

1 Research Overview and Interests

Current work. Since summer 2006, I have worked as a research assistant for Alexander Gray’s FASTlab. I have explored fast algorithms for problems in computational geometry, including convex hulls, path finding, and minimum spanning trees. I have spent much of my time investigating scientific applications for these algorithms, especially in biology and chemistry. As a graduate student, I am investigating ways to apply my earlier techniques to designing faster and more accurate computational chemistry algorithms.

Past research. Previously, during Summer 2005, I participated in the NSF Research Experience for Undergraduates program at Georgia Tech. I studied online graph coloring algorithms with Prof. Trotter in the School of Mathematics. This experience gave me valuable experience with reading papers, organizing a complex research topic, and exploring new ideas. However, I found mathematics too far removed from the scientific applications I was interested in, which led me to my current work.

2 Euclidean Minimum Spanning Trees

Context. Euclidean Minimum Spanning Trees are fundamental structures in computational geometry and are widely applied in network design, optimization, and computer vision. Additionally, the EMST can be used to find a hierarchical clustering of the underlying points. This method of clustering is commonly applied throughout science and is particularly popular for clustering gene microarray data [2] and objects in deep-space surveys [1].

The EMST problem is well-studied, and there are effective algorithms for many cases, such as in two dimensions. However, none of the existing algorithms are truly scalable to massive data sets and arbitrary dimensions. New algorithms are necessary to handle large-scale data gathering efforts, such as the Sloan Digital Sky Survey and the Human Genome Project. Working with Alexander Gray and his FASTlab, I was able to develop and implement the fastest existing algorithm for EMST in arbitrary dimensions [4].

Teamwork. My first work with the FASTlab was on extending some of the lab’s previous work to design a fast algorithm for the EMST problem. Since I was new, I relied on discussions with the lab’s graduate students to become familiar with existing approaches and learn the code base. Throughout the project, I often discussed ideas and difficulties with Dr. Gray and the other students in the lab. These interactions often allowed me to find ways around difficulties.

Individual work. Although others helped as needed, I did most of the work individually. I searched the existing literature to determine the best existing algorithms. I determined the link between the previous work and the EMST problem and worked out the details of my own EMST algorithm. I wrote the code for my algorithm and the competitors, then organized and carried out the comparisons. I have submitted a conference paper on my algorithm which is currently under review.

Algorithm details. I applied an existing framework, known as *generalized N-body problems* [3], to finding EMST’s. This framework has been previously used to create fast algorithms for many problems in statistics and machine learning including the N -point correlation, all-nearest-neighbors, and kernel density estimation. Using these ideas, I was able to design and implement the fastest existing algorithm for EMST’s in any dimension.

My algorithm, DUALTREEBORUVKA, applies multi-tree recursion to Borůvka’s

algorithm for finding MST's. Borůvka's algorithm maintains a spanning forest and connects each component to its nearest neighbor in each step. I applied the space-partitioning trees and multi-tree recursion ideas to quickly compute the nearest neighbor of each component. In particular, the algorithm uses a *kd-tree*, which maintains a hyper-rectangle or bounding-box around each node.

Dual-tree Recursion. The simplest way to compute the nearest neighbor of a component Q is to find the nearest neighbor of all points $q \in Q$ and keep the closest. For each q , we store the distance to the nearest neighbor found so far, $d^u(q)$, as an upper bound on the true distance. Using the *kd-tree*, the algorithm can compare a point q with an entire node R . Using the bounding box of R , the algorithm can compute a lower bound $d^l(q, R) < d(q, r)$ for all $r \in R$. Then, if $d^l(q, R) > d^u(q)$ for some R , it can *prune* any further consideration of points in R . Otherwise, recursively consider the two children of R .

This idea can be extended even further. Since Borůvka's algorithm needs the nearest neighbor of each component Q , my algorithm exploits the fact that these points are also grouped in the *kd-tree*. It maintains an upper bound $d^u(Q)$ on $d^u(q)$ for all points $q \in Q$. It compares nodes Q and R , computes the minimum distance between their bounding boxes: $d^l(Q, R)$, and prunes if $d^l(Q, R) > d^u(Q)$. Otherwise, it recursively compares the children of Q with the children of R . By accounting for points in the same component, this method can quickly and efficiently solve the problem of finding each component's nearest neighbor.

Results This method proved to be the fastest known algorithm for finding EMST's in general metric spaces. I compared it to several well known MST algorithms, including GEOMST2 [5], the previous fastest algorithm. Since this research was motivated by scientific applications, especially cosmology and

computational biology, I also tested my algorithm on two large data sets from these areas. One is three-dimensional spectral data taken from the SDSS, and the other is a compressed (12-dimensional) representation of 16,000 steps of folding simulations for 20 proteins. These experiments clearly demonstrated the scalability of my algorithm in terms of running time and storage requirements.

Experience. Through my research involvements, particularly designing my EMST algorithm, I gained experience with all aspects of algorithm design. I searched and compared the existing methods and explored the scientific literature for applications. I gained experience with powerful algorithmic techniques and the challenges associated with

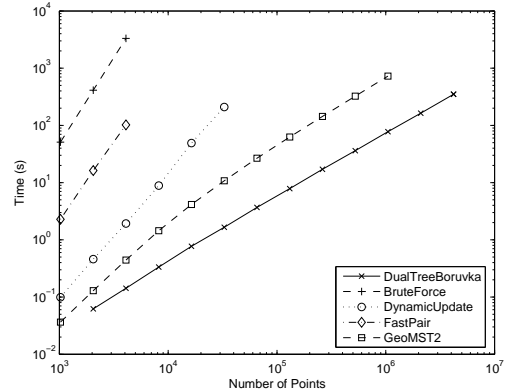


Figure 1: Running-time comparison on clustered synthetic data.

dim	N	GEOMST2	DTB	Speedup
3	389354	78.0	16.9	4.6
12	320000	702	62.3	11.3

Table 1: Running times on SDSS data and protein data.

applying them to new problems. I also have worked on communicating my research through a conference submission. These experiences, particularly with these algorithmic methods, are crucial to my current work with fast algorithms for computational chemistry problems.

References

- [1] J. D. Barrow, S. P. Bhavsar, and D. H. Sonoda. Minimal spanning trees, filaments and galaxy clustering. *MNRAS*, 216:17–35, Sept. 1985.
- [2] M. B. Eisen, P. T. Spellman, P. O. Brown, and D. Botstein. Cluster analysis and display of genome-wide expression patterns. *PNAS*, 95(25):14863–14868, 1998.
- [3] A. Gray and A. W. Moore. N-body problems in statistical learning. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 13*. MIT Press, 2001.
- [4] W. March and A. Gray. Large scale euclidean mst and hierarchical clustering. In *SIAM International Conference on Data Mining 2008*. Under review, 2008.
- [5] G. Narasimhan, J. Zhu, and M. Zachariasen. Experiments with Computing Geometric Minimum Spanning Trees. In *Proceedings of ALENEX'00*, Lecture Notes in Computer Science, pages 183–196. Springer-Verlag, January 2000.