

A Semantic Map of the last.fm Music Folksonomy

Joseph Biberstine, Russell J. Duhon, Elisha Allgood, Katy Börner

Cyberinfrastructure for Network Science Center

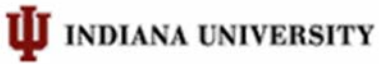
School of Library and Information Science

Indiana University

André Skupin

Department of Geography

San Diego State University

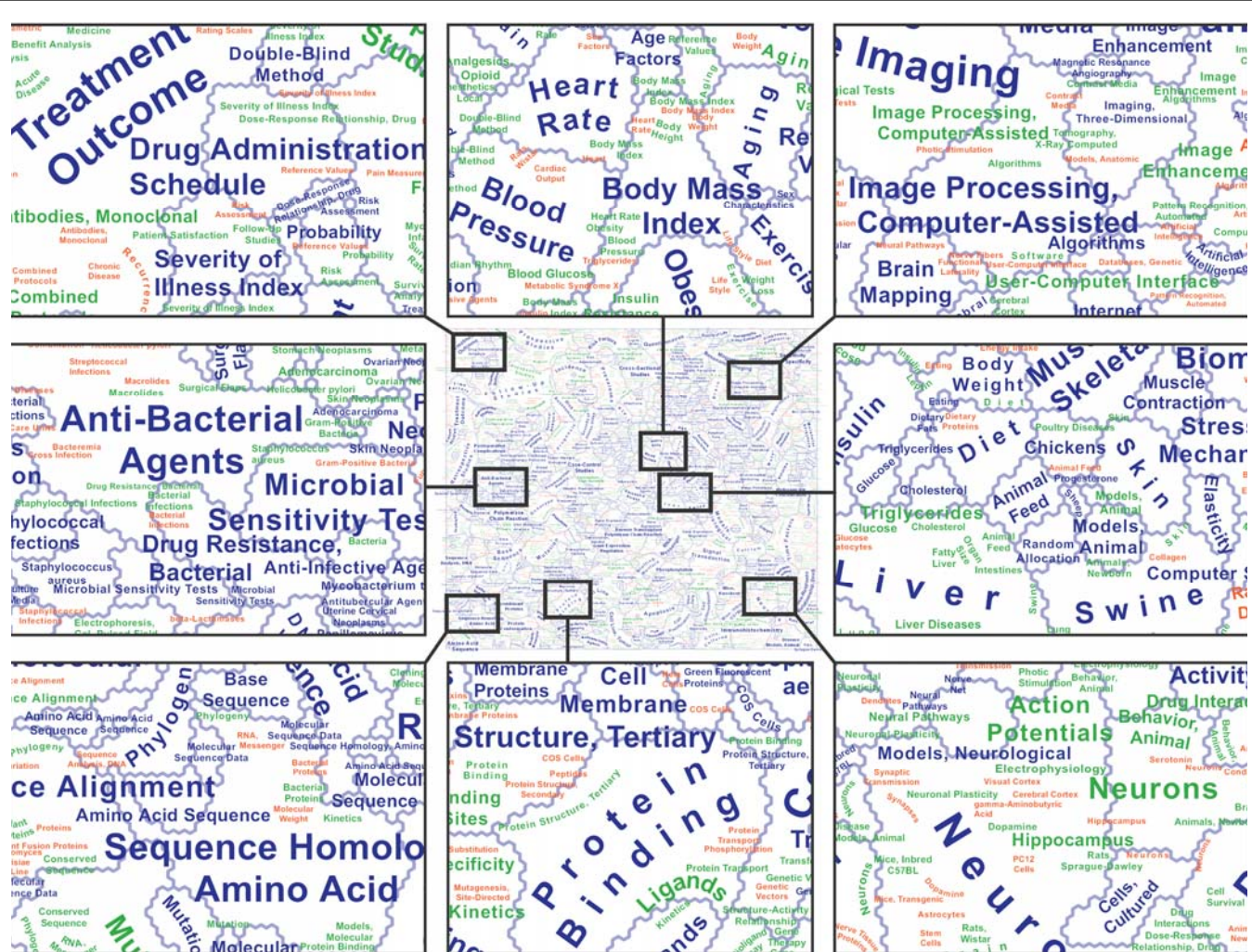


Finalist, 2010 NSF/Science Visualization Challenge

Finalist, 2011 Intl. Institute for Information Design (IIID) Awards

Background: Preceding Project

- “Accuracy of Models for Mapping the Medical Sciences”
 - NIH SBIR (HHSN268200900053C) & James S. McDonnell Foundation
 - Boyack, K.W., Newman, D., Duhon, R.D., Klavans, R., Patek, M., Biberstine, J.R., Schijvenaars, B., Skupin, A., Ma, N. and Börner, K. (2011) Clustering More Than Two Million Biomedical Publications: Comparing the Accuracies of Nine Text-Based Similarity Approaches. *PLoS ONE*. 6(3): e18029.
- Collaborators on SOM portion:
 - Katy Börner, Joseph Biberstine, Russell Duhon (Indiana Univ.)
- Data: 2.15 Mio documents (Medline, 2004-2008)
- Methods:
 - vector-space model (VSM) → self-organizing map (SOM) → GIS
- Challenge:
 - serial run time (estim.): 4 years (!)
- Solution:
 - parallelization & supercomputing: runtime reduced to 6 days



By Joseph Biberstine, Russell J. Duhon, Jennifer R. S. Coffey, Katy Börner, Cyberinfrastructure for Network Science Center, Indiana University, Bloomington and André Skupin, San Diego State University

This landscape is a 275 by 275 grid of hexagonal neurons. Regions on the landscape are labeled by the MeSH terms with which their constituent neurons associate most strongly. Light purple borders separate regions defined by each neuron's single strongest term association; those regions are marked with purple labels. Green labels depict regions defined by the second-strongest term association of each neuron, and so on, as shown in the legend below.



Acknowledgments

