media and, as a result, there is an ever-larger amount of available music from an ever-increasing array of artists. Given this environment of content, how can we best use all of the available information to discover new music? Can both social metadata and content based comparisons be exploited to improve discovery of new material? Can this crowd-sourced tangle of social networking ties provide insights into the dynamics of popular music? Does the structure of a network of artists have any relevance to music-related studies such as music recommendation or musicology?

Motivated by this, we examine the Myspace artist network. Though there are a number of music oriented social networking websites (*e.g.* Soundcloud¹, Jamendo², etc.), Myspace³ has become the *de facto* standard for web-based music artist promotion. For the purpose of this paper, *artist* and *artist page* are used interchangeably to refer to the collection of media and social relationships found at a specific Myspace page residing in Myspace's artist subnetwork. Although exact figures are not made public, recent estimates suggest there are well over 8 million artist pages⁴ on Myspace.

The Myspace social network, like most social networks, is based upon relational links between *friends* designating some kind of association. Within each Myspace user's friends there is a subset of between 8 and 40 *top friends*. While generic friends are mutually confirmed, individual users unilaterally elevate top friends from the generic friends set. Additionally, pages by *artists* contain streaming and downloadable media of some kind either audio, video or both.

B. Fields*, C. Rhodes and M. Casey are with the Intelligent Sound and Music Systems group, Department of Computing, Goldsmiths, University of London, New Cross, London SE14 6NW, United Kingdom. *e-mail: b.fields@gold.ac.uk

K. Jacobson and M. Sandler are with the Centre for Digital Music, School of Electronic Engineering and Computer Science, Queen Mary University of London, Mile End Road, London E1 4NS, United Kingdom.

M. Casey is also with the Bregman Music Audio Research Studio, Music Department, Dartmouth College, Hanover, New Hampshire, 03755, United States

Manuscript received November 15, 2009; revised February 16, 2010.

¹<http://www.soundcloud.com/>

²<http://www.jamendo.com/>

³<http://www.myspace.com/>

⁴<http://techradar1.wordpress.com/2008/01/11/facebookmyspace-statistics/>

To work towards answering the questions posed above, we explore the relationship between the connectivity of pairs of artists on the Myspace top friends artist network and measures of acoustic dissimilarity between these artists. Furthermore, we identify communities of artists based on the Myspace network topology and attempt to relate these community structures to musical genre. Finally, we present a prototype system of music playlist generation, paying attention to means for its evaluation.

Immediately following this section is a review of relevant literature from complex network theory and signal-based music analysis. Section III then provides a detailed discussion of the sample network's properties and our methods. The initial experiments into the relationship between the social connectivity and the acoustic feature space and their results are presented and discussed in Section IV. The implications and direction for future work are discussed in Section V.

II. BACKGROUND

We begin the background with a discussion of existing tools for the analysis and manipulation of networks in Section II-A. This subsection covers complex network analysis, network flow analysis, particular issues pertaining to networks of musicians and community structure. In Section II-B we examine highlights of past work in audio content-based music similarity. Lastly, Section II-C provides background in the independence measures we use in later sections.

A. Existing Tools for Networks

1) *Complex Networks*: Complex network theory deals with the structure of relationships in complex systems. Using the tools of graph theory and statistical mechanics, physicists have developed models and metrics for describing a diverse set of real-world networks – including social networks, academic citation networks, biological protein networks, and the World-Wide Web. It has been shown that these diverse networks often exhibit several unifying characteristics such as small worldness, scale-free degree distributions, and community structure [1]. We will discuss some of the characteristics of the Myspace artist network in III-B. For a more in depth discussion of complex network analysis techniques the reader is referred to [1], [2]

2) *Network Flow Analysis*: The basic premise in network flow analysis is to examine a network's nodes as sources and sinks of some kind of *traffic* [3]. Typically, though not exclusively, flow networks are directed, weighted graphs. Many useful measures for determining the density of edge connectivity between sources and

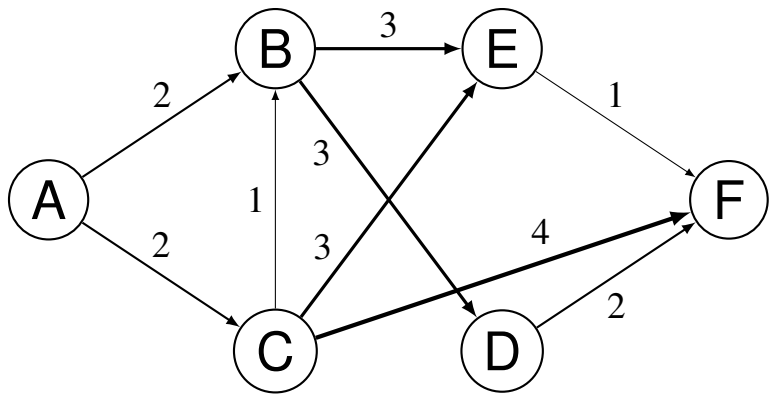


Fig. 1: A simple flow network with directed weighted edges. Edge width is representative of node capacity, which is also labeled on each edge. Treating node A as the source and node F as the sink, the maximum flow is 4.

sinks can be found in this space [4]. One of the most common among them is the Maximum Flow, which is a means of measuring the maximum capacity for fluid to flow between a source node to a sink node or, equivalently, the smallest sum of edge weights that must be *cut* from the network to create exactly two subgraphs, one containing the source node and one containing the sink node. This equivalency is the Maximum Flow/Minimum Cut Theorem [5]. If the edges in a graph are unweighted, this value is also equivalent to the number of paths from the source to the sink which share no common edges. Mature algorithms, incorporating a number of optimization strategies, are available for computing the maximum flow between nodes [3], [6].

An example of Maximum Flow can be seen on the network in figure 1. The narrowest flow capacity from node A to node F are the edges $E(a, b)$ and $E(a, c)$, where $E(a, b) + E(a, c) = 4$. The maximum flow can simply be found by taking the sum of the magnitude of each edge in the minimum cut set.

The few examples of network flow analysis being applied in music informatics deal primarily with constructing playlists using segments of a complete solution to the Traveling Salesman Problem [7] or use exhaustive and explicit textual metadata [8] without comparisons to content-based metrics.

3) *Musician Networks*: Quite naturally, networks of musicians have been studied in the context of complex network theory – typically viewing the artists as nodes in the network and using either collaboration, influence, or similarity to define network edges. These networks of musicians exhibit many of the properties expected in social networks [9]–[11]. However, these studies all

examine networks created by experts (*e.g.* All Music Guide⁵) or via algorithmic means (*e.g.* Last.fm⁶) as opposed to the artists themselves, as is seen in Myspace and other similar networks. Networks of music listeners and bipartite networks of listeners and artists have also been studied [12], [13].

