

We explore the trade-offs of performing linear algebra in Apache Spark versus the traditional C and MPI approach by examining three widely-used matrix factorizations: NMF (for physical plausibility), PCA (for its ubiquity) , and CX (for model interpretability). We apply these methods to TB-scale problems in particle physics, climate modeling, and bioimaging using algorithms that map nicely onto Spark’s data-parallelism model. We perform scaling experiments on up to 1600 Cray XC40 nodes, describe the sources of slowdowns, and provide tuning guidance to obtain high performance.