This research is funded by the Cyberinfrastructure for Network Science Center at the School of Library and Information Science, Indiana University; the National Science Foundation under award SBE-0738111; and NIH Awards R21DA024259 and IU4ARR029822-01. Indiana University's Big Red supercomputer used in this study is supported by the National Science Foundation under Grant No. ACI-0338618I, OCI-0451237, OCI-0535258, and OCI-0504075, a Shared University Research grants from IBM, Inc., and the Indiana METACyt Initiative supported in part by Lilly Endowment, Inc.

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.



Figure 1 consists of six conceptual maps (a-f) illustrating the structure of research domains in biology and medicine. The maps are arranged in a 3x2 grid. Map (a) shows broad domains: Epidemiology, Social Aspects, Tools and Diagnosis, Human Factors, Env. Factors, Genetics and Proteomics, and Cell Biology. Map (b) highlights 'Modelling & Research Methodology, Dig.' in red, with associated topics: Students/Biomed, Logistic models for Ed?, binary/multi-v outcomes, Attention which subsumes visual attention, Body Systems, Environmental Monitoring/Issues, and Heat?. Map (c) shows Research Methods, Research Methods, Research Methods, Visual Perception, Cancer/Oncology, Research Methods, Well-being, Environmental Medicine/Epidemiology Genetics, Nutrition, Plant Science, Cancer/Oncology, and Genetics. Map (d) shows Clinical Trials, Population Studies, Imaging, Informatics, Environmental Models, Computational Animal, Genomics, and Animal Models. Map (e) shows Clinical Research, Experimental Psychopathology Posture, Genetics Basic Research Non-clinical, Electrophysiology, Cellular, and Basic Animal Research. Map (f) shows Cancer, Genetics, Proteomics, Obesity, Environment Health, Neuro-Cognition, Neuro-Cognition, Neurophysiol/ Psychophysiol, Neurophysiol/ Psychophysiol, and Proteomics.

Raw Data - Source

- Last.fm is a social Internet radio site
 - users share information about songs they are listening to
 - they can also tag songs
 - with any strings of text they like



Need new music?

Last.fm lets you effortlessly **keep a record of what you listen to*** from any player. Based on your taste, Last.fm recommends you more music and concerts!

