

Large-scale linear regression: Development of high-performance routines

Alvaro Frank, Diego Fabregat-Traver and Paolo Bientinesi

Aachen Institute
for Advanced Study in
Computational Engineering Science

Financial support from the
Deutsche Forschungsgemeinschaft (German Research Foundation)
through grant GSC 111 is gratefully acknowledged.

Large-scale linear regression: Development of high-performance routines

Alvaro Frank, Diego Fabregat-Traver, and Paolo Bientinesi

AICES, RWTH Aachen, 52062 Aachen, Germany

alvaro.frank@rwth-aachen.de,
fabregat@aices.rwth-aachen.de,
pauldj@aices.rwth-aachen.de

Abstract. In statistics, series of ordinary least squares problems (OLS) are used to study the linear correlation among sets of variables of interest; in many studies, the number of such variables is at least in the millions, and the corresponding datasets occupy terabytes of disk space. As the availability of large-scale datasets increases regularly, so does the challenge in dealing with them. Indeed, traditional solvers—which rely on the use of “black-box” routines optimized for one single OLS—are highly inefficient and fail to provide a viable solution for big-data analyses. As a case study, in this paper we consider a linear regression consisting of two-dimensional grids of related OLS problems that arise in the context of genome-wide association analyses, and give a careful walkthrough for the development of OLS-GRID, a high-performance routine for shared-memory architectures; analogous steps are relevant for tailoring OLS solvers to other applications. In particular, we first illustrate the design of efficient algorithms that exploit the structure of the OLS problems and eliminate redundant computations; then, we show how to effectively deal with datasets that do not fit in main memory; finally, we discuss how to cast the computation in terms of efficient kernels and how to achieve scalability. Importantly, each design decision along the way is justified by simple performance models. OLS-GRID enables the solution of 10^{11} correlated OLS problems operating on terabytes of data in a matter of hours.

Keywords: Linear regression, ordinary least squares, grids of problems, genome wide association analysis, algorithm design, out-of-core, parallelism, scalability

1 Introduction

Linear regression is an extremely common statistical tool for modeling the relationship between two sets of data. Specifically, given a set of “independent variables” x_1, x_2, \dots, x_p , and a “dependent variable” y , one seeks the correlation terms $\beta_i, i = 1, \dots, p$ in the linear model

$$y = \beta_1 x_1 + \dots + \beta_p x_p. \tag{1}$$

In matrix form, Eq. (1) is expressed as $\mathbf{y} = X\beta + \epsilon$, where $\mathbf{y} \in R^n$ is a vector of n “observations”, the columns of $X \in R^{n \times p}$ are “predictors” or “covariates”, the vector $\beta = [\beta_1, \dots, \beta_p]^T$ contains the “regression coefficients”, and ϵ is an error term that one wishes to minimize. In many disciplines, linear regression is used to quantify the relationship between one or more \mathbf{y} ’s from the set \mathcal{Y} , and each of many x ’s from the set \mathcal{X} . The computational challenges raise from the all-to-all nature of the problem (estimate how strongly each of the covariates is related to each of the observations), and from the sheer size of the datasets \mathcal{Y} and \mathcal{X} , which often cannot be stored directly in main memory.

One of the standard approaches to fit the model (1) to given \mathbf{y} and X is by solving an ordinary least squares (OLS) problem; in linear algebra terms, this corresponds to computing the vector β such that

$$\beta = (X^T X)^{-1} X^T \mathbf{y}.$$

In typical datasets, \hat{m} , the number of available covariates ($\hat{m} = |\mathcal{X}|$), is much larger than p , the number of variables actually used in the model. In this case, a group of $l < p$ covariates is kept fixed, and the remaining $p - l$ slots are filled from \mathcal{X} , in a rotating fashion; it is not uncommon that the value $p - l$ is very small, often just one, thus originating $m \geq \hat{m}$ distinct OLS problems. Mathematically, this means computing a series of β_i 's such that

$$\beta_i = (X_i^T X_i)^{-1} X_i^T \mathbf{y}, \quad \text{where } i = 1, \dots, m; \quad (2)$$

here X_i consists of two parts: X_L , which contains l columns and is fixed across all m OLS problems, and X_{R_i} , which instead contains $p - l$ columns taken from \mathcal{X} . In many applications, m can be of the order of millions or even more.

When $t > 1$ dependent variables ($t = |\mathcal{Y}|$) are to be studied against \mathcal{X} , the problem (2) assumes the more general form

$$\beta_{ij} = (X_i^T X_i)^{-1} X_i^T \mathbf{y}_j, \quad (3)$$

where $i = 1, \dots, m$, and $j = 1, \dots, t$, indicating that one has to compute a two-dimensional grid of β_{ij} 's, each one corresponding to an OLS problem. This is for instance the case in genomics (multi-trait genome-wide association analyses) [1] and econometrics (explanatory variable exploration) [2].

Despite the fact that OLS solvers are provided by many libraries and languages (e.g., LAPACK, NAG, MKL, Matlab, R), no matter how optimized those are, any approach that aims at computing the 2D grid (3) via $t \times m$ invocations of a “black-box” routine is entirely unfeasible. The main limitations come from the fact that this approach leads to the execution of inefficient and redundant operations, lacks a mechanism to effectively manage data transfers from and to hard disk, and underutilizes the resources on parallel architectures.

In this paper, we consider an instance of Eq. (3) as it arises in genomics, and develop OLS-GRID, a parallel solver tailored for this application. Specifically, we focus on the study of *omics* data¹ in the context of genome wide association analyses (GWAA).² Omics GWAA study the relation between m groups of genetic markers and t phenotypic traits in populations of n individuals. In terms of OLS, each trait is represented by a vector y_j containing the trait measurements (one per individual); each matrix $X_i = [X_L | X_{R_i}]$ is composed of a set of l fixed covariates such as sex, age, and height (X_L), and one of the groups of $r = p - l$ markers (X_{R_i}). A positive correlation between markers X_{R_i} and trait y_j indicates that the markers may have an impact in the expression of the trait.

Typical problem sizes in omics GWAA are roughly $10^3 \leq n \leq 10^5$, $2 \leq p \leq 20$ (with $r = 1$ or 2), $10^6 \leq m \leq 10^8$, and $10^2 \leq t \leq 10^5$. An exemplary analysis with size $n = 30,000$, $p = 10$ ($l = 8$, $r = 2$), $m = 10^7$, and $t = 10^4$, poses three considerable challenges. First, it requires the computation of 10^{11} OLS problems, which, if tackled by a traditional “black-box” solver, would perform $O(10^{18})$ floating point operations (flops). Despite the fact that the problem lends itself to a lower computational cost and efficient solutions, a black-box solver ignores the structure of the problem and requires large clusters to obtain a solution. The second challenge is posed by the size of the datasets to be processed: assuming single-precision data (4 bytes per element), a GWAA solver reads as input about 2.4 TBs of data and produces as output 4 TBs of data. If the data movement is not analyzed properly, the time spent in I/O transfers might render the computation unfeasible.