4) *Community Structure*: Recently, as more data heavy complex networks have been created across many domains, there has been a significant amount of interest in algorithms for detecting community structures in these networks. These algorithms are meant to find dense subgraphs (communities) in a larger sparse graph. More formally, the goal is to find a partition $\mathcal{P} = \{C_1, \dots, C_c\}$ of the nodes in graph G such that the proportion of edges inside C_k is high compared to the proportion of edges between C_k and other partitions.

Because our network sample is moderately large, we restrict our analysis to use more scalable community detection algorithms. We make use of the greedy modularity optimization algorithm [14] and the walktrap algorithm [15]. These algorithms are described in detail in Section III-D.

B. Content-Based Music Analysis

Many methods have been explored for content-based music analysis, attempting to characterize a music signal by its timbre, harmony, rhythm, or structure. One of the most widely used methods is the application of Mel-frequency cepstral coefficients (MFCC) to the modeling of timbre [16]. While a number of other spectral features have been used with success [17], when used in combination with various statistical techniques MFCCs have been successfully applied to music similarity and genre classification tasks [18]–[21].

A simple and fairly prevalent means to move from the high dimensional space of MFCCs to single similarity measure is to calculate the mean and covariance of each coefficient across an entire song and take the Euclidean distance between these mean and covariance sets (*e.g.* [19]). In the Music Information Retrieval Evaluation eXchange (MIREX) [22], [23] competitions of both 2007⁷ and 2009⁸, this method was shown to do a reasonable job of approximating human judgments of content-based similarity. A slightly more complex approach for computing timbre-based similarity between two songs

or collections of songs creates Gaussian Mixture Models (GMM) describing the MFCCs and comparing the GMMs using a statistical distance measure. Often the Earth Mover’s Distance (EMD), a technique first used in computer vision [24], is the distance measure used for this purpose [21], [25]. The EMD algorithm finds the minimum work required to transform one distribution into another.

C. Measuring Independence Between Distributions

When comparing social and acoustic similarity in this work, in addition to examining linear correlation via Pearson correlation, we will also find the mutual information contained across the social and acoustic similarity distributions. Taken from information theory, mutual information is the amount of dependence (usually measured in bits) that one distribution has on another (see for example [26]). Given two distributions X and Y , Mutual information $I(X; Y)$ can be defined as

$$I(X; Y) = H(X) - H(X|Y) \quad (1)$$

where $H(X)$ is the marginal entropy of the distributions X and Y and $H(X|Y)$ is the conditional entropy of X given Y .

All mutual information and related entropy calculations in this work are calculated using pyentropy⁹, a python library for performing information theoretic analysis on data distributions [27].

III. COMPLEX NETWORKS AND AUDIO CONTENT ANALYSIS

In this section, we report on our sampling of the Myspace network, describing our method in Section III-A and properties of our sample in Section III-B. We describe how to treat songs (rather than artists) as nodes in the graph in Section III-C, and describe our methods for community structure detection in Section III-D.

A. Sampling Myspace

The Myspace social network presents a variety of challenges. Firstly, its size prohibits analyzing the graph in its entirety, even when considering only the artist pages: therefore we sample a small yet sufficient portion of the network. Secondly, the Myspace social network is filled with noisy data – plagued by spammers and orphaned accounts: we limit the scope of our sampling in a way that minimizes this noise. Finally, there currently is no published interface for easily collecting the network

⁵<http://www.allmusic.com/>

⁶<http://www.lastfm.com/>

⁷see http://www.music-ir.org/mirex/2007/index.php/Audio_Music_Similarity_and_Retrieval_Results entry by G. Tzanetakis

⁸see http://www.music-ir.org/mirex/2009/index.php/Audio_Music_Similarity_and_Retrieval_Results entry by G. Tzanetakis

⁹<http://code.google.com/p/pyentropy/>

data from Myspace. Our data is collected using web crawling and HTML document scraping techniques¹⁰.

1) *Artist Pages*: It is important to note we are only concerned with a subset of the Myspace social network – the Myspace *artist* network. Myspace artist pages are different from standard Myspace pages in that they include a distinct audio player application containing material uploaded by that user. Standard practice (and a requirement of the End User License Agreement) is that this material has been generated by this user. We therefore use the presence or absence of this player to determine whether or not a given page is an artist page where, as stated in Section I, *artist page* is used to refer to the collection of social links and audio material assumed to be generated by the same person or group of people.

A Myspace page will include a top friends list. This is a hyperlinked list of other Myspace accounts explicitly specified by the user and, unlike generic friends, need not be a reciprocal relationship. The top friends list is limited in length with a maximum length of 40 friends (the default length is 16 friends). In constructing our sampled artist network, we use the top friends list to create a set of directed edges between artists. Only top friends who also have artist pages are added to the sampled network; standard Myspace pages are ignored. We also ignore the remainder of the friends list (*i.e.* friends that are not specified by the user as top friends), assuming these relationships are not as relevant. Our sampling method is based on the assumption that artists specified as top friends have some meaningful musical connection for the user – whether through collaboration, stylistic similarity, friendship, or artistic influence. This artificially limits the outdegree of each node in such a way as to only track social connections that have been selected by the artist to stand out, beyond the self-promoting noise of their complete friend list. Further, it is also a practical reduction as top friends can be scraped from the same single html document as all the other artist metadata. 50 friends are displayed per page, so gathering a full friend list would require $\frac{N}{50}$ pages to be scraped¹¹, significantly increasing the number of page requests required to sample the same number of artists.

Aside from these social connections metadata about the artist is gathered. This includes the name of the artist, number of page views and genre labels for the artists. The audio files associated with each artist page in the sampled network are also collected for feature extraction.

¹⁰Myspace scraping is done using tools from the MyPySpace project available at <http://mypyspace.sourceforge.net>

¹¹Where N is the number of friends, typically 10^3 but in some cases of the order 10^7 .

	n	m	$\langle k \rangle$	l	d_{max}
undirected	15478	91326	11.801	4.479	9
directed	15478	120487	15.569	6.426	16

TABLE I: The network statistics for the Myspace artist network sample where n is the number of nodes, m is the number of edges, $\langle k \rangle$ is the average degree, l is the mean geodesic distance, and d_{max} is the diameter, as defined in Section II-A1.

Note that while genre tags collected are at the level of artist; therefore all audio files associated with that artist will have the same genre labels applied.

2) *Snowball Sampling*: There are several network sampling methods; however, for the networks like the Myspace artist network, snowball sampling is the most appropriate method [28], [29]. In this method, the sample begins with a seed node (artist page), then the seed node’s neighbors (top friends), then the neighbors’ neighbors, are added to the sample. This breadth-first sampling is continued until the fraction of nodes in the sample reaches the target or *sampling ratio*. Here, we randomly select a seed artist¹² and collect all artist nodes within 6 edges to collect 15,478 nodes. If the size of the Myspace artist network is around 7 million, then this is close to the 0.25% sampling ratio suggested for accurate degree distribution estimation in sampled networks. Note that the sampling ratio is not sufficient for estimating other topological metrics such as the clustering coefficient and assortativity [30]; such global measures are not required for this paper.

