



# Parallel accelerated vector similarity calculations for genomics applications<sup>☆</sup>

Wayne Joubert<sup>a,\*</sup>, James Nance<sup>a</sup>, Deborah Weighill<sup>a,b</sup>, Daniel Jacobson<sup>a,b</sup>

<sup>a</sup> Oak Ridge National Laboratory, 1 Bethel Valley Road, Oak Ridge, TN 37831, United States

<sup>b</sup> The Bredeben Center for Interdisciplinary Research and Graduate Education, University of Tennessee, Knoxville, 444 Greve Hall, 821 Volunteer Blvd, Knoxville, TN 37996-3394, United States

## ARTICLE INFO

### Article history:

Received 22 May 2017  
Revised 3 March 2018  
Accepted 25 March 2018  
Available online 27 March 2018

### MSC:

65Y05 [Computer aspects of numerical algorithms: Parallel computation]  
68W10 [Algorithms: Parallel algorithms]

### Keywords:

High performance computing  
Parallel algorithms  
NVIDIA<sup>®</sup> GPU  
Intel<sup>®</sup> Xeon Phi  
Comparative genomics  
Vector similarity metrics  
Proportional Similarity metric

## ABSTRACT

The surge in availability of genomic data holds promise for enabling determination of genetic causes of observed individual traits, with applications to problems such as discovery of the genetic roots of phenotypes, be they molecular phenotypes such as gene expression or metabolite concentrations, or complex phenotypes such as diseases. However, the growing sizes of these datasets and the quadratic, cubic or higher scaling characteristics of the relevant algorithms pose a serious computational challenge necessitating use of leadership scale computing. In this paper we describe a new approach to performing vector similarity metrics calculations, suitable for parallel systems equipped with graphics processing units (GPUs) or Intel Xeon Phi processors. Our primary focus is the Proportional Similarity metric applied to Genome Wide Association Studies (GWAS) and Phenome Wide Association Studies (PheWAS). We describe the implementation of the algorithms on accelerated processors, methods used for eliminating redundant calculations due to symmetries, and techniques for efficient mapping of the calculations to many-node parallel systems. Results are presented demonstrating high per-node performance and parallel scalability with rates of more than five quadrillion ( $5 \times 10^{15}$ ) elementwise comparisons achieved per second on the ORNL Titan system. In a companion paper we describe corresponding techniques applied to calculations of the Custom Correlation Coefficient for comparative genomics applications.

© 2018 Elsevier B.V. All rights reserved.

## 1. Introduction

The measurement of the similarity of pairs of vectors is a computation required in many science domains including chemistry, image processing, linguistics, ecology, document processing and genomics. To satisfy domain-specific requirements, many different similarity measures have been developed [1,2].

<sup>☆</sup> This manuscript has been authored by UT-Battelle, LLC under Contract No. DE-AC05-00OR22725 with the U.S. Department of Energy. The United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United States Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes. The Department of Energy will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

\* Corresponding author.

E-mail address: [joubert@ornl.gov](mailto:joubert@ornl.gov) (W. Joubert).

The focus of the present study is the use of similarity measures in the analysis of GWAS and PheWAS results. GWAS analyses involve the statistical association of genetic variants with measured phenotypes [3]. These can be complex phenotypes such as disease states, or molecular phenotypes, such as the concentration of a particular metabolite or the expression level of a particular gene. While GWAS analyses are generally considered to involve testing association of variants with a single or limited number of phenotypes, PheWAS analyses involve testing the association of variants with a large number of different phenotypes [4]. The results of GWAS and PheWAS studies can be represented as a matrix of significant associations between variants and phenotypes, and profile vectors variants and phenotypes can be extracted from the rows and columns of this matrix. Pairwise comparisons of these vectors can allow for the discovery of phenotypes affected by similar genetic elements, or of groups of variants which affect similar phenotypes. (For example, see [5]). These studies, however, are computationally expensive, insofar as the computational work required for pairwise comparison grows as the square of the number of vectors. Even more challenging is the execution of higher-order studies which consider three or more vectors at a time—a technique required in order to discover relationships not discoverable by means of 2-way methods alone [6]—for which the computational complexity is even higher. In the past, such studies could be performed efficiently on workstations or small compute clusters. However, because of the large quantities of data involved, it is now necessary to employ large-scale high performance computing to execute scientific campaigns at the largest scales.