¹ With the term *omics* we refer to large-scale analyses involving at least hundreds of traits [3,4,5].

² GWAA are often also referred to as genome wide association studies (GWAS) and whole genome association studies (WGAS).

Finally, the computation needs to be parallelized and organized so that the potential of the current multi-core and many-core architectures is fully exploited.

Contributions. This paper is concerned with the design and the implementation of OLS-GRID, a high-performance algorithm for large-scale linear regression. While we use omics GWAA as a case study, the discussion is relevant to a range of OLS-based applications. Specifically, we 1) illustrate how to take advantage of the specific structure in the grid of OLS problems to design specialized algorithms, 2) analyze the data transfers from and to disk to effectively deal with large datasets that do not fit in main memory, and 3) discuss how to cast the computation in terms of efficient kernels and how to attain scalability on many-core architectures. Moreover, by making use of simple performance models, we identify the performance bottlenecks with respect to the problem size. OLS-GRID, available as part of the GenABEL suite [6], allows one to execute an analysis of the aforementioned size in less than 7 hours on a 40-core node.

Related work. Genome-wide association analyses received a lot of attention in the last decade [7]. Numerous high-impact findings have been reported, including but not limited to the identification of genetic variations associated to a common form of blindness, type 2 diabetes, and Parkinson’s disease [8,9,10,11]. A popular approach to GWAA is the so called Variance Components model, which boils down to a set of equations similar to Eq. (3). The main difference with the present work lies on the core equation, where one has to solve grids of generalized least squares (GLS) problems instead of grids of OLSs. A number of libraries have been developed to carry out GLS-based GWAA, the most relevant being FaST-LMM, GEMMA, GWFGLS, and OmicABEL [6,12,13,5].

OmicABEL, developed within our research group, showed remarkable performance improvements with respect to the other existing methods [14,15]. Clearly, the same library can be used, by setting the covariance matrix to the identity, to solve Eq. (3). While possible, this is not advisable: OmicABEL reduces the two-dimensional grid of GLS problems to the grid of OLS problems

$$\tilde{b}_{ij} = (\tilde{X}_{ij}^T \tilde{X}_{ij})^{-1} \tilde{X}_{ij}^T \tilde{y}_j,$$

which is deceptively similar to Eq. (3); however, the subtle differences in the dependencies (subindices) of the design matrix X lead to a more expensive and less efficient algorithm. In Sec. 5 we show how the new OLS-GRID outperforms OmicABEL-Eig considerably.

A number of tools are focused on OLS-based linear regression analyses for GWAA; among them, we mention ProbABEL (also part of the GenABEL suite), GWASP, and BOSS [16,17,18]. GWASP stands out for its elegant algorithmic approach and a performance-oriented design; a more detailed discussion is given in Sec. 5.4.

Organization of the paper. The remainder of the paper is organized as follows. In Sec. 2 we describe the design of an efficient algorithm that exploits problem-specific knowledge. In Sec. 3 we analyze the data transfers and discuss possible limitations inherent to the problem at hand, while in Sec. 4 we focus on the high performance and scalability of our software. Section 5 presents performance results. Finally, Sec. 6 draws conclusions and discusses future work.

2 From the problem specification to an efficient algorithm

In this section we discuss the steps leading to the core algorithm behind OLS-GRID. Starting from an elementary and generic algorithm, we incrementally refine it into an efficient algorithm tailored specifically for grids of OLS arising in GWAA studies. The focus is on the reduction of the computational complexity.

As a starting point to solve Eq. (3), we consider a most naive black-box solver consisting of two nested loops (for each i and j) around a QR-based algorithm for OLS [19]:

$$\begin{aligned}\{Q, R\} &:= qr(X) \\ b &:= R^{-1}Q^T y.\end{aligned}$$

Since the design matrix X only depends on the index i , the loop ordering i, j makes it possible to compute the QR factorization once and reuse it across the j loop. Even with this simple code motion optimization, this first algorithm performs substantial redundant calculations due to the structure of X_i .

Recall that each matrix X_i can be logically seen as consisting of two parts $(X_L|X_{R_i})$, where X_L is constant, and X_{R_i} varies across different instances of X_i . Aware that the QR factorization of X_i has to be computed for each i , the question is whether or not such structure is exploitable. The answer lies in the ‘‘Partitioned Matrix Expression’’ (PME) of the QR factorization: for a given operation, the PME reveals how (portions of) the output operands can be expressed in terms of (portions of) the input operands [20].

In this specific case, we wonder if and how the vertical partitioning of X propagates to the computation of Q and R . Consider

$$\left(Q_L \mid Q_R \right) \left(\begin{array}{c|c} R_{TL} & R_{TR} \\ \hline 0 & R_{BR} \end{array} \right) = \left(X_L \mid X_R \right), \quad (4)$$

where $Q_L \in R^{n \times l}$, $Q_R \in R^{n \times r}$, $R_{TL} \in R^{l \times l}$, $R_{TR} \in R^{l \times r}$, and $R_{BR} \in R^{r \times r}$. By rewriting Eq. (4) as

$$\left(Q_L R_{TL} = X_L \mid Q_L R_{TR} + Q_R R_{BR} = X_R \right),$$

and using the orthogonality of Q ($Q_L^T Q_R = 0$), one derives the assignments

$$\begin{aligned}\{Q_L, R_{TL}\} &:= qr(X_L) \\ R_{TR} &:= Q_L^T X_R \\ \{Q_R, R_{BR}\} &:= qr(X_R - Q_L R_{TR}),\end{aligned} \quad (5)$$

which indicate that the factorization of X_L can be computed only once and re-used as X_R varies.

In Alg. 1 all the observations made so far are incorporated. In particular, the loop ordering is set to i, j , and code motion is applied whenever possible, to avoid redundant calculations. Each line of the algorithm is annotated with the corresponding BLAS or LAPACK kernel and its computational cost. Realizing that not only X_L , but also Q_L and R_{TL} are constant (they only depend on X_L), we can now deliver a more sophisticated algorithm which saves even more flops.

As a direct effect of the partitioning of Q_i in $(Q_L|Q_{R_i})$ in line 7, the vector z_{ij} can be decomposed as

$$\begin{pmatrix} z_{T_j} \\ z_{B_{ij}} \end{pmatrix} := \begin{pmatrix} Q_L^T y_j \\ Q_{R_i}^T y_j \end{pmatrix},$$

suggesting that the top part ($z_{T_j} := Q_L^T y_j$) may be precomputed, once per vector y_j , and then reused.