With snowball sampling there is a tendency to over-sample hubs because they have many links and are easily picked up early in the breadth-first sampling. This effect reduces the degree distribution exponent by introducing a higher proportion of node with high connectivity then are seen in the complete network. This produces a heavier tail but preserving the overall power-law nature of the network [29].

B. Network Analysis of the Myspace Artist Network Sample

The Myspace artist network sample exhibits many of the network characteristics common to social networks and other real-world networks. Some of the network’s statistics are summarized in Table I.

We see that the MySpace artist network is like many other social networks in its “small world” characteristics - having a small diameter and geodesic distance.

¹²The artist is *Karna Zoo*, Myspace url: <http://www.myspace.com/index.cfm?fuseaction=user.viewProfile&friendID=134901208>

Additionally, in previous work, it has been shown that the Myspace artist network is assortative with respect to genre labels – that is, artists preferentially form connections with other artists that have the same genre labels [31].

Although the network is constructed as a directed network, for some of our experiments we convert to an undirected network to simplify analysis. This conversion is done to reduce complexity for analysis and to better examine the reflexive properties of our network. Each edge is considered bi-directional, that is $(i, j) = (j, i)$, and if a reflexive pair of edges existed in the directed graph, only one bi-directional edge exists in the undirected graph.

The degree distribution for this undirected reduction network is plotted in Figure 2 on a log-log scale. As mentioned earlier, it is common to find a power-law degree distribution in social networks [1]. However, exponential degree distributions have been reported previously in some types of music recommendation networks [9]. This is especially true for networks with imposed degree limits. For moderate degree values ($35 < k < 200$), our sample shows a power-law distribution. For lower degree values, the distribution is closer to exponential. This may be related to the fact that our network has an out degree limit imposed by Myspace restricting the maximum number of top friends ($k_{out} \leq 40$). The power-law fit also breaks down for high values of k – most likely due to the limited scope of our sample. Similar “broad-scale” degree distributions have been reported for citation networks and movie actor networks [32]. A more detailed analysis of this Myspace artist network can be found in [31].

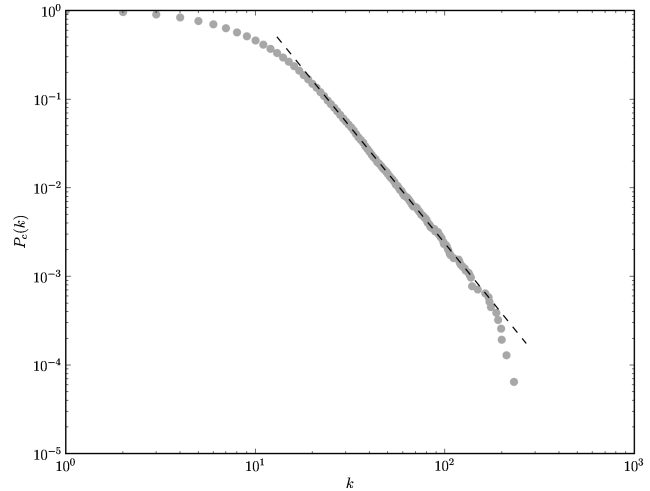


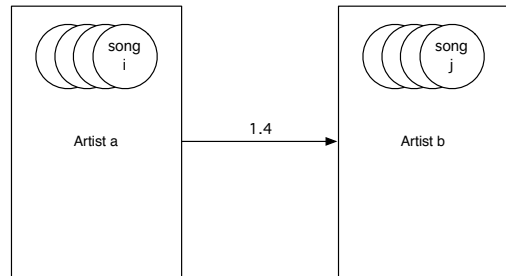
Fig. 2: The cumulative degree distributions for the Myspace artist network sample. For moderate values of k , the distribution follows a power-law (indicated by the dotted line), but for low and high values the decay is exponential.

C. Artists or Songs?

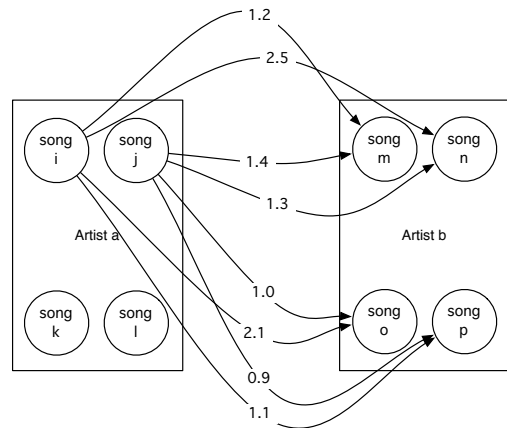
An unweighted graph between songs can be constructed by simply applying the artist connections to their associated songs; weights can be assigned to these song-to-song edges individually, for example based on acoustic dissimilarity between pairs of songs computed with the methods described in Section IV. These node relationships are illustrated in Figure 3.

D. Community Structure

We apply two community detection algorithms to our network sample – the greedy optimization of modularity [14] and the walktrap algorithm [15]. Both of these algorithms are reasonably efficient for networks of our size and both algorithms can be easily adapted to incorporate audio-based similarity measures. The work of this section discusses and extends the community structural analysis of [33] with further audio analysis.



(a) The sampled artist to artist relationship



(b) The expanded artist relationship, with songs as nodes. Note that the connections of song k and song l have been omitted for clarity.

Fig. 3: A comparison of sampled and song expanded means of representing the relationship between artists.

1) *Greedy Modularity Optimization*: Modularity is a network property that measures the appropriateness of a network division with respect to network structure. Modularity can be defined in several different ways [2]. In general, modularity Q is defined as the number of edges within communities minus the expected number of such edges. Let A_{ij} be an element of the network's adjacency matrix and suppose the nodes are divided into communities such that node i belongs to community C_i . We define modularity Q as the fraction of edges within communities minus the expected value of the same quantity for a random network. Then Q can be calculated as follows:

$$Q = \frac{1}{2m} \sum_{ij} \left[A_{ij} - \frac{d_i d_j}{2m} \right] \delta_{C_i C_j} \quad (2)$$

where the $\delta_{C_i C_j}$ function is 1 if $C_i = C_j$ and 0 otherwise, m is the number of edges in the graph, and d_i is the *degree* of node i – that is, the number of edges incident on node i . The sum of the term $\frac{d_i d_j}{2m}$ over all node pairs in a community represents the expected fraction of edges within that community in an equivalent random network where node degree values are preserved.

If we consider Q to be a benefit function we wish to maximize, we can then use an agglomerative approach to detect communities – starting with a community for each node such that the number of partitions $|\mathcal{P}| = n$ and building communities by amalgamation. The algorithm is greedy, finding the changes in Q that would result from the merge of each pair of communities, choosing the merge that results in the largest increase of Q , and then performing the corresponding community merge. It can be proven that if no community merge will increase Q the algorithm can be stopped because no further modularity optimization is possible [14]. Using efficient data structures based on sparse matrices, this algorithm can be performed in time $\mathcal{O}(m \log n)$.