*We had to invent a word for this, it's called scrobbling.



Raw Data - Summary

- Gathered during first half of 2009
- 99,405 registered users
 - 52,452 active
- 281,818 tags
- 1,393,559 songs
- 10,936,545 annotations
 - an annotation is a (user, tag, song) triple, a tagging event
- data originally collected for:
 - Schifanella, R., Barrat, A., Cattuto, C., Markines, B., and Menczer, F. (2010). Folks in Folksonomies: Social Link Prediction from Shared Metadata. *Proc. 3rd ACM International Conference on Web Search and Data Mining (WSDM)*.



Top Tags

<i>rock</i>	<i>singer-songwriter</i>	<i>heavy metal</i>
<i>electronic</i>	<i>80s</i>	<i>chillout</i>
<i>seen live</i>	<i>folk</i>	<i>dance</i>
<i>indie</i>	<i>hard rock</i>	<i>british</i>
<i>alternative</i>	<i>progressive rock</i>	<i>90s</i>
<i>pop</i>	<i>indie rock</i>	<i>psychedelic</i>
<i>female vocalists</i>	<i>electronica</i>	<i>blues</i>
<i>jazz</i>	<i>punk</i>	<i>hip-hop</i>
<i>classic rock</i>	<i>instrumental</i>	<i>post-rock</i>
<i>experimental</i>	<i>soul</i>	<i>new wave</i>
<i>ambient</i>	<i>black metal</i>	<i>soundtrack</i>
<i>metal</i>	<i>industrial</i>	<i>classical</i>
<i>alternative rock</i>	<i>death metal</i>	<i>00s</i>

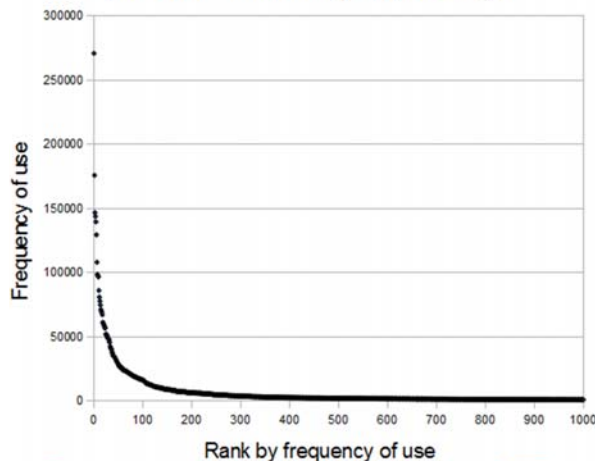
Tags Are More Than Just Genres

- Intensional
 - From recognized genres to simple objective facts
 - rock (rank 1)
 - electronic (2)
 - ..
 - female vocalists (7)
 - female vocalist (64)
 - acoustic (51)
- Extensional
 - A mix of social signals, properties of the user-song experience, and aides to personal categorization
 - seen live (3)
 - beautiful (48)
 - favorites (54)
 - albums i own (97)
 - altar of the metal gods (58)

Raw Data - Thresholding

- SOM method will not scale to 280,000+ tags/dims in raw form
- we consider only the 1,000 most frequently applied tags
- keep only songs annotated with any of these tags by some user
- characterize each song as a vector over each tag dimension summed across users

1,000 Most Frequently Applied Tags



	Raw	Preprocessed
Tags	281,818	1,000
Songs	1,393,559	1,088,761
Avg Tags/Song	7.8	6.8

Self-Organizing Map Algorithm - Parallelized Implementation

- Task completely intractable using typical SOM software
- Preceding project trained on twice as many data and twice as many dimensions, using our own implementation
 - Divide the training data among multiple processes
 - Each process holds a complete copy of the map
 - Periodically synchronize process-local copies of the map to create a new process-global map
- Adapted with several project-specific optimizations from:
 - Lawrence, R.D., Almasi, G.S., Rushmeier, H.E. (1999). A Scalable Parallel Algorithm for Self-Organizing Maps with Applications to Sparse Data Mining Problems. *Data Mining and Knowledge Discovery*.