Similarly, the structure of R_i can be exploited to avoid redundant computation within the triangular system in line 8. By using the same top-bottom splitting for z_{ij} , and partitioning b_j accordingly, we obtain the expression

$$\left(\begin{array}{c|c} R_{TL} & R_{TR_i} \\ \hline 0 & R_{BR_i} \end{array} \right) \begin{pmatrix} b_T \\ b_B \end{pmatrix} = \begin{pmatrix} z_{T_j} \\ z_{B_{ij}} \end{pmatrix},$$

Algorithm 1 : Structure of X exposed and exploited

```
1:  $\{Q_L, R_{TL}\} := qr(X_L)$  (QR)  $2nl^2$ 
2: for  $i := 1$  to  $m$  do
3:    $R_{TR_i} := Q_L^T X_{R_i}$  (GEMM)  $2mlnr$ 
4:    $T_i := X_{R_i} - Q_L R_{TR_i}$  (GEMM)  $2mlnr$ 
5:    $\{Q_{R_i}, R_{BR_i}\} := qr(T_i)$  (QR)  $2mnr^2$ 
6:   for  $j := 1$  to  $t$  do
7:      $z_{ij} := \begin{pmatrix} Q_L^T \\ Q_{R_i}^T \end{pmatrix} y_j$  (GEMV)  $2mtpn$ 
8:      $b_j := \begin{pmatrix} R_{TL} & R_{TR_i} \\ 0 & R_{BR_i} \end{pmatrix}^{-1} z_{ij}$  (TRSV)  $mt\mathbf{p}^2$ 
9:   end for
10: end for
```

which can be rewritten as

$$\left(\begin{array}{l} R_{TL}b_T + R_{TR_i}b_B = z_{T_j} \\ R_{BR_i}b_B = z_{B_{ij}} \end{array} \right).$$

Straightforward manipulation suffices to show that b_T and b_B can be computed as

$$\begin{aligned} b_{B_{ij}} &:= R_{BR_i}^{-1} z_{B_{ij}} \\ b_{T_{ij}} &:= R_{TL}^{-1} (z_{T_j} - R_{TR_i} b_{B_{ij}}). \end{aligned} \quad (6)$$

Given the dependencies with the loop indices, the top part of b can be partially precomputed by first moving the operation $k_j := R_{TL}^{-1} z_{T_j}$ to the same initial loop where $z_{T_j} := Q_L^T y_j$ is precomputed, and then moving the operation $h_i := R_{TL}^{-1} R_{TR_i}$ out of the innermost loop. The resulting algorithm is displayed in Alg. 2. As discussed, z_{T_j} and k_j are precomputed in an initial loop (lines 2–5), and k_j is kept in memory (a few megabytes at most) for later use within the nested double-loop (line 14). Also, the computation of h_i is taken out of the innermost loop (line 10) and reused within the innermost loop (line 14). The cost of Alg. 2 is dominated by the term $2mtrn$, that is, a factor of p/r fewer operations with respect to Alg. 1.

3 Out-of-core algorithm: Analysis of data streaming

As mentioned in Sec. 1, the second challenge to be addressed when dealing with series and grids of least squares problems is the management of large datasets. An example clarifies the situation; in the characteristic scenario in which $n = 30,000$, $p = 10$ ($r = 2$), $m = 10^7$, and $t = 10^4$, the input dataset is of size 2.4 TBs, and the computation generates 4 TBs as output. Since such datasets exceed the main memory capacity of current shared-memory nodes and therefore reside on disk, one has to resort to an out-of-core algorithm [21]. The main idea behind our design is to effectively *tile* (block) the computation to amortize the time spent in moving data.

To emphasize the need for tiling, we commence by discussing the I/O requirements of a naive out-of-core implementation of Alg. 2, as sketched in Alg. 3. First, the algorithm requires the loading of the entire set of vectors y_j (line 2) for the computations in the initial loop. Then, it loads once the entire set of matrices X_{R_i} (line 6), loads m times the entire set of vectors y_j (line 9), and finally requires the storage of the resulting $m \times t$ vectors b_{ij} (line 11). The reading of y_j for a total of m times (line 9) is the clear I/O bottleneck. For the aforementioned example, the algorithm generates 12 petabytes of disk-to-memory traffic, which at a (rather optimistic) transfer rate of 2 GBytes/sec,

Algorithm 2 : Structure of Q and R exposed and exploited

```
1:  $\{Q_L, R_{TL}\} := qr(X_L)$  (QR)  $2nl^2$ 
2: for  $j := 1$  to  $t$  do
3:    $z_{T_j} := Q_L^T y_j$  (GEMV)  $2tln$ 
4:    $k_j := R_{TL}^{-1} z_{T_j}$  (TRSV)  $tl^2$ 
5: end for
6: for  $i := 1$  to  $m$  do
7:    $R_{TR_i} := Q_L^T X_{R_i}$  (GEMM)  $2mlnr$ 
8:    $T_i := X_{R_i} - Q_L R_{TR_i}$  (GEMM)  $2mlnr$ 
9:    $\{Q_{R_i}, R_{BR_i}\} := qr(T_i)$  (QR)  $2mnr^2$ 
10:   $h_i := R_{TL}^{-1} R_{TR_i}$  (TRSM)  $m l^2 r$ 
11:  for  $j := 1$  to  $t$  do
12:     $z_{B_{ij}} := Q_{R_i}^T y_j$  (GEMV)  $2mtrn$ 
13:     $b_{B_{ij}} := R_{BR_i}^{-1} z_{B_{ij}}$  (TRSV)  $mtr^2$ 
14:     $b_{T_{ij}} := k_j - h_i b_{B_{ij}}$  (GEMV)  $2mtlr$ 
15:  end for
16: end for
```

would take 70 days of I/O. It is thus imperative to reorganize I/O and computation to greatly reduce the amount of generated I/O traffic.