2) *Random Walk: Walktrap*: The walktrap algorithm uses random walks on G to identify communities. Because communities are more densely connected, a random walk will tend to be ‘trapped’ inside a community – hence the name “walktrap”.

At each time step in the random walk, the walker is at a node and moves to another node chosen randomly and uniformly from its neighbors. The sequence of visited nodes is a *Markov chain* where the states are the nodes of G . At each step the transition probability from node i to node j is $P_{ij} = \frac{A_{ij}}{d_i}$ which is an element of the transition matrix P for the random walk. We can also write $P = D^{-1}A$ where D is the diagonal matrix of the degrees ($\forall i, D_{ii} = d_i$ and $D_{ij} = 0$ where $i \neq j$).

The random walk process is driven by powers of P : the probability of going from i to j in a random walk of length t is $(P^t)_{ij}$ which we will denote simply as P_{ij}^t . All of the transition probabilities related to node i are contained in the i^{th} row of P^t denoted as $P_{i\bullet}^t$. We then define an inter-node distance measure:

$$r_{ij} = \sqrt{\sum_{k=1}^n \frac{(P_{ik}^t - P_{jk}^t)^2}{d_k}} = \|D^{-\frac{1}{2}} P_{i\bullet}^t - D^{-\frac{1}{2}} P_{j\bullet}^t\| \quad (3)$$

where $\|\cdot\|$ is the Euclidean norm of \mathbb{R}^n . This distance can also be generalized as a distance between communities: $r_{C_i C_j}$ or as a distance between a community and a node: $r_{C_i j}$.

We then use this distance measure in our algorithm. Again, the algorithm uses an agglomerative approach, beginning with one partition for each node ($|\mathcal{P}| = n$). We first compute the distances for all adjacent communities (or nodes in the first step). At each step k , two communities are chosen based on the minimization of the mean σ_k of the squared distances between each node and its community.

$$\sigma_k = \frac{1}{n} \sum_{C_i \in \mathcal{P}_k} \sum_{i \in C_i} r_{i C_i}^2 \quad (4)$$

Direct calculation of this quantity is known to be NP-hard [15], so instead we calculate the variations $\Delta\sigma_k$. Because the algorithm uses a Euclidean distance, we can efficiently calculate these variations as

$$\Delta\sigma(C_1, C_2) = \frac{1}{n} \frac{|C_1||C_2|}{|C_1| + |C_2|} r_{C_1 C_2}^2 \quad (5)$$

The community merge that results in the lowest $\Delta\sigma$ is performed. We then update our transition probability matrix

$$P_{(C_1 \cup C_2)\bullet}^t = \frac{|C_1| P_{C_1\bullet}^t + |C_2| P_{C_2\bullet}^t}{|C_1| + |C_2|} \quad (6)$$

and repeat the process updating the values of r and $\Delta\sigma$ then performing the next merge. After $n - 1$ steps, we get one partition that includes all the nodes of the network $\mathcal{P}_n = \{N\}$. The algorithm creates a sequence of partitions $(\mathcal{P}_k)_{1 \leq k \leq n}$. Finally, we use modularity to select the best partition of the network, calculating $Q_{\mathcal{P}_k}$ for each partition and selecting the partition that maximizes modularity.

Because the value of t is generally low (we use $t = 4$), this community detection algorithm is quite scalable. For most real-world networks, where the graph is sparse, this algorithm runs in time $\mathcal{O}(n^2 \log n)$ [15]. Note though, the optimized greedy modularity algorithm is considerably faster than the walktrap algorithm – $\mathcal{O}(m \log n)$ versus $\mathcal{O}(n^2 \log n)$

IV. HYBRID METHODS OF SIMILARITY ANALYSIS

In this section we will discuss and extend the work of [34], [35], applying additional acoustic similarity measures and expanding our analysis of resulting distributions of similarity measurements. The geodesic distance between all pairs of artists within the sample are compared to the acoustic similarity of songs associated with each artist. Maximum flow analysis is subsequently used to further explore this artist social space. This measure is compared to the same artist based acoustic similarity and an additional song to song acoustic metric is used. Lastly, community segmentation and structural analysis are explored as a further means of understanding the interaction between these two spaces.

MFCCs are extracted from each audio signal using a Hamming window on 8192 sample FFT windows with 4096 sample overlap. All MFCCs are created with the `fftExtract` tool¹³. For each artist node a GMM is built from the concatenation of MFCC frames for all songs found on each artist's Myspace page. Generally artists have between 1 and 4 songs, although some artists have many more. The mean number of songs is slightly more than 3.5 per artist. An $n \times n$ matrix is populated with the earth mover's distance λ_{ij} between the GMMs corresponding to each pair of nodes in the sample. As a second acoustic dissimilarity measure, the software suite Marsyas¹⁴ is used in the exact configuration that was used in the MIREX 2009 Audio Similarity and Retrieval¹⁵ task to generate MFCC-based average value vectors per song and then to generate an $N \times N$ euclidean distance matrix of these songs. These distance matrices are used to draw λ values to compare against the song expanded graph as detailed in Section III-C.

A. Geodesic Distance

The relation between audio signal dissimilarity and the geodesic path length is first examined using a box and whisker plot. The plot is shown in Figure 4. These dissimilarities are grouped according to the geodesic distance in the undirected network between the artist nodes i and j , d_{ij} . There appears to be no clear correlation between these λ values and geodesic distance. The Pearson product-moment correlation coefficient confirms this giving a ρ of -0.0016 . This should be viewed in the context of the number of pairwise relationships used, implying it is stable, at least for the community of artists found via this sample of the network.

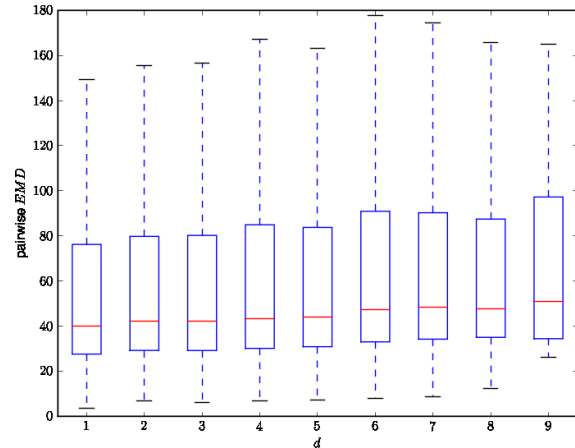


Fig. 4: The box and whisker plot showing the spread of pair-wise artist dissimilarity grouped by geodesic distance as found on the artist graph. The whiskers cover the second and seventh eighths beyond the inner quartiles covered in each box.

B. Maximum Flow

In our Myspace top friends graph, the maximum flow is measured on the directed and undirected reduction of the unweighted graph from the source artist node to the sink artist node.

1) *Experiment*: The maximum flow value is calculated, using the snowball sample entry point as the fixed source against every other node in turn as a sink, yielding the number of edges connecting each sink node to the entry point node at the narrowest point of connection. The acoustic distances are then be compared to these maximum flow values.