Training the Map

- 2D hexagonal lattice of 32,400 neurons (180x180)
- Input space metric: cosine similarity
 - Induced interpretation: Each training vector (and so consequently each neuron vector) represents a direction in the 1,000-dimensional tag space

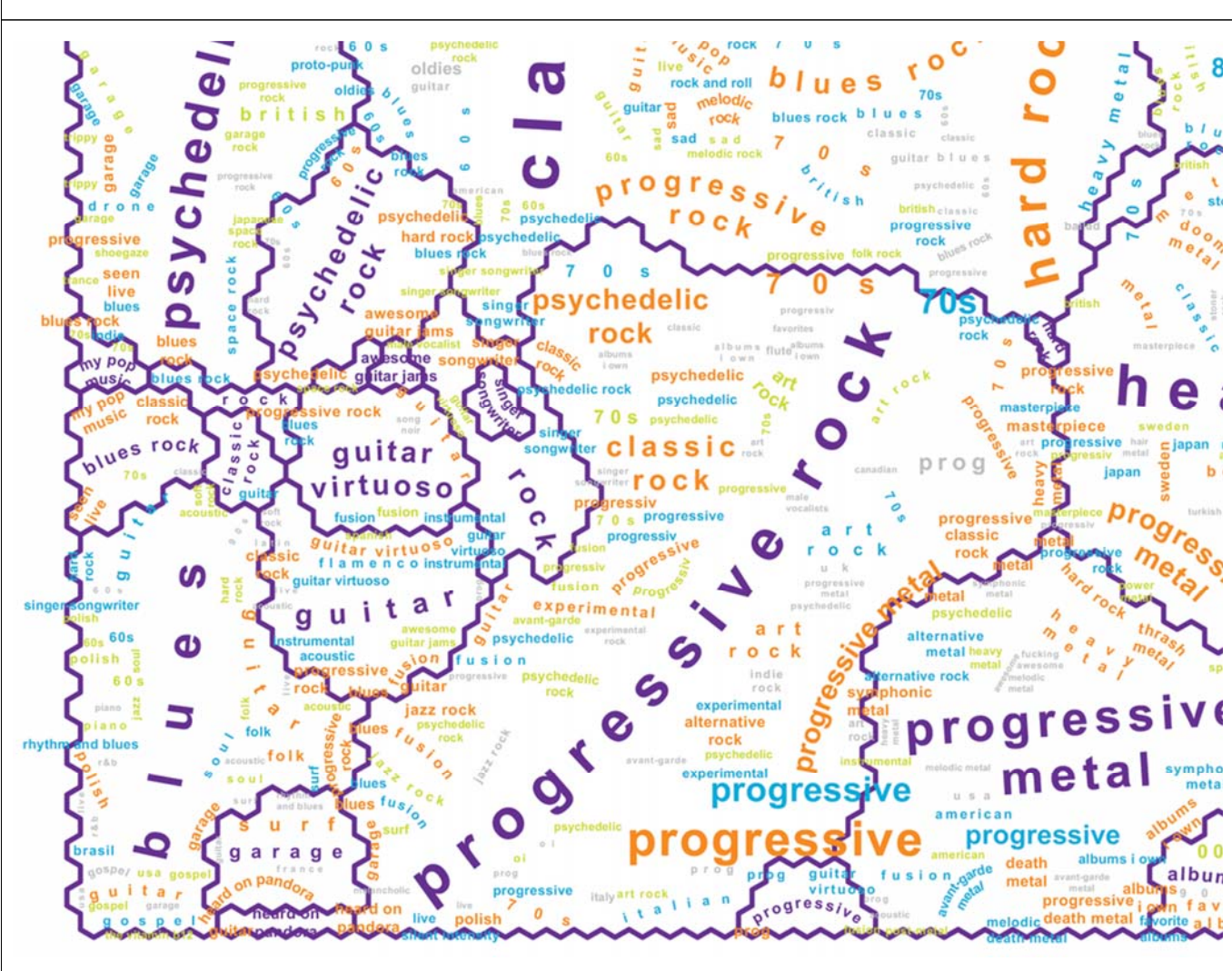
$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