Algorithm 3 : Naive out-of-core approach

```
1: Load matrix  $X_L$   $4ln \rightarrow < 1$  MBs
2: for  $j := 1$  to  $t$  do
3:   Load vector  $y_j$   $4tn \rightarrow 1.2$  GBs
4:   Compute with  $y_j$ 
5: end for
6: for  $i := 1$  to  $m$  do
7:   Load matrix  $X_{R_i}$   $4mnr \rightarrow 2.4$  TBs
8:   Compute with  $X_{R_i}$ 
9:   for  $j := 1$  to  $t$  do
10:    Load vector  $y_j$   $4mtn \rightarrow 12$  PBs
11:    Compute with  $X_{R_i}, y_j$ 
12:    Store vector  $\beta_{ij}$   $4mtp \rightarrow 4$  TBs
13:   end for
14: end for
```

3.1 Analysis of computational cost over data movement

A *tiled* algorithm decomposes the computation of the 2D grid of problems into subgrids or *tiles*. Instead of loading one single matrix X_{R_i} and vector y_j , and storing one single vector b_{ij} , the idea is to load and store multiple of them in a *slab*. The tiled algorithm presented in Alg. 4 loads slabs of t_b vectors y_j (\hat{Y}) and m_b matrices X_{R_i} (\hat{X}_R) in lines 4, 11, and 19; and stores slabs of computed $m_b \times t_b$ vectors b_{ij} (\hat{B}) in line 27.

With tiling, the amount of I/O required by line 19 of Alg. 4 is constrained to

$$t \times n \times \frac{m}{m_b},$$

Algorithm 4 : Tiled algorithm: Data loaded and stored in slabs

```

1: Load  $X_L$ 
2:  $\{Q_L, R_{TL}\} := qr(X_L)$  (QR)  $2nl^2$ 
3: for  $j := 1$  to  $t/t_b$  do
4:   Load slab  $\hat{Y}_j$ 
5:   for  $l := 1$  to  $t_b$  do
6:      $\hat{Z}_{T_{jl}} := Q_L^T \hat{Y}_{jl}$  (GEMV)  $2tln$ 
7:      $\hat{K}_{jl} := R_{TL}^{-1} \hat{Z}_{T_{jl}}$  (TRSV)  $tl^2$ 
8:   end for
9: end for
10: for  $i := 1$  to  $m/m_b$  do
11:   Load slab  $\hat{X}_{R_i}$ 
12:   for  $k := 1$  to  $m_b$  do
13:      $\hat{R}_{TR_{ik}} := Q_L^T \hat{X}_{R_{ik}}$  (GEMM)  $2mlnr$ 
14:      $\hat{T}_{ik} := \hat{X}_{R_{ik}} - Q_L \hat{R}_{TR_{ik}}$  (GEMM)  $2mlnr$ 
15:      $\{\hat{Q}_{R_{ik}}, \hat{R}_{BR_{ik}}\} := qr(\hat{T}_{ik})$  (QR)  $2mnr^2$ 
16:      $\hat{H}_{ik} := R_{TL}^{-1} \hat{R}_{TR_{ik}}$  (TRSM)  $m1^2r$ 
17:   end for
18:   for  $j := 1$  to  $t/t_b$  do
19:     Load slab  $\hat{Y}_j$ 
20:     for  $k := 1$  to  $m_b$  do
21:       for  $l := 1$  to  $t_b$  do
22:          $\hat{Z}_{B_{ik}jl} := \hat{Q}_{R_{ik}}^T \hat{Y}_{jl}$  (GEMV)  $2mtrn$ 
23:          $\hat{B}_{B_{ik}jl} := R_{BR_{ik}}^{-1} \hat{Z}_{B_{ik}jl}$  (TRSV)  $mtr^2$ 
24:          $B_{T_{ik}jl} := \hat{K}_{jl} - \hat{H}_{ik} \hat{B}_{B_{ik}jl}$  (GEMV)  $2mtlr$ 
25:       end for
26:     end for
27:     Store slab  $\hat{B}_{ij}$ 
28:   end for
29: end for

```

and, most importantly, can be adjusted by setting the parameter m_b . We note that, in terms of I/O overhead, there is no need to set m_b to the largest possible value allowed by the available main memory: it suffices to choose m_b large enough so that the I/O bottleneck shifts to the loading of X_R and writing of B :

$$t \times n \times \frac{m}{m_b} \ll m \times n \times r + m \times t \times p.$$

After choosing a sufficiently large value for m_b , the ratio of computation over data movement is

$$\frac{O(mtrn)}{O(mnr + mtp)} \equiv \frac{O(trn)}{O(nr + tp)}. \quad (7)$$

Therefore, beyond the freedom in parameterizing m_b , *it is the actual instance of the equation, that is, the problem sizes, which determines whether the computation is compute-bound or IO-bound.* We illustrate the different scenarios (memory-bound vs IO-bound) in Sec. 5, where we present experimental results.

3.2 Overlapping I/O with computation: double buffering

For problems that are not largely dominated by I/O, it is possible to reduce or even completely eliminate the overhead due to I/O operations by overlapping I/O with computation. In the development of OLS-GRID, we adopted the well-known double buffering mechanism to asynchronously load and store data: the main memory is logically split into two buffers; while operations within an iteration of the innermost loop are performed on the slabs previously loaded in one of the buffers, a dedicated I/O thread operating on the other buffer downloads the results of the previous iteration and uploads the slabs for the next one. The experiments in Sec. 5 confirm that this mechanism mitigates the negative effect of data transfers and, for compute-bound scenarios, it prevents I/O from limiting the scalability of our solver.

4 High performance and scalability

In the previous sections, we designed an algorithm that avoids redundant computations, and studied the impact of data transfers to constrain (or even eliminate) the overhead due to I/O. One last issue remains to be addressed: how to exploit shared-memory parallelism. In this section, we illustrate how to reorganize the computation so that it can be cast in terms of efficient BLAS-3 operations, and discuss how to combine different types of parallelism to attain scalability.

The computational bottleneck of Alg. 4 is the operation $Q_{R_i}^T y_j$ in line 22, which is executed $m \times t$ times. This is confirmed visually by Fig. 1, which presents, for $n = 1,000$ and varying m and t , the weight of that operation as a percent of the total computation (flops) performed by the algorithm. Already for small values of m and t (in the hundreds), the operation accounts for 90% of the computation. When m and t take more realistic values (in the thousands or more), the percent attains values close to 100%.

The operation $Q_{R_i}^T y_j$ is implemented by the inefficient BLAS-2 GEMV kernel, which on a single-core only attains about 10% of the peak performance, and on multi-cores suffers from poor scalability. The idea to overcome this bottleneck is to combine the m_b matrices $Q_{R_i}^T$ into a block $\hat{Q}_{R_i}^T$, and the t_b vectors y_j into a block \hat{Y}_j . This transformation effectively recasts the small matrix-vector products (GEMV's) into the large and much more efficient matrix-matrix operation $\hat{Q}_{R_i}^T \hat{Y}_j$ (GEMM), which achieves both close-to-optimal performance and high scalability.