In order to better understand a result from analysis of our Myspace sample, a baseline for comparison must be used. To that end, random permutations of the node locations are examined. In order to preserve the overall topology present in the network, this randomization is performed by randomizing the artist label and associated music attached to a given node on the network. This is done ten fold, creating a solid baseline to test the null hypothesis that the underlining community structure is not responsible for any correlation between maximum flow values and λ_{ij} from either of the two acoustic dissimilarity measures.

2) *Results*: The results of this experiment show no simple relationship between the sampled network and the randomized network. This can be seen in Table II and in Figures 5 and 6. There is an increase in the median EMD for the less well connected (*i.e.* lower maximum flow value) node pairs in the Myspace sample graph, though

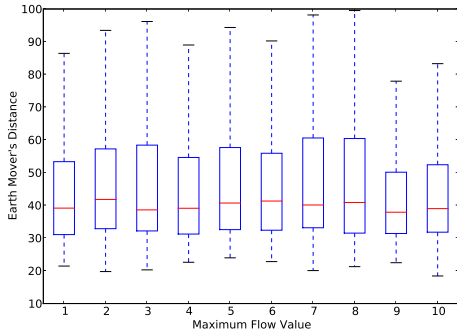
¹³source code at <http://omras2.doc.gold.ac.uk/software/fftextextract/>

¹⁴<http://marsyas.info/>

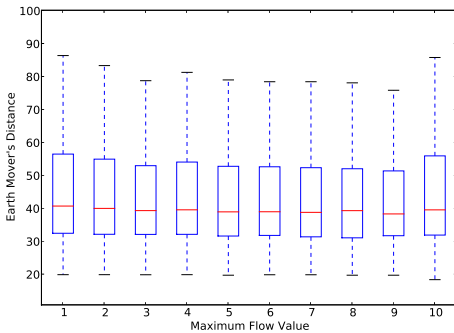
¹⁵<http://music-ir.org/mirex/2009/results/abs/GTfinal.pdf>

	H-value	P-value
From sample	12.46	0.19
Random permutations	9.11	0.43

TABLE III: The Kruskal-Wallis one-way ANOVA test results of EMD against maximum flow for both the sampled graph and its random permutations. The H-values are drawn from a chi-square distribution with 10 degrees of freedom.



(a) The EMD distribution on the sampled graph

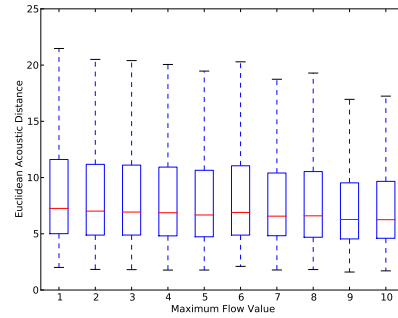


(b) The EMD distribution on the random permutations of the graph, maintaining the original edge structure.

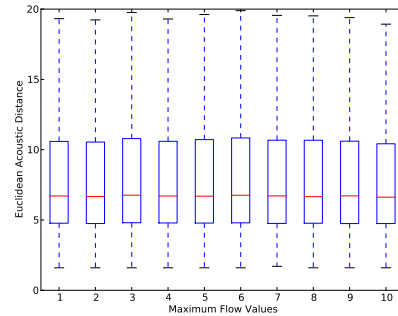
Fig. 5: The box and whisker plots showing the distribution of EMD grouped by maximum flow value between artists on the Myspace social graph and the randomized permutations of the graph.

this is not significant enough to indicate a correlation, while the randomized permutations are near flat. Perhaps the easiest way to examine the relationship between the sampled graph and randomized one is through the deltas of each group’s median from the entire dataset median. This data is shown in the second and fourth column in Table II and Figure 7. Further, the Kruskal-Wallis one-way ANOVA results for both the sample graph and averaged across the 10 fold permutations are shown in Table III.

Additionally, we calculated the mutual information be-



(a) The Euclidean distance distribution on the sampled graph



(b) The Euclidean distance distribution on the random permutations of the graph, maintaining the original edge structure.

Fig. 6: The box and whisker plots showing the distribution of Euclidean distance grouped by maximum flow value between artists on the Myspace social graph and the randomized permutations of the graph.

audio distance type	$H(X)$	$H(X Y)$	$H(Y)$	$I(X; Y)$
Euclidean distance	3.100	3.00	8.65	0.100
GMM/EMD	3.098	2.723	8.65	0.375

TABLE IV: Entropy values for the acoustic distances and maximum flow values. X is the set of audio distance measurements, Y is the set of maximum flow values.

tween the flow network and both of the acoustic distance measures. These can be seen, along with the entropy of each set in Table IV. Here we can beyond simply looking at an implied near independence. The mutual information between the maximum flow values and either of two acoustic distance measure is a small fraction of the entropy of either respective set of distances.

C. Using Audio in Community Detection

Both community detection algorithms described in Section III-D are based on the adjacency matrix A of the graph. This allows us to easily extend these algorithms

Max Flow	Earth Movers Distance				Marsyas generated Euclidean Distance			
	median	deviation	randomized	deviation	median	deviation	randomized	deviation
1	40.80	1.26	39.10	-0.43	7.256	0.571	6.710	0.025
2	45.30	5.76	38.34	-1.19	7.016	0.331	6.668	-0.016
3	38.18	-1.35	38.87	-0.66	6.932	0.247	6.764	0.079
4	38.21	-1.32	38.64	-0.89	6.872	0.187	6.707	0.022
5	40.00	0.47	39.11	-0.42	6.673	-0.011	6.695	0.010
6	41.77	2.25	39.02	-0.51	6.896	0.211	6.761	0.076
7	39.94	0.41	39.24	-0.29	6.568	-0.116	6.714	0.029
8	39.38	-0.15	38.76	-0.77	6.597	-0.087	6.660	-0.023
9	38.50	-1.03	38.87	-0.66	6.270	-0.414	6.717	0.032
10	39.07	-0.46	40.85	1.32	6.253	-0.431	6.623	-0.061

TABLE II: Node pairs of median acoustic distance values grouped by actual minimum cut values and randomized minimum cut values, shown with deviations from the global medians of 39.53 for EMD and 6.6848 for Euclidean distance. EMD weights are on the left and Euclidean distances as generated by Marsyas are on the right.