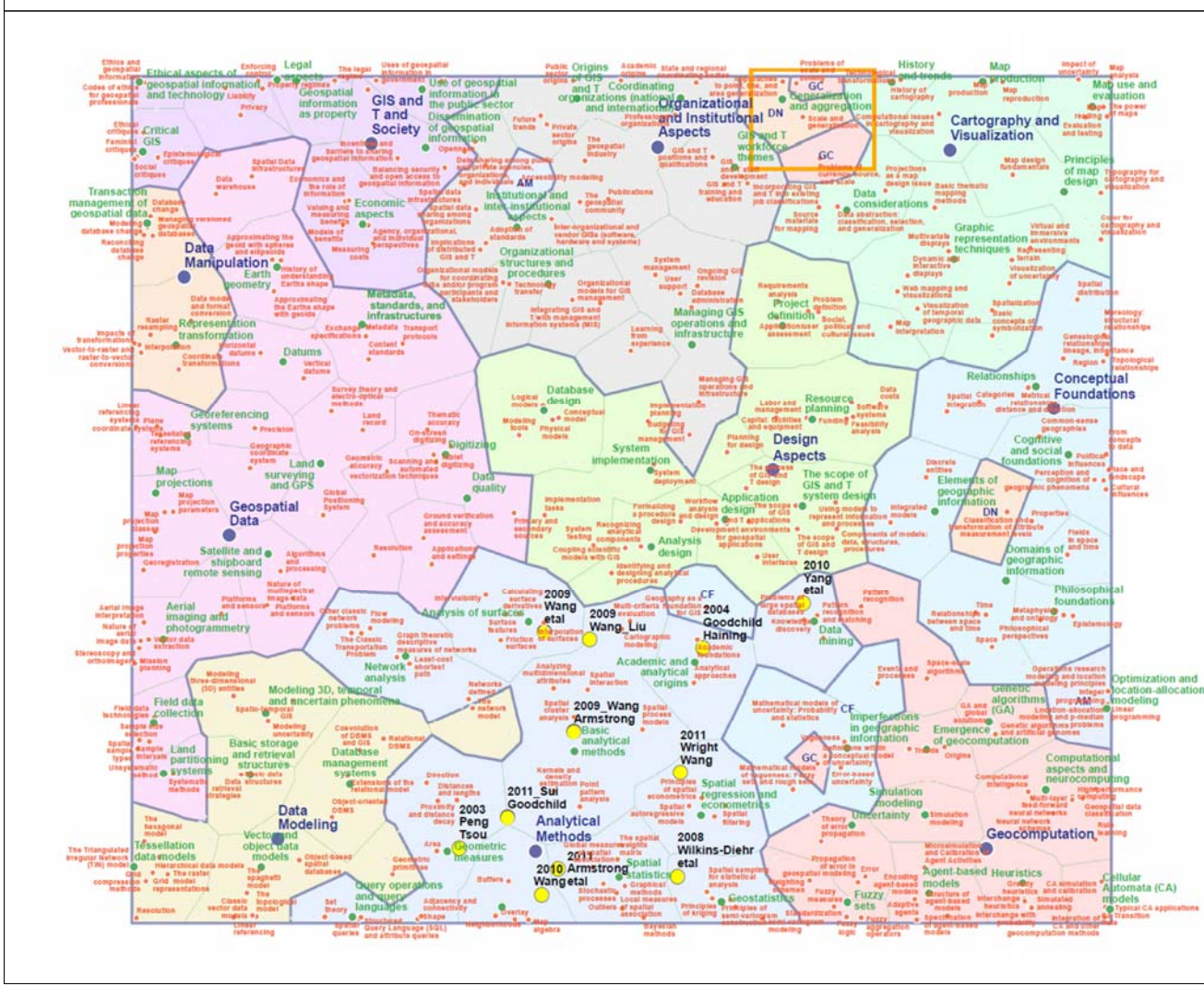
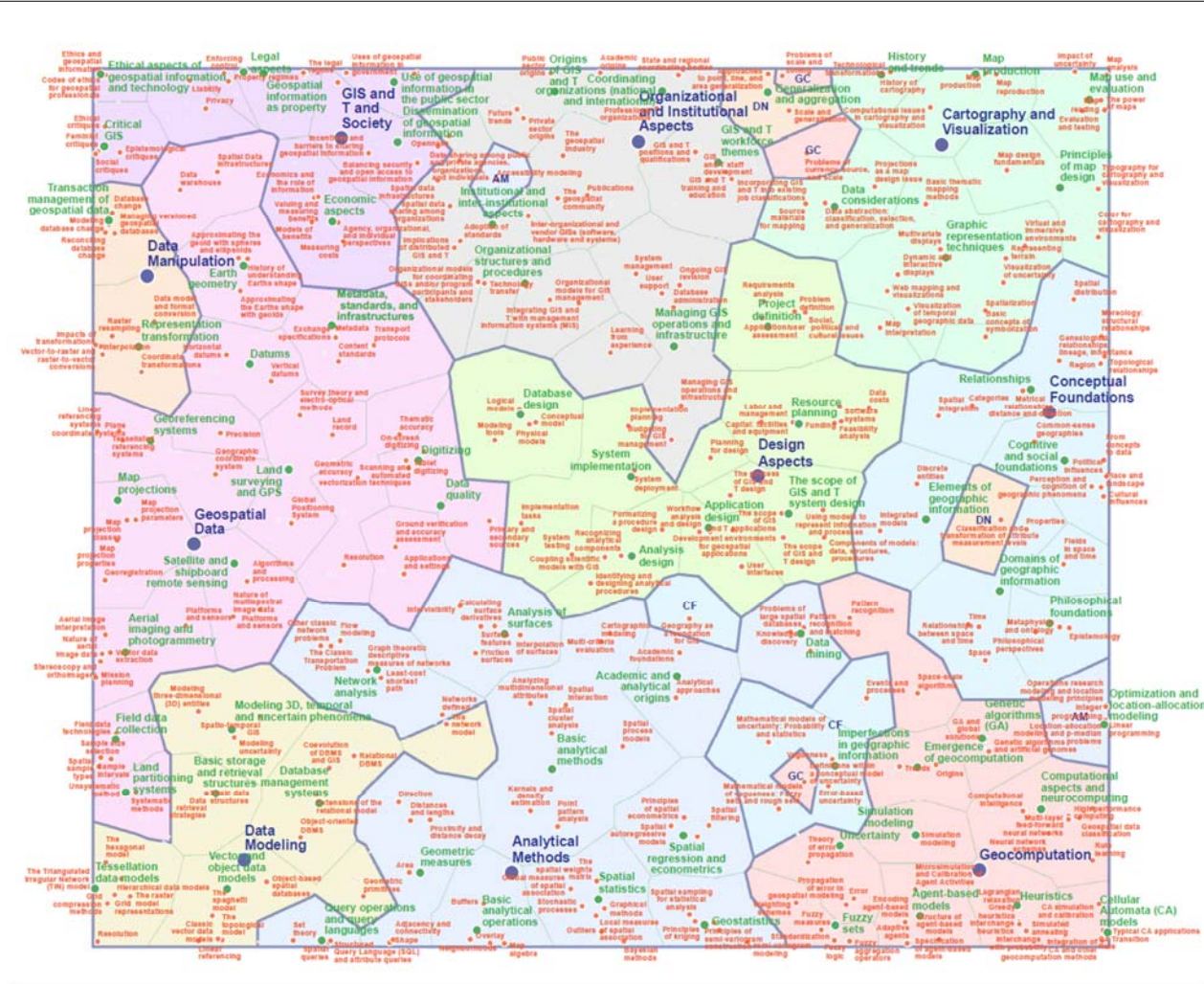
- 50 complete passes over the training data
- 300 processes across 100 compute nodes of Big Red, a supercomputer at Indiana University
 - Parallel runtime = 13 hours
 - Serial equivalent runtime = 5 months

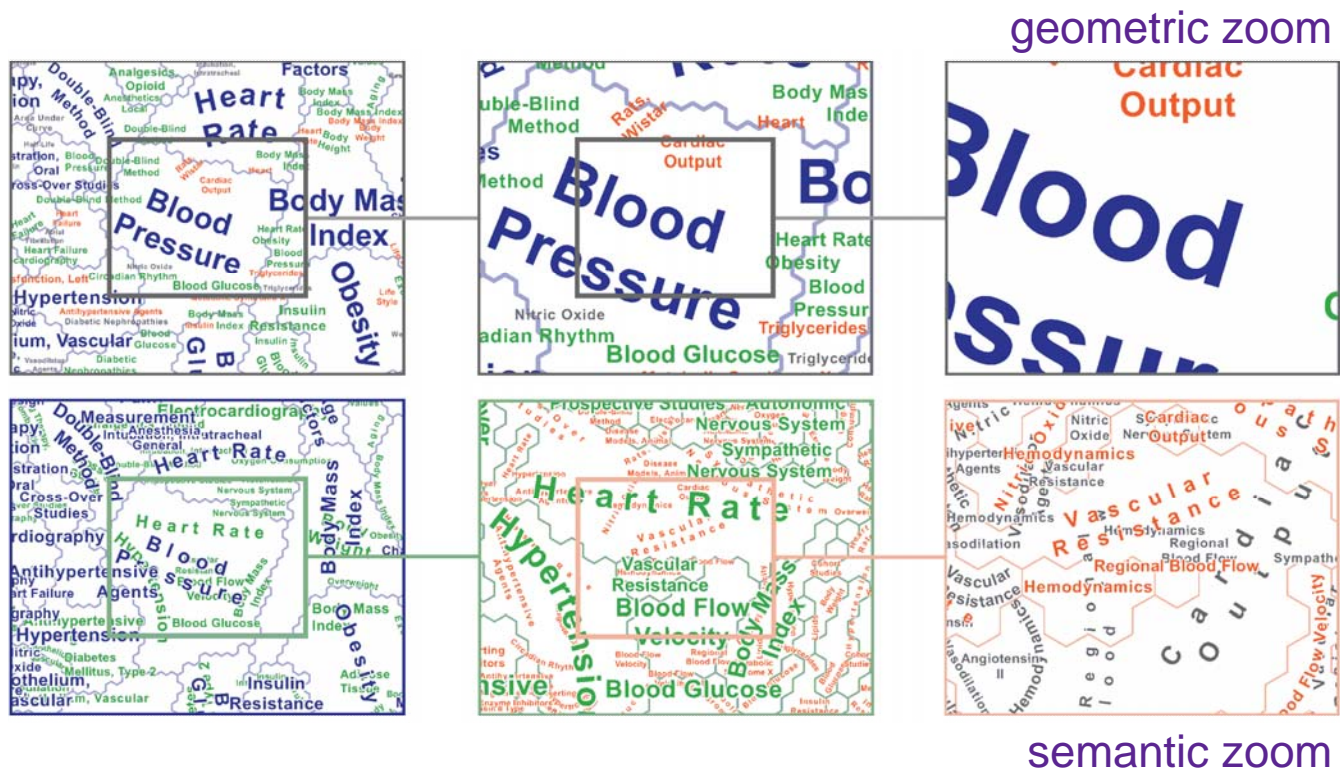


Visualization









Potential Applications

- Interactive music navigator and playlist generator
- Mapping portfolios as fields of neuronal activation
 - for the set of songs associated with any entity, we can see where in the world they belong
 - user: her/his favorite songs
 - band: their complete work
 - group of users: What is their turf?
 - ... or look at derivatives of these fields
 - What is the difference between *The Who* and *The Guess Who*?
 - How has this entity moved through the world of music over time?
 - Where have listeners like me headed next?

Contributors

- Joseph Biberstine, Indiana University
- André Skupin, San Diego State University
 - contact: skupin@mail.sdsu.edu
- Russell J. Duhon, IU
- Elisha Allgood, IU
- Katy Börner, IU

Funding

- School of Library and Information Science, Indiana University
- Cyberinfrastructure for Network Science Center, Indiana University
- National Science Foundation