While the use of a multi-threaded version of the BLAS library for the GEMM operation $\hat{Q}_{R_i}^T \hat{Y}_j$ enables high scalability, it is a poor choice for other sections of the algorithm, which do not scale

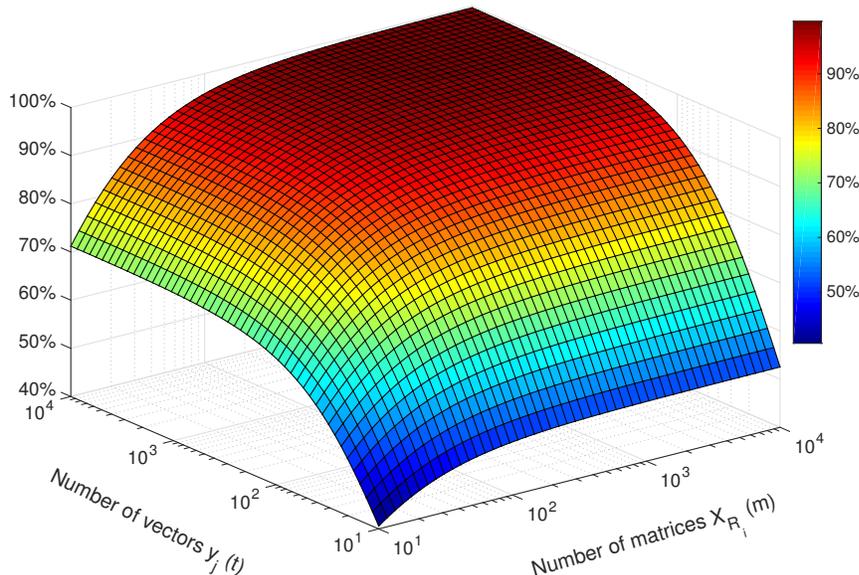


Fig. 1. Percentage of the computation performed by the operation $Q_{R_i}^T y_j$ from the total operations required by Algorithm 4.

as nicely. In fact, when using a large number of cores (as is the case in our experiments), the weight of these other sections increases to the point that it affects the overall scalability. As a solution to mitigate this problem we turn to a hybrid parallelism combining multi-threaded BLAS with OpenMP. Lines 13, 14, and 16 in Alg. 4 are examples of operations that do not scale with a mere use of a multi-threaded BLAS. Let us illustrate the issue with the matrix product in line 13. One of these GEMM's in isolation multiplies $Q_L^T \in R^{l \times n}$ with $X_{R_i} \in R^{n \times r}$; that is, it multiplies a wide matrix with only a few rows times a thin matrix with only a few columns. Even when m_b X_{R_i} 's are combined into the block \hat{X}_{R_i} , one of the matrices (Q_L) is still rather small. In terms of number of operations per element, the ratio is very low, which results in an inefficient and poorly scalable operation. In this case, using a multi-threaded BLAS is not sufficient (in our experiments we observed a performance of about 5 GF/s, i.e., below 1% efficiency). Instead, we decide to explicitly split the operation among the compute threads by means of OpenMP directives. Each thread takes a proportional number of columns from \hat{X}_{R_i} and computes the corresponding GEMM using a single-threaded BLAS call. This second alternative results in a speedup of 5x to 10x, sufficient to mitigate the impact in the overall performance and scalability.

The resulting algorithm is displayed in Alg. 5. Among the computational bottlenecks, we highlight in green the operations parallelized by means of a multi-threaded version of BLAS (line 18) and in blue the operations parallelized using OpenMP (lines 10–12). As we show in the next section, the recasting of small operations into larger ones and the use of a hybrid form of parallelism that combines multi-threaded BLAS and OpenMP leads to satisfactory performance and scalability signatures.

Algorithm 5 : Efficient and scalable algorithm: OLS-GRID

```

1: Load  $X_L$ 
2:  $\{Q_L, R_{TL}\} := qr(X_L)$  (QR)  $2nl^2$ 
3: for  $j := 1$  to  $t/tb$  do
4:   Load slab  $\hat{Y}_j$ 
5:    $\hat{Z}_{T_j} := Q_L^T \hat{Y}_j$  (GEMM)  $2tln$ 
6:    $\hat{K}_j := R_{TL}^{-1} \hat{Z}_{T_j}$  (TRSM)  $tl^2$ 
7: end for
8: for  $i := 1$  to  $m/mb$  do
9:   Load slab  $\hat{X}_{R_i}$ 
10:   $\hat{R}_{TR_i} := Q_L^T \hat{X}_{R_i}$  (GEMM)  $2mlnr$ 
11:   $\hat{T}_i := \hat{X}_{R_i} - Q_L \hat{R}_{TR_i}$  (GEMM)  $2mlnr$ 
12:   $\hat{H}_i := R_{TL}^{-1} \hat{R}_{TR_i}$  (TRSM)  $ml^2r$ 
13:  for  $k := 1$  to  $m_b$  do
14:     $\{\hat{Q}_{R_{i_k}}, \hat{R}_{BR_{i_k}}\} := qr(\hat{T}_{i_k})$  (QR)  $2mnr^2$ 
15:  end for
16:  for  $j := 1$  to  $t/tb$  do
17:    Load slab  $\hat{Y}_j$ 
18:     $\hat{Z}_{B_{i_j}} := \hat{Q}_{R_i}^T \hat{Y}_j$  (GEMM)  $2mtrn$ 
19:    for  $k := 1$  to  $m_b$  do
20:       $\hat{B}_{B_{i_k j}} := R_{BR_{i_k}}^{-1} \hat{Z}_{B_{i_k j}}$  (TRSM)  $mtr^2$ 
21:       $B_{T_{i_k j}} := \hat{K}_j - \hat{H}_{i_k} \hat{B}_{B_{i_k j}}$  (GEMM)  $2mtlr$ 
22:    end for
23:    Store slab  $\hat{B}_{ij}$ 
24:  end for
25: end for

```

5 Experimental results

In this section, OLS-GRID (the implementation of Alg. 5) is tested in a variety of settings to provide evidence that it makes a nearly optimal use of the available resources; it is also compared and contrasted with a number of available alternatives, namely a naive solver (the linear regression solver from ProbABEL), a specialized solver for grids of GLS (generalized least squares) problems (OmicABEL-Eig), and a specialized solver for grids of OLS problems (GWASP).

All our tests were run on a system consisting of 4 Intel(R) Xeon(R) E7-4850 Westmere-EX multi-core processors. Each processor comprises 10 cores operating at a frequency of 2.00 GHz, for a combined peak performance in single precision of 640 GFlops/sec. The system is equipped with 256 GBs of RAM and 8 TBs of disk as secondary memory; our measurements indicate that the Lustre file system attains a maximum bandwidth of about 300 MBs/sec and 1.7 GBs/sec for writing and reading operations, respectively. The solver was compiled using Intel’s icc compiler (v14.0.1), and linked to Intel’s MKL multi-threaded library (v14.0.1), which provides for BLAS and LAPACK functionality. The routine makes use of the OpenMP parallelism provided by the compiler through a number of *pragma* directives. All computations were performed in single precision, and the correctness of the results was assessed by direct comparison with the OLS solver from LAPACK (SGELS).