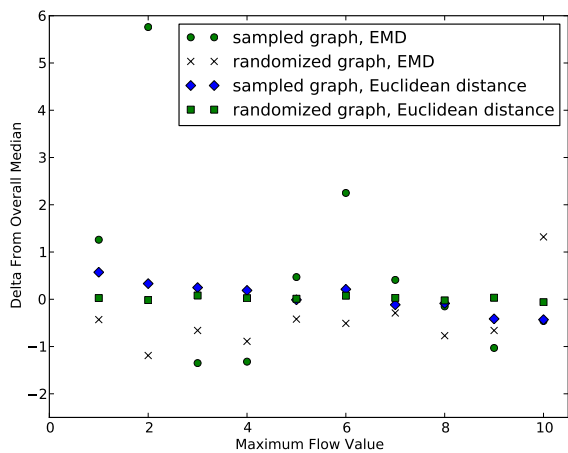


Fig. 7: The deltas from the global median for each maximum flow value group of acoustic distance values, from the sampled graph and the randomized graph.

to include audio-based similarity measures. We simply insert an inter-node similarity value for each non-zero entry in A . We calculate these similarity values using both audio-based analysis methods detailed in Section IV.

These dissimilarity values must be converted to similarity values to be successfully applied to the community detection algorithms. We do this by taking the reciprocal of each dissimilarity.

$$A_{ij} = \begin{cases} \lambda_{ij}^{-1} & \text{if nodes } i \text{ and } j \text{ are connected,} \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

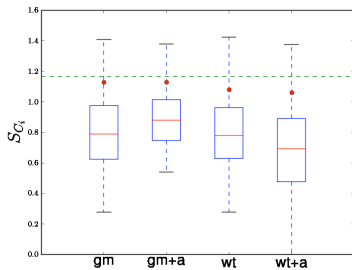
1) *Genre Entropy*: Now that we have several methods for detecting community structures in our network, we need a means of evaluating the relevance of these structures in the context of music. Traditionally, music

and music artists are classified in terms of *genre*. If the structure of the Myspace artist network is relevant to music, we would expect the communities identified within the network to be correlated with musical genres. That is, communities should contain nodes with a more homogenous set of genre associations than the network as a whole.

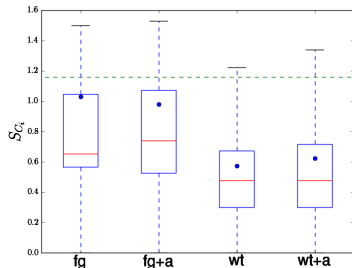
As mentioned in our description of our sampling of the Myspace network (Section III-A), we have collected genre tags that are associated with each artist. In order to measure the diversity of each community with respect to genre we use a variant of Shannon entropy we call *genre entropy* S . This approach is similar to that of Lambiotte [13]. For a given community C_k we calculate genre entropy as:

$$S_{C_k} = - \sum_{\gamma \in C_k} P_{\gamma|C_k} \log P_{\gamma|C_k} \quad (8)$$

where $P_{\gamma|C_k}$ is the probability of finding genre tag γ in community C_k . As the diversity of genre tags in a community C_k increases, the genre entropy S_{C_k} increases. As the genre tags become more homogenous, the value of S_{C_k} decreases. If community C_k is described entirely by one genre tag then $S_{C_k} = 0$. We can calculate an overall genre entropy S_G by including the entire network sample. In this way, we can evaluate each community identified by comparing S_{C_k} to S_G . If the community structures in the network are related to musical genre, we would expect the communities to contain more homogenous mixtures of genre tags. That is, usually, we would expect $S_{C_k} \leq S_G$. However, as community size decreases the genre entropy will tend to decrease because fewer tags are available. To account for this, we create a random partitioning of the graph that results in the same number of communities with the same number of nodes in each community and calculate



(a) Audio weights are Earth Mover's Distance



(b) Audio weights are Euclidean distance generated by Marsyas

Fig. 8: Box and whisker plots showing the spread of community genre entropies for each graph partitioning method where gm is greedy modularity, gm+a is greedy modularity with audio weights, wt is walktrap, and wt+a is walktrap with audio weights. The horizontal line represents the genre entropy of the entire sample. The circles represent the average value of genre entropy for a random partition of the network into an equivalent number of communities. (a) uses the Earth Mover's Distance for audio weight, (b) uses Euclidean distance from Marsyas.

the corresponding genre entropies S_{rand} to provide a baseline.

If an artist specified no genre tags, this node is ignored and makes no contribution to the genre entropy score. In our data set, 2.6% of artists specified no genre tags.

2) *Results*: The results of the various community detection algorithms are summarized in Figure 1 and Table 1. When the genre entropies are averaged across all the detected communities, we see that for every community detection method the average genre entropy is lower than S_G as well as lower than the average genre entropy for a random partition of the graph into an equal number of communities. This is strong evidence that the community structure of the network is related to musical genre.

It should be noted that even a very simple examination of the genre distributions for the entire network sample

algorithm	c	$\langle S_C \rangle$	$\langle S_{rand} \rangle$	Q
none	1	1.16	-	-
gm	42	0.81	1.13	0.61
gm+a	33	0.90	1.13	0.64
wt	195	0.80	1.08	0.61
wt+a	271	0.70	1.06	0.62

TABLE V: Results of the community detection algorithms where c is the number of communities detected, $\langle S_C \rangle$ is the average genre entropy for all communities, $\langle S_{rand} \rangle$ is the average genre entropy for a random partition of the network into an equal number of communities, and Q is the modularity for the given partition.

suggests a network structure that is closely related to musical genre. Of all the genre associations collected for our data set, 50.3% of the tags were either “Hip-Hop” or “Rap” while 11.4% of tags were “R&B”. Smaller informal network samples, independent of our main data set, were also dominated by a handful of similar genre tags (*i.e.* “Alternative”, “Indie”, “Punk”). In context, this suggests our sample was essentially “stuck” in a community of Myspace artists associated with these particular genre inclinations. However, it is possible that these genre distributions are indicative of the entire Myspace artist network. Regardless, given that the genre entropy of our entire set is so low to begin with it is an encouraging result that we could efficiently identify communities of artists with even lower genre entropies.

Without audio-based similarity weighting, the greedy modularity algorithm (gm) and the walktrap algorithm (wt) result in genre entropy distributions with no statistically significant differences. However the walktrap algorithm results in almost five times as many communities which we would expect to result in a lower genre entropies because of smaller community size. Also note that as discussed in Section III-D the optimized greedy modularity algorithm is considerably faster than the walktrap algorithm.

With audio-based similarity weighting, we see mixed results. Applying audio weights to the greedy modularity algorithm (fg+a) actually increased genre entropies but the differences between fg and fg+a genre entropy distributions are not statistically significant. Audio-based weighting applied to the walktrap algorithm (wt+a) results in a statistically significant decrease in genre entropies compared to the un-weighted walktrap algorithm ($p = 4.2 \times 10^{-4}$). It should be noted that our approach to audio-based similarity results in dissimilarity measures that are mostly orthogonal to network structure [34].

V. CONCLUSIONS AND FUTURE WORK

We have presented an analysis of the community structures found in a sample of the Myspace artist network. The communities detected have lower entropy over genre labels than a graph with randomly permuted labels. We have applied two efficient algorithms to the task of partitioning the Myspace artist network sample into communities and we have shown how to include audio-based similarity measures in the community detection process. We have evaluated our results in terms of genre entropy - a measure of genre tag distributions - and shown the community structures in the Myspace artist network are related to musical genre.

