

# A Semantic Map of the last.fm Music Folksonomy

Biberstine, J.<sup>1</sup>, Börner, K.<sup>1</sup>, Duhon, R. J.<sup>1</sup>, Hardy, E.<sup>1</sup>, Skupin, A.<sup>2</sup>

<sup>1</sup>School of Library and Information Science, Indiana University, Bloomington, IN 47405

<sup>2</sup>Department of Geography, San Diego State University, San Diego, CA 92040, USA  
Email: skupin@mail.sdsu.edu

## 1. Introduction

What does the world of music look like? We present a visualization derived from a repository of user-generated tags attached to more than one million items within the social music website last.fm (<http://www.last.fm/>). The site enables users to discover new music based on their listening history and – crucial for our study – users can annotate music-related items such as artists and songs with arbitrary tags, ranging from categories like “rock” or “jazz” to event-related attributes (e.g. “seen live”) and affective utterances (e.g., “songs I absolutely love”). Tags also vary in scope, from rather broad categories, like “classical”, to finer distinctions, like “britpop” or “female fronted metal”.

The map offers viewers a mix of recognition, surprise, and discovery. Viewers have reported appreciation of the coherent patterns of hierarchical relationships among musical styles. The map is also notable for offering opportunities for discovering new musical categories, from the various flavors of “metal” to such niche areas as “shoegaze” or “drone”.

The project presented here advances the visual analysis of social media on a number of fronts. Together with a recent study visualizing two million biomedical documents (Boyack et al. 2011), which generated several key technology elements used in our project, this is to our knowledge one of the largest self-organizing maps ever created in a single process (e.g., not hierarchically built, unlike Kohonen et al. 1999), and certainly one of the most comprehensive semantic depictions of the world of music. A project of this scale would not have been possible without parallelization and supercomputing resources.

GIScience concepts are at the core of how the project advances the analysis of social media data, which has traditionally been dominated by a network paradigm. Instead, our use of SOM builds on a view of social media items as existing in a high-dimensional attribute space, which calls for very different conceptual and computational approaches. Such geographic and mapping metaphors as *scale*, *region*, or *base map* feature prominently in the conceptual design of this visualization, while its rendering leverages commercial off-the-shelf GIS technology.

## 2. Data and Methods

The original data set, collected during the first half of 2009 (Schifanella et al. 2010) contains almost 1.4 million different music-related items, to which over 280,000 different tags were attached. In a vector-space type framework (Salton 1989), this would amount to a space that is extremely high-dimensional and yet occupied by exceedingly sparse vectors. That provides a justification for reducing the dimensionality by removing all but the 1,000 most frequently used tags. Items not annotated with any of those remaining tags are then deleted. This reduces the number of items from 1,393,559 to 1,088,761, with the average item having 6.8 overall tags and 3.8 unique tags associated with it. The top twenty most popular tags were: rock, electronic, seen live, indie, alternative, pop, female vocalists, jazz, classic rock, experimental, ambient, metal, alternative rock, singer-songwriter, 80s, folk, hard rock, progressive rock, indie rock, electronica, punk.





Efforts like the one described here can benefit from further integration of GIScience concepts and principles. An example is the introduction of a continuous field view of social media items and activities, as an alternative to current conceptualizations that seem predominantly informed by an information retrieval tradition (Skupin 2009).

## Acknowledgements

This work was funded by the Cyberinfrastructure for Network Science Center at Indiana University, the James S. McDonnell Foundation, and the National Science Foundation under grant SBE-0738111. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation. We also gratefully acknowledge the feedback received from three anonymous reviewers whose comments helped to improve this extended abstract.

## References

- Blei, D.M., Ng, A.Y. and Jordan, M.I., 2003, Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3, pp. 993-1022.
- Boyack, K.W., Newman, D., Duhon, R.J., Klavans, R., Patek, M., Biberstine, J.R., Schijvenaars, B., Skupin, A., Ma, N.A.L. and Börner, K., 2011, Clustering More than Two Million Biomedical Publications: Comparing the Accuracies of Nine Text-Based Similarity Approaches. In *PLoS ONE*, p. e18029.
- Kohonen, T., Kaski, S., Lagus, K., Salojärvi, J., Honkela, T., Paatero, V. and Saarela, A., 1999, Self Organization of a Massive Text Document Collection. In *Kohonen Maps*, E. Oja and S. Kaski (Eds.), pp. 171-182 (Amsterdam: Elsevier).
- Kohonen, T., Hynninen, J., Kangas, J. and Laaksonen, J., 1996, *SOM\_PAK: The Self-Organizing Map Program Package* (Espoo, Finland: Helsinki University of Technology, Laboratory of Computer and Information Science).
- Lawrence, R.D., Almasi, G.S. and Rushmeier, H.E., 1999, A scalable parallel algorithm for self-organizing maps with applications to sparse data mining problems. *Data Mining and Knowledge Discovery*, 3, pp. 171-195.
- Salton, G., 1989, *Automated Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer* (Reading MA: Addison-Wesley Publishing Company).
- Schifanella, R., Barrat, A., Cattuto, C., Markines, B. and Menczer, F., 2010, Folks in Folksonomies: Social Link Prediction from Shared Metadata. In *Proc. 3rd ACM International Conference on Web Search and Data Mining (WSDM '10)*, pp. 271-280 (New York, NY, USA: ACM).
- Skupin, A., 2004, The World of Geography: Visualizing a knowledge domain with cartographic means. *Proceedings of the National Academy of Sciences*, 101, pp. 5274-5278.
- Skupin, A., 2009, Discrete and continuous conceptualizations of science: Implications for knowledge domain visualization. *Journal of Informetrics*, 3, pp. 233-245.