For the asynchronous I/O transfers, we initially tested the AIO library (available in all Unix systems). We observed that I/O operations had a considerable impact in performance by limiting the flops attained by GEMM. Since AIO does not offer the means to specify the amount of threads spawned for I/O and to pin threads to cores, we developed our own light-weight library which uses one single thread for I/O operations and allows thread pinning.³

5.1 Compute-bound vs IO-bound scenarios

We commence the study of our own solver by showing the practical implications of the ratio computation over I/O transfers

$$\frac{O(trn)}{O(nr + tp)}.$$

We collected the time spent in computations and the time spent in I/O operations for a range of problems varying the size of n . Figure 2 presents these results. The horizontal red line corresponds to a ratio of 1, while the other three lines represent the ratio compute time over I/O time for different values of n and for an increasing number of compute threads.

While the bandwidth is fixed, the computational power increases with the number of compute threads. For $n = 5,000$ and $n = 20,000$, the I/O time eventually overtakes computation time (the ratio becomes < 1). At that point, the problem becomes IO-bound and the use of additional compute threads does not reduce the time-to-solution. For $n = 30,000$ instead, the time spent in I/O never grows larger than the time spent in computation and I/O is perfectly overlapped.

In the next experiment we show how the results in Fig. 2 translate into scalability results.

5.2 Scalability

We study now the scalability of OLS-GRID. Figure 3 presents scalability results for the same problem sizes discussed in the previous experiment. The red diagonal line represents perfect scalability. As

³ The library is publicly available at <http://hpc.rwth-aachen.de/~fabregat/software/lwaio-v0.1.tgz>.

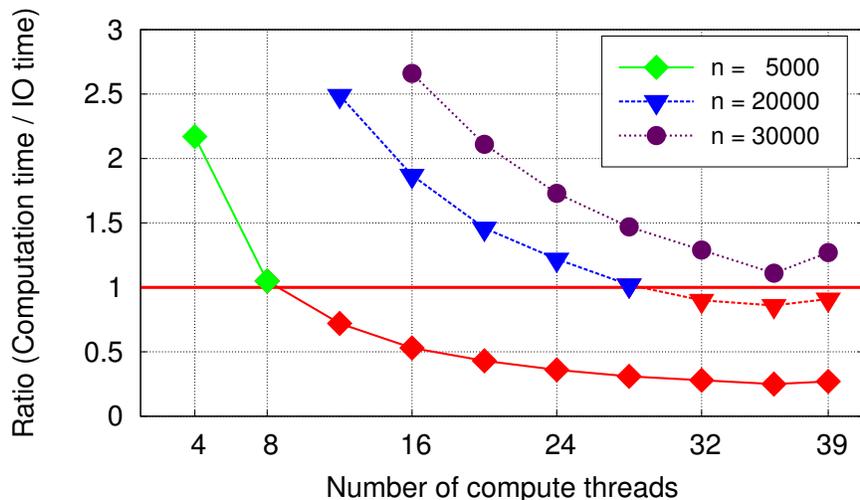


Fig. 2. Ratio of compute time over I/O time. Results for a varying value of n and an increasing number of compute threads. The other problem sizes are fixed: $l = 4$, $r = 1$, $m = 10^6$, $t = 10^4$, $m_b = 10^4$, and $t_b = 5 \times 10^3$.

the ratio in Fig. 2 suggested, the line for $n = 5,000$ shows perfect scalability for up to 8 threads; from that point on, the I/O transfers dominate the execution time and no larger speedups are possible. Similarly, for $n = 20,000$, almost perfect scalability is attained with up to 28 compute threads. Again, beyond that number, the I/O dominates and the scalability plateaus. Instead, for $n = 30,000$, the I/O never dominates execution time and the solver attains speedups of around 34x when 36 compute threads are used. In all cases, we observe a drop in scalability when 40 compute threads are used; this is because one compute thread and the I/O thread share one core. We attribute the drop at 39 compute threads for the experiment with $n = 30,000$ —and, less apparently, at $n = 5,000$ and $n = 20,000$ —to potential memory conflicts, since in the used architecture each pair of threads share the L1 and L2 caches.

An important message to extract from these results is the need to understand the characteristics of the problem at hand to decide which architecture fits best our needs. In the case of omics GWAA, the size of n plays an important role in the decision of whether investing in further computational power or larger bandwidth.

5.3 Efficiency

As a final result, we quantify how efficiently OLS-GRID uses the available resources. To this end, we measured the time-to-solution for a variety of scenarios using 36 compute threads and compared the performance with that of the practical peak performance, that is, the best performance attained by GEMM. The results are presented in Fig. 4; the tested scenarios include all combinations of $n = (5,000, 30,000)$, $p = (5, 10)$, $m = (10^6, 10^7)$, and $t = (10^3, 10^4)$. While for small sizes of n the I/O transfers limit the efficiency, with a maximum of 0.3, for large values of n the attained efficiency ranges from 0.64 to 0.93.

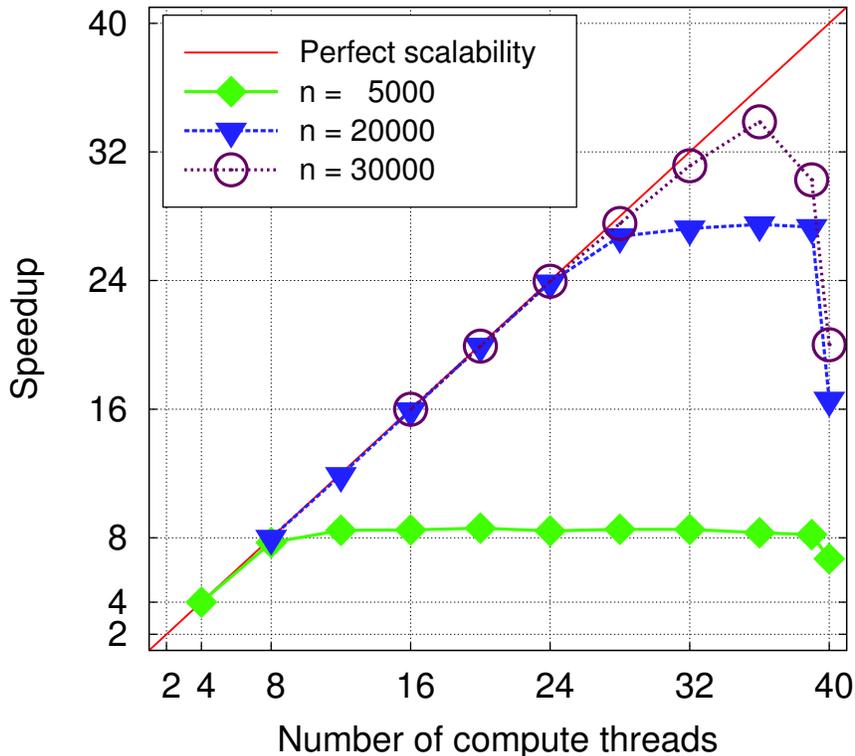


Fig. 3. Scalability results for a varying value of n . The other problem sizes are fixed: $l = 4$, $r = 1$, $m = 10^6$, $t = 10^4$, $m_b = 10^4$, and $t_b = 5 \times 10^3$.