We compared social space of the Myspace sample with content-based acoustic space in two ways in Section IV. First the geodesic distances of pairs of artists were compared to the acoustic distance between these pairs of artists. Then maximum flow between pairs of artists was compared to both the acoustic distance between the artists and amongst the artists' songs. While not perfectly orthogonal, the artist social graph and the acoustic dissimilarity matrix clearly encode different relational aspects between artists. This can be clearly seen in the small amount of mutual information seem between the sets of distances. The implication is that using both of these spaces in applications driven by similarity measures will result in much higher entropy in the data available to such an application. This suggests that a recommendation or discovery system that can use both domains well has the potential to perform much better than a similar system that relies on only one domain in isolation.

Furthermore, while an inverse relationship between Earth Mover's Distance and the maximum flow value might be expected on the basis of the conventional wisdom that a community of artists tend to be somehow aurally similar, this does not appear to be the case. The evidence, at least in this sample set, does not support this relationship. However, based upon the difference in result from the Kruskal-Wallis one-way ANOVA test and simple observation of the deviation from the global median the maximum flow values and Earth Mover's Distances do seem affected by the artist created social links, though it is not a simple relationship and describing it precisely is difficult.

Because the Myspace artist network might be of interest to other researchers, we have converted our graph data to a more structured format. We have created a Web service¹⁶ that describes any Myspace page in a machine-

readable Semantic Web format. Using FOAF¹⁷ and the Music Ontology¹⁸ [36], the service describes a Myspace page in XML RDF. This will allow future applications to easily make use of Myspace network data (*e.g.* for music recommendation).

While it is unclear how to best to use all the available information from the wide range of artists and musicians, what this work makes clear is that there are advantages to complex multi-domain notions of similarity in music. By using both acoustic and social data recommender systems have more avenues to pursue to present new material to users in a transparent way. Whether either of these spaces can provide insight into the other remains an open question, though our work tend to show the likely predictability of one space from the other is low. In spite or perhaps because of this separation, and given the sheer quantity of data available on the web, it seems inevitable that these domains will be used in tandem in future music recommendation and musicological study.

A. Further Investigation into Artist Networks

In future work we plan to examine community detection methods that operate locally, without knowledge of the entire network. We also plan to address further directed artist graph analysis, bipartite networks of artists and listeners, different audio analysis methods, and the application of these methods to music recommendation.

Many of these tasks require the expansion of our sample network. The goal of any effort to expand the sample size of a network such as Myspace is best focused on ways to make the sample set more indicative of the whole. While it is impossible to assess this without capturing the entire graphs some assumptions can be made. Snowball sampling has a tendency to oversample hubs. Given this, a better expanded network is likely to result through the selections of new starting seed artist (most likely at random) and proceeding via a breadth-first crawl until that crawl results in overlap with the known network. It is reasonable to assume that this method, when used over multiple hubs, will produce a lower proportion of high centrality hubs than simply continuing further with the existing breadth first crawl. With a lower proportion of these over-sampled hubs, the social structure of the sample would better match that of the whole.

B. Playlist-Based Applications

1) *The Max Flow Playlist*: In order to build playlists using both acoustic and social network data, the Earth

¹⁶available at <http://dbtune.org/myspace>

¹⁷<http://www.foaf-project.org/>

¹⁸<http://musicontology.com/>

Mover's Distance is used between each pair of neighbors as weights on the Myspace sample network. Two artists are then selected, a starting artist as the source node and a final artist as the sink node. One or more paths are then found through the graph via the maximum flow value, generating the list and order of artists for the playlist. The song used for each artist is the most popular at the time of the page scrape. In this way playlists are constructed that are influenced both by timbre similarity and bound by social context, regardless of any relationship found between these two spaces found via the work discussed in Section IV. Playlists generated using this technique were informally auditioned, and were found to be reasonable on that basis.

There is clearly potential in the idea of the maximum flow playlist. When using either audio similarity measure as a weight the results appear to be quite good, at least from a qualitative perspective. The imposed constraint of the social network alleviates to some extent short comings of a playlist built purely through the analysis of acoustic similarity by moving more toward the balance between completely similar works and completely random movement.

2) *Steerable Optimized Self-Organizing Radio*: Using the song-centric graph the following system is proposed as a means of deployment. This system will be designed to play a continuous stream of songs via an internet radio stream. The playback system will begin with an initial, seed song and destination song, then construct a playlist. While this playlist is being broadcast, anyone tuning into the broadcast will also be able to vote via a web based application on the next song to serve as the destination. In order to produce a usable output the vote system would present a list of *nominees* selected as a representative track from various communities as segregated via means discussed in Section IV-C.

Once the current destination song begins broadcast, the voting for the next cycle will cease. The current destination song will be considered the seed song for the next cycle and the song with a plurality of votes will become the new destination, then the next playlist will be calculated and its members broadcast. This process will continue for the duration of the broadcast. If this automatic playlist creation system is allowed to run for a sufficient amount of time, a great deal of user data will be recorded. This would include voting behavior, average length of time continuously listened and whether listeners (or at least IP addresses) return. This provides a built-in means of human listener evaluation for these playlists.

It is hoped that this system, or one like it, will provide an application driven means to evaluate the usability of

the measures explored in this work in task of music discovery and recommendation.

ACKNOWLEDGMENT

This work is supported in part by the Engineering and Physical Sciences Research Council via the Online Music Recognition And Searching II (OMRAS2) project, reference numbers EP/E017614/1 and EP/E02274X/1. Additional support as part of the Networked Environments for Music Analysis (NEMA) project, funded by The Andrew W. Mellon Foundation. Also, thanks are due to our anonymous reviewers, their comments have helped strengthen this work significantly.

REFERENCES

- [1] M. E. J. Newman, "The structure and function of complex networks," *SIAM Review*, vol. 45, p. 167, 2003. [Online]. Available: <http://www.citebase.org/abstract?id=oai:arXiv.org:cond-mat/0303516>
- [2] L. F. Costa, F. A. Rodrigues, G. Travieso, and P. R. V. Boas, "Characterization of complex networks: A survey of measurements," *Advances In Physics*, vol. 56, p. 167, 2007. [Online]. Available: doi:10.1080/00018730601170527
- [3] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin, *Network flows: theory, algorithms, and applications*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1993.
- [4] H. Nagamochi and T. Ibaraki, "Computing edge-connectivity in multigraphs and capacitated graphs," *SIAM J. Discret. Math.*, vol. 5, no. 1, pp. 54–66, 1992.
- [5] P. Elias, A. Feinstein, and C. Shannon, "A note on the maximum flow through a network," *Information Theory, IEEE Transactions on*, vol. 2, no. 4, pp. 117–119, Dec 1956.
- [6] A. V. Goldberg and R. E. Tarjan, "A new approach to the maximum-flow problem," *J. ACM*, vol. 35, no. 4, pp. 921–940, 1988.
- [7] P. Knees, T. Pohle, M. Schedl, and G. Widmer, "Combining audio-based similarity with web-based data to accelerate automatic music playlist generation," in *Proc. 8th ACM international workshop on Multimedia information retrieval*, 2006, pp. 147 – 154.
- [8] M. Alghoniemy and A. Tewfik, "A network flow model for playlist generation," in *Multimedia and Expo, 2001. IEEE International Conference on*, 2001. [Online]. Available: citeseer.ist.psu.edu/alghoniemy01network.html
- [9] P. Cano, O. Celma, M. Koppenberger, and J. M. Buldu, "The topology of music recommendation networks," *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 2006, <http://arxiv.org/abs/physics/0512266v1>. [Online]. Available: <http://www.citebase.org/abstract?id=oai:arXiv.org:physics/0512266>
- [10] P. Gleiser and L. Danon, "Community structure in jazz," *Advances in Complex Systems*, vol. 6, no. 4, pp. 565–573, 2003.
- [11] J. Park, O. Celma, M. Koppenberger, P. Cano, and J. M. Buldu, "The social network of contemporary popular musicians," *Int. J. of Bifurcation and Chaos*, vol. 17, pp. 2281–2288, 2007. [Online]. Available: <http://www.citebase.org/abstract?id=oai:arXiv.org:physics/0609229>
- [12] A. Anglade, M. Tiemann, and F. Vignoli, "Virtual communities for creating shared music channels," in *Proc. of Int. Symposium on Music Information Retrieval*, 2007.

- [13] R. Lambiotte and M. Ausloos, "On the genre-fication of music: a percolation approach (long version)," *The European Physical Journal B*, vol. 50, p. 183, 2006. [Online]. Available: doi:10.1140/epjb/e2006-00115-0
- [14] A. Clauset, M. E. J. Newman, and C. Moore, "Finding community structure in very large networks," *Phys. Rev. E*, vol. 70, no. 6, p. 066111, Dec 2004.
- [15] P. Pons and M. Latapy, "Computing communities in large networks using random walks (long version)," arXiv:physics/0512106v1, 2005. [Online]. Available: <http://www.citebase.org/abstract?id=oai:arXiv.org:physics/0512106>
- [16] B. Logan, "Mel frequency cepstral coefficients for music modeling," in *Int. Symposium on Music Information Retrieval*, 2000.
- [17] M. Casey, R. Veltkamp, M. Goto, M. Leman, C. Rhodes, and M. Slaney, "Content-based music information retrieval: Current directions and future challenge," *Proc. of the IEEE*, vol. 96, no. 4, pp. 668–696, 2008.
- [18] A. Berenzweig, B. Logan, D. P. W. Ellis, and B. P. W. Whitman, "A large-scale evaluation of acoustic and subjective music-similarity measures," *Computer Music J.*, vol. 28, no. 2, pp. 63–76, 2004.
- [19] G. Tzanetakis, *Marsyas: a case study in implementing Music Information Retrieval Systems*. Information Science Reference, 2007. [Online]. Available: http://marsyas.sness.net/pdfs/0000/0007/imis_bookchapter.pdf
- [20] B. Logan and A. Salomon, "A music similarity function based on signal analysis," *Multimedia and Expo, 2001. ICME 2001. IEEE International Conference on*, pp. 745–748, 22-25 Aug. 2001.
- [21] E. Pampalk, "Computational models of music similarity and their application in music information retrieval," Ph.D. dissertation, Technischen Universität Wien, May 2006.
- [22] J. S. Downie, "The music information retrieval evaluation exchange (mirex)," *D-Lib Magazine*, Dec. 2006. [Online]. Available: <http://dlib.org/dlib/december06/downie/12downie.html>
- [23] —, "The music information retrieval evaluation exchange (2005–2007): A window into music information retrieval research," *Acoustical Science and Technology*, vol. 29, no. 4, pp. 247–255, 2008. [Online]. Available: http://www.jstage.jst.go.jp/article/ast/29/4/29_247/_article
- [24] Y. Rubner, C. Tomasi, and L. J. Guibas, "The earth mover's distance as a metric for image retrieval," *International Journal of Computer Vision*, vol. 40, no. 2, pp. 99–121, 2000.
- [25] J. Aucouturier and F. Pachet, "Improving timbre similarity : How high's the sky ?" *J. Negative Results in Speech and Audio Sciences*, 2004.
- [26] R. Steuer, J. Kurths, C. Daub, J. Weise, and J. Selbig, "The mutual information: Detecting and evaluating dependencies between variables," *Bioinformatics*, vol. 18 Suppl.2, pp. S231–S240, 2002. [Online]. Available: <http://direct.bl.uk/bld/PlaceOrder.do?UIN=132115618&ETOC=RN&from=searchengine>
- [27] R. A. Ince, R. S. Petersen, D. C. Swan, and S. Panzeri, "Python for information theoretic analysis of neural data," *Front Neuroinformatics*, vol. 3, pp. 4–4, 2009. [Online]. Available: <http://www.hubmed.org/display.cgi?uids=19242557>
- [28] Y.-Y. Ahn, S. Han, H. Kwak, S. Moon, and H. Jeong, "Analysis of topological characteristics of huge online social networking services," in *Proceedings of the 16th international conference on World Wide Web*. Alberta, Canada: IW3C2, May 2007.
- [29] S. H. Lee, P.-J. Kim, and H. Jeong, "Statistical properties of sampled networks," *Physical Review E*, vol. 73, pp. 102–109, January 2006.
- [30] H. Kwak, S. Han, Y.-Y. Ahn, S. Moon, and H. Jeong, "Impact of snowball sampling ratios on network characteristics estimation: A case study of cyworld," KAIST, Tech. Rep. CS/TR-2006-262, November 2006.
- [31] K. Jacobson and M. Sandler, "Musically meaningful or just noise, an analysis of on-line artist networks," in *Proc. of CMMR*, 2008, pp. 306–314.
- [32] L. A. N. Amaral, A. Scala, M. Barthélemy, and H. E. Stanley, "Classes of small-world networks," in *Proceeding of the National Academy of Sciences*, 2000.
- [33] K. Jacobson, B. Fields, and M. Sandler, "Using audio analysis and network structure to identify communities in on-line social networks of artists," in *Proc. of Int. Symposium on Music Information Retrieval*, 2008.
- [34] B. Fields, K. Jacobson, M. Casey, and M. Sandler, "Do you sound like your friends? exploring artist similarity via artist social network relationships and audio signal processing," in *Int. Computer Music Conference*, August 2008.
- [35] B. Fields, K. Jacobson, C. Rhodes, and M. Casey, "Social playlists and bottleneck measurements : Exploiting musician social graphs using content-based dissimilarity and pairwise maximum flow values," in *Proc. of Int. Symposium on Music Information Retrieval*, September 2008.
- [36] Y. Raimond, S. Abdallah, M. Sandler, and F. Gaisson, "The music ontology," in *Int. Conference on Music Information Retrieval*, 2007.