This paper describes advances in the development of algorithms and software to address this need. We present vector similarity measure calculation techniques for large datasets run on one of the world's largest compute systems, scaled to thousands of compute nodes equipped with GPU accelerators. The primary contributions of this paper are implementations of similarity calculation methods which: 1) achieve high absolute performance on GPUs as a result of careful mapping of calculations to the memory hierarchy and exploiting of the highly computationally intense BLAS-3-like structure of the targeted algorithms; 2) use asynchronous internode communication, data transfers and computations to ameliorate the costs of data motion; 3) strategically arrange the computations to avoid the potential 2X-6X performance loss factor arising from the redundant calculations due to symmetry; 4) carefully parallelize the algorithms to enable near-perfect scalability to thousands of compute nodes on leadership-class systems.

In this paper we focus on the 2-way and 3-way variants of the Proportional Similarity metric, also known as the Czekanowski metric [6,7], using an approach that is generalizable to other metrics. In the companion paper [8] we describe corresponding work on Custom Correlation Coefficient (CCC) [9] calculations with applications to comparative genomics.

Improving computational throughput for performing comparisons between pairs, triples or larger subsets of a set of vectors has been the focus of significant recent work centering around the use of parallelism, accelerated GPU or Intel Xeon Phi processing, or both. A broad overview of epistasis detection in comparative genomics including computational issues pertaining to parallelism and GPU acceleration is given in [10]. The GBOOST code, discussed in [11], is a gene-gene interaction code for 2-way studies optimized for single GPUs using encoding of gene data into bit strings with avoidance of redundant computations. Wang et al. [12] describes GWISFI, a single-GPU code for 2-way GWAS calculations. Gonzalez-Dominguez et al. [13] develops a UPC++ code for gene-gene interaction studies for small numbers of GPUs and Intel Phi processors exploiting vector hardware and hardware population count instructions. Gonzalez-Dominguez and Schmidt [14] considers 3-way interactions on a node with 4 GPUs. Solomonik et al. [15] develops parallel tensor computation methods, structurally similar to 3-way metrics computations, with particular attention to avoiding redundant computations; however, the work does not consider GPUs or shaping of the computational regions to accommodate processors with long vector lengths. Haque et al. [16] discusses similarity metric calculations for chemical informatics applications on single GPUs using space filling curve methods and hardware population count instructions; it recognizes the correspondence of these calculations to BLAS-3 matrix-matrix product computations and pays close attention to optimizing memory accesses. Wang et al. [17] considers 2-way studies on compute clouds using MapReduce on conventional CPUs. Yang et al. [18] adapts existing packages to perform 2-way CPU and GPU studies and 3-way CPU studies on as many as 200 cores in parallel. Goudey et al. [19] performs  $k$ -way GWAS studies for arbitrary  $k$  with consideration of load balancing and elimination of redundancies on a 4096-node IBM Blue Gene/Q system; results for a single GPU are also presented. Luecke et al. [20] performs 2-way analyses on up to 126 nodes of the Intel Phi-based Stampede system (cf. [21]). Koesterke et al. [22] considers 2-way computations on thousands of compute cores with good scalability and good absolute performance on conventional CPUs. Finally, recent work in [23] considers  $k$ -selection similarity search methods with applications to image data with results for small numbers of GPUs; that work however focuses primarily on the  $k$ -selection problem for nonexhaustive inexact similarity search, a different problem from what is considered here.

The present study is to our knowledge the first work bringing together all the required ingredients for high performance 2-way and 3-way comparative genomics studies on modern leadership-class systems: use of accelerated processors at high absolute performance; optimization of calculations for use with complex memory hierarchies; elimination of redundant computations; algorithm and code design to minimize costs of I/O; and careful arrangement of communications for near-ideal scaling to many thousands of compute nodes.

The remainder of this paper is structured as follows. After describing the 2-way and 3-way Proportional Similarity metrics in Section 2, we describe the techniques used to map these methods to GPUs and other manycore accelerated processors in Section 3. Then we describe the parallelization techniques applied to these methods in Section 4, followed by implementation details in Section 5. Computational results on the 27 petaflop Oak Ridge National Laboratory (ORNL) Cray XK7 Titan system are presented in Section 6, and conclusions are given in Section 7.

Download English Version:

<https://daneshyari.com/en/article/6935017>

Download Persian Version:

<https://daneshyari.com/article/6935017>

[Daneshyari.com](https://daneshyari.com)