5.4 Comparison with alternative solvers

We conclude the performance study with a detailed discussion on how OLS-GRID compares to other existing solvers. All presented timings were obtained in the same experimental setup as the previous sections, and making use of the available 40 cores.

ProbABEL is a collection of tools for linear and logistic regression, as well as Cox proportional hazards models [16]. While in this discussion we concentrate on its deficiencies in terms of performance for large-scale linear regression analyses, we want to clarify that this is a versatile tool offering functionality well beyond the core linear regression solver. ProbABEL’s OLS solver does not take advantage of the structure of the problem at hand, and thus delivers poor performance. For instance, an analysis with sizes $n = 30,000$, $p = 5$ ($l = 4$, $p = 1$), $m = 10,000$, and $t = 100$ takes almost half an hour, while OLS-GRID completes in less than two seconds. The gap in performance between a traditional black-box solver and a carefully designed solver is enormous: if the previous analysis were extended to $m = 10^7$ and $t = 10^4$, ProbABEL’s algorithm would become unusable (estimated completion time of more than 13 years), while our OLS-GRID completes the analysis in 6.9 hours.

OmicABEL-Eig is the multi-trait ($t > 1$) solver for GLS-based analyses in OmicABEL [5,15]. The concepts behind OmicABEL-Eig are similar to those discussed in this paper; in particular, the algorithm relies on the reduction of GLS problems to OLS ones. However, the dependencies in the

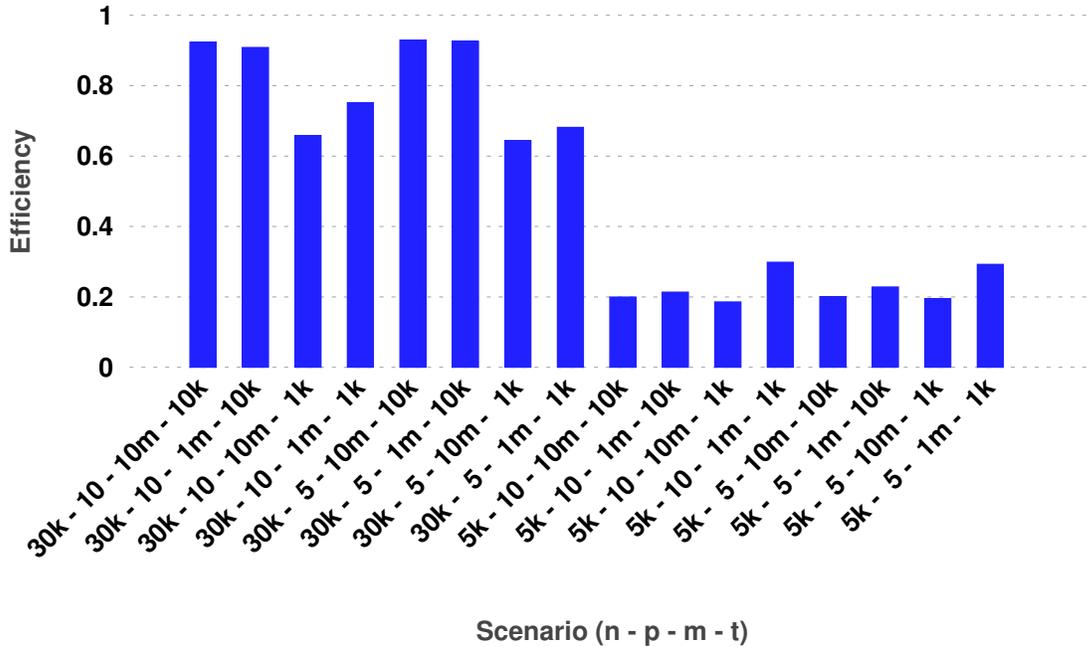


Fig. 4. Efficiency of our solver for a range of scenarios using 36 compute threads. GEMM’s peak performance is used as practical peak performance of the architecture.

grid to be solved are different, and computationally the consequences are huge. On the one hand, while the asymptotic complexity of both OmicABEL-Eig and OLS-GRID is the same, the constants for the former are large. On the other hand, even though OmicABEL-Eig shows great scalability, its efficiency compared to the present solver is low because it cannot take advantage of GEMM. As a result, the same analysis that took OLS-GRID 6.9 hours to complete, would require a month of computation for OmicABEL-Eig. In short, even the use of an existing fully optimized solver for a similar regression analysis may not be the best choice, and it is worth the effort of designing a new one tailored to the needs of the specific problem at hand.

GWASP is a recently developed solver tailored for the computation of a grid of OLS problems similar to that addressed in this paper [17]. Thanks to techniques such as the aggregation of vector-vector and matrix-vector operations into matrix-vector and matrix-matrix operations, respectively, GWASP delivers an efficient algorithm. Unfortunately, a direct comparison with OLS-GRID is not possible, since 1) this solver only computes the last r entries of β_{ij} , 2) the focus is on single-trait ($t = 1$) analyses, and 3) only R code is provided.

In order to study GWASP’s performance, we implemented the algorithm in C using the BLAS and LAPACK libraries, and ran a number of single-trait experiments. For a problem of size $n = 5,000$, $p = 5$ ($l = 4, r = 1$), $m = 10^6$, and $t = 1$, GWASP took 61 seconds, while OLS-GRID completed in 14 seconds. Similarly, if n is increased to $n = 30,000$, the analysis completes in 314 and 84 seconds for GWASP and OLS-GRID, respectively. Since the computational complexity of both solvers is similar and these scenarios are IO-bound, the differences mainly reside in the careful hybrid parallelization of OLS-GRID and the overlapping of data transfers.

