



LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

Decomposition of Large Scale Semantic Graphs via an Efficient Communities Algorithm

Y. Yao

February 11, 2008

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

Auspices Statement

This work performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. This work was funded by the Laboratory Directed Research and Development Program at LLNL under project tracking code 06-ERD-038.

Decomposition of Large Scale Semantic Graphs via an Efficient Communities Algorithm

Project Overview

Semantic graphs have become key components in analyzing complex systems such as the Internet, or biological and social networks. These types of graphs generally consist of sparsely connected clusters or “communities” whose nodes are more densely connected to each other than to other nodes in the graph. The identification of these communities is invaluable in facilitating the visualization, understanding, and analysis of large graphs by producing subgraphs of related data whose interrelationships can be readily characterized. Unfortunately, the ability of LLNL to effectively analyze the terabytes of multisource data at its disposal has remained elusive, since existing decomposition algorithms become computationally prohibitive for graphs of this size. We have addressed this limitation by developing more efficient algorithms for discerning community structure that can effectively process massive graphs.

Project Goals

Current algorithms for detecting community structure, such as the high quality algorithm developed by Girvan and Newman [1], are only capable of processing relatively small graphs. The cubic complexity of Girvan and Newman, for example, makes it impractical for graphs with more than approximately 10^4 nodes. Our goal for this project was to develop methodologies and corresponding algorithms capable of effectively processing graphs with up to 10^9 nodes. From a practical standpoint, we expect the developed scalable algorithms to help resolve a variety of operational issues associated with the productive use of semantic graphs at LLNL.

Relevance to LLNL Mission

In recent years, LLNL has developed semantic graph technologies capable of fusing disparate facts from diverse sources into massive semantic graphs to facilitate inference of complex and anomalous behaviors embedded within the data. A critical challenge in effectively applying this technology to the Lab’s mission space is to decompose massive graphs into meaningful subgraphs that an analyst can efficiently interrogate to identify these behaviors. This research represents a significant contribution to the counterterrorism, biodefense, and nonproliferation missions of LLNL, because efficient decomposition methodologies will provide the foundation for information analysis environments enabling large-scale data mining, information discovery and visualization.

FY07 Accomplishments and Results

During FY07, we completed a graph clustering implementation that leverages a dynamic graph transformation to more efficiently decompose large graphs. In essence, our approach dynamically transforms the graph (or subgraphs) into a tree structure consisting of biconnected components interconnected by bridge links. This isomorphism allows us to compute edge betweenness, the chief source of inefficiency in Girvan and Newman’s decomposition algorithm, much more efficiently, leading to significantly reduced computation time. Test runs on a desktop computer have shown reductions of up to 89% (see Table 1).

Our focus this year has been on the implementation of parallel graph clustering on one of LLNL’s supercomputers. In order to achieve efficiency in parallel computing, we have exploited the fact that large semantic graphs tend to be sparse, comprising loosely connected dense node clusters. When implemented on distributed memory computers, our approach performed well on several large graphs with up to one billion nodes, as shown in Table 2. The rightmost column of Table 2 contains the associated Newman’s modularity [1], a metric that is widely used to assess the quality of community structure. (See [4] for more details.)

Existing algorithms produce results that merely approximate the optimal solution, i.e., maximum modularity. We have developed a verification tool for decomposition algorithms, based upon a novel integer linear programming (ILP) approach, that computes an exact solution. We have used this ILP methodology to find the maximum modularity and corresponding optimal community structure for several well-studied graphs in the literature (e.g., Figure 1) [3].

The above approaches assume that modularity is the best measure of quality for community structure. In an effort to enhance this quality metric, we have also generalized Newman’s modularity based upon an insightful random walk interpretation that allows us to vary the scope of the metric. Generalized modularity has enabled us to develop new, more flexible versions of our algorithms.

In developing these methodologies, we have made several contributions to both graph theoretic algorithms and software engineering. We have written two research papers for refereed publication [3-4] and are working on another one [5]. In addition, we have presented our research findings at three academic and professional conferences.

Table 1: Computation time reduction over Girvan & Newman’s method

T_{gn} : Girvan & Newman’s time (min); T: our time (min); $R = 100(T - T_{gn}) / T_{gn}$.

Graph	Nodes	Links	T_{gn}	T	R(%)
Erdos972	5488	8972	1180.1	129.6	89.0
Hep-Th	7610	15751	1524.5	670.0	56.1
Kohonen	4470	12720	22.8	8.7	62.0
Power	4941	6594	723.0	550.4	23.9

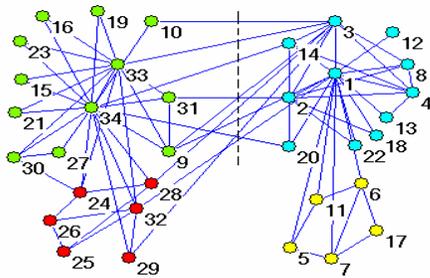


Figure 1: Optimal community structure for Zachary’s karate club [2]

Table 2: Computation time for parallel graph clustering

T: computation time (min); Q: modularity

Graph	Nodes	Links	CPUs	T	Q
G10m	10000000	43749984	4	28.9	0.40
G100m	100000000	298437392	32	706.0	0.72
G1000m	1000000000	2624753446	512	710.0	0.78

Related References

1. M. E. J. Newman, and M. Girvan, Finding and evaluating community structure in networks, *Phys. Rev. E* 69, 026113, 2004
2. W. Zachary, “An information flow model for conflict and fission in small groups,” *Journal of Anthropological Research*, 33, 452-473, 1977

3. Y. Yao, T. L. Hickling, W. G. Hanley, and J. S. Lenderman, Graph Clustering Evaluation via Integer Linear Programming, *Proceedings of the 2007 International Conference on Data Mining*, pp 369-375, 2007 (UCRL-PROC-228037)
4. Y. Yao and J. S. Lenderman, Efficient Betweenness Computation for Finding Community Structure in Graphs, *Computers and Operations Research*, Submitted (UCRL-JRNL-235157)
5. J. S. Lenderman and Y. Yao, A Generalization of Motivated by a Random Walk Interpretation, *Phys. Rev. E*, To submit