

# Computational Science Graduate Fellowship Write Up

Koby Hayashi

January 7, 2019

## 1 Field of Interest

2250 character limit TOTAL for this section.

### 1.1 Problem Statement

My research focuses on scalable algorithms and software for computing constrained low-rank approximations for the mining, analysis and compression of data that may be modeled by tensors (including matrices) and hypergraphs. These arise in many areas of interest to the DoE such as robust clustering, anomaly detection, imaging analysis, and scientific simulation. For instance tensor decompositions have recently been explored as a means of compressing data generated from combustion simulations at Sandia. The increasing size and complexity of data sets, along with the machines used to analyze them, requires new methods for effectively utilizing and scalable algorithms for computing low-rank approximations.

I am pursuing this research from 3 perspectives: developing scalable implementations, exploring applications, and formulating methods for computing constrained low-rank approximations. Specifically, I will focus on Non-negative Matrix Factorization (NMF) and its variants, tensor factorizations, and joint factorizations (factorizations utilizing multiple data sources). Demonstrating the usefulness of these methods combined with developing efficient implementations will allow researchers in various domains to utilize them in their work.

As an example, there is little work on applying NMF to clustering hypergraph data. Hypergraphs are generalizations of graphs where edges may connect multiple vertices. Many data sets can naturally be treated as hypergraphs. NMF has demonstrated competitiveness with spectral clustering and other state of the art methods for analyzing standard graph data. Spectral hypergraph clustering has been widely explored but many open questions remain such as how to appropriately and efficiently represent the hypergraph. Just as NMF proved to be a viable alternative in the graph case I believe it can show strong results for hypergraphs. I will investigate ways to efficiently leverage the hypergraph structure for a variety of data sets. Success will be measured with empirical results and mathematical reasoning.

### 1.2 Potential Impact: HPC

**Discuss the potential impact of your research on your field. How would it advance high performance computing in general?** High performance computing will play a key role in my research as I seek to implement and deploy software that can run on large real world data sets. I have worked on the supercomputers Rhea, Eos, and Titan at ORNL. Tools I have experience

with include OpenMP and MPI for parallel programming. Furthermore, through course work and research experience I have gained experience in gathering performance benchmarks and evaluating the results.

HPC, as many fields are, currently suffers from a number of big data problems. For example, applications that produce too much or too complex data to store and/or analyze. Tools developed by this research can help alleviate such issues. Tensor decompositions have already been show to be a useful tool for the lossy compression of simulation data at Sandia. Now there is interest within the DoE in using low-rank approximations in the areas of inverse problems, linear unmixing, and reduced order modeling. Different data sets and problems have varying requirements of methods used to analyze them. For example, non-negative data often benefits from a non-negativity constraint if interpretability is desired in the returned factorization. However, knowing what constraints to impose, solving the derived optimization problem, and implementing or finding a code capable of handling a large problem size is often difficult. While there exist fast HPC implementations of approximations such as the SVD or Eigenvalue Decomposition this is rarely the case for other constrained low-rank approximations.

For example, during my MS degree my colleagues and I published results discovered by applying a tensor decomposition, the CP decomposition, to a 4-way Neuro-Imaging tensor. The data set we worked with was on the order of tens of gigabytes while Neuro-Imaging data sets can be hundred of gigabytes into the terabytes. Previously, there was no scalable software readily available that met our specific requirements. Our implementation is more general and up to 2x faster than existing implementations. In conclusion, developing software for computing constrained low-rank approximations that can handle large problem sizes and deal with diverse data sets will benefit the HPC and general scientific community.



Name: Koby Hayashi

Year 0 of 4

## Program of Study

Listed are the courses in science and engineering, applied mathematics, and computer science that you agreed to take on your most recent proposed Program of Study. Please sign this page and return it to the address below.

### University: Georgia Institute of Technology

Course number	Course Title	Credit hours	Term and Year	Grade	Academic Level	Fulfill POS
<b>Science/Engineering</b>						
PSY6042	Neuroimaging	3S	Fall 2019		G	Fulfill
PSY6090	Cognitive Neuroscience	3S	Spring 2020		G	Fulfill
<b>Mathematics and Statistics</b>						
CSE6643	Numerical Linear Algebra	3S	Spring 2019		G	Fulfill
CSE6644	Iterative Methods for Systems of Equations	3S	Fall 2019		G	Fulfill
<b>Computer Science</b>						
CSE6140	Computational Science and Engineering Algorithms	3S	Fall 2018	A	G	Fulfill
CSE6220	High Performance Computing	3S	Spring 2019		G	Fulfill
CSE6230	High Performance Parallel Computing	3S	Fall 2018	A	G	Fulfill