For the general case of $t > 1$, GWASP may only be used one trait at a time. In the last scenario with $t = 1,000$, it would require days to complete. By contrast, OLS-GRID finishes in 138 seconds,

thanks to the optimizations enabled when considering the 2D grid of problems as a whole. Despite this relatively large gap, we believe that GWASP can be extended to multi-trait analyses and to the computation of entire β 's, while achieving performance comparable to that of our solver.

6 Conclusions

We addressed the design and implementation of efficient solvers for large-scale linear regression analyses. As case study, we focused on the computation of a two-dimensional grid of ordinary least squares problems as it appears in the context of genome-wide association analyses. The resulting routine, OLS-GRID, showed to be highly efficient and scalable.

Starting from the mathematical description of the problem, we designed an incore algorithm that exploits the available problem-specific knowledge and structure. Next, to enable the solution of problems with large datasets that do not fit in today multi-core's main memory, we transformed the incore algorithm into an out-of-core one, which, thanks to tiling, constrains the amount of required data movement. By incorporating the double-buffering technique and using an asynchronous I/O library, we completely eliminated I/O overhead whenever possible.

Finally, by reorganizing the calculations, the computational bottleneck was cast in terms of the highly efficient GEMM routine. Combining multi-threaded BLAS and OpenMP parallelism, OLS-GRID also attains high performance and high scalability. More specifically, for large enough analyses, the solver achieves single-core efficiency beyond 90% of peak performance, and speedups of up to 34x with 36 cores.

While previously existing tools allow for multiple trait analysis, these are limited to small datasets and require considerable runtimes. We enable the analysis of previously intractable datasets and offer shorter times to solution. Our OLS-GRID solver is already integrated in the GenABEL suite, and adopted by a number of research groups.

6.1 Future work

Two main research directions remain open. On the one hand, support for distributed-memory architectures is desirable, allowing for further reduction in time-to-solution. This step will require a careful distribution of workload among nodes and the use of advanced I/O techniques to prevent data movement from becoming a bottleneck.

On the other hand, we are already working on the support for so-called *missing* and *erroneous data*. Often, large datasets are collected by combining data from different sources. Therefore, subsets of data for certain individuals (for instance, entries in y_j vectors) may not be available and labeled as missing values (NaNs). Erroneous data may origin, for instance, from errors in the acquisition. Accepting input data with NaN values (and preprocessing it accordingly) will make our software of broader appeal and will facilitate its adoption by the broad GWAA community.

Acknowledgments

Financial support from the Deutsche Forschungsgemeinschaft (German Research Association) through grant GSC 111 is gratefully acknowledged. The authors thank Yurii Aulchenko for fruitful discussions on the biological background of GWAA.

References

1. L. A. Hindorff, P. Sethupathy, H. A. Junkins, E. M. Ramos, J. P. Mehta, F. S. Collins, T. A. Manolio, Potential etiologic and functional implications of genome-wide association loci for human diseases and traits, Proc. Natl. Acad. Sci. USA 106 (23) (2009) 9362–9367. doi:10.1073/pnas.0903103106.

2. X. Sala-I-Martin, I just ran two million regressions, *The American Economic Review* 87 (2) (1997) 178–183.
3. C. Gieger, L. Geistlinger, E. Altmaier, M. Hrabé de Angelis, F. Kronenberg, T. Meitinger, H.-W. W. Mewes, H.-E. E. Wichmann, K. M. Weinberger, J. Adamski, T. Illig, K. Suhre, Genetics meets metabolomics: a genome-wide association study of metabolite profiles in human serum., *PLoS Genet* 4 (11) (2008) e1000282+.
4. A. Demirkan, C. M. Van Duijn, P. Ugoicsai, A. Isaacs, P. P. Pramstaller, G. Liebisch, J. F. Wilson, . Johansson, I. Rudan, Y. S. Aulchenko, et al., Genome-wide association study identifies novel loci associated with circulating phospho- and sphingolipid concentrations., *PLoS Genetics* 8 (2) (2012) e1002490.
5. D. Fabregat-Traver, S. Sharapov, C. Hayward, I. Rudan, H. Campbell, Y. Aulchenko, P. Bientinesi, Big-Data, High-Performance, Mixed Models Based Genome-Wide Association Analysis with omicABEL software, *F1000Research* 3 (200).
6. Y. S. Aulchenko, S. Ripke, A. Isaacs, C. M. van Duijn, GenABEL: an R library for genome-wide association analysis., *Bioinformatics* 23 (10) (2007) 1294–6.
7. L. Hindorff, J. MacArthur, A. Wise, H. Junkins, P. Hall, A. Klemm, T. Manolio, A catalog of published genome-wide association studies, available at: www.genome.gov/gwastudies. Accessed April 12th (2015).
8. G. S. Hageman, et al., A common haplotype in the complement regulatory gene factor H (HF1/CFH) predisposes individuals to age-related macular degeneration, *Proc. Natl. Acad. Sci. USA* 102 (20) (2005) 7227–7232.
9. A. O. Edwards, R. Ritter, K. J. Abel, A. Manning, C. Panhuysen, L. A. Farrer, Complement factor H polymorphism and age-related macular degeneration, *Science* 308 (5720) (2005) 421–424.
10. T. M. Frayling, Genome-wide association studies provide new insights into type 2 diabetes aetiology, *Nat Rev Genet* 8 (9) (2007) 657–662.
11. M. A. Nalls, et al., Large-scale meta-analysis of genome-wide association data identifies six new risk loci for Parkinson’s disease, *Nat Genet* 46 (9) (2014) 989–993, letter.
12. C. Lippert, J. Listgarten, Y. Liu, C. M. Kadie, R. I. Davidson, D. Heckerman, Fast linear mixed models for genome-wide association studies, *Nat. Methods* 8 (10) (2011) 833–835.
13. X. Zhou, M. Stephens, Genome-wide efficient mixed-model analysis for association studies, *Nat. Genet.* 44 (7) (2012) 821–824.
14. D. Fabregat-Traver, Y. S. Aulchenko, P. Bientinesi, Solving sequences of generalized least-squares problems on multi-threaded architectures, *Applied Mathematics and Computation (AMC)* 234 (2014) 606–617.
15. D. Fabregat-Traver, P. Bientinesi, Computing petaflops over terabytes of data: The case of genome-wide association studies, *ACM Trans. Math. Software* 40 (4) (2014) 27:1–27:22.
16. Y. S. Aulchenko, M. V. Struchalin, C. M. van Duijn, ProbABEL package for genome-wide association analysis of imputed data, *BMC Bioinformatics* 11 (2010) 134.
17. K. Sikorska, E. Lesaffre, P. F. J. Groenen, P. H. C. Eilers, GWAS on your notebook: fast semi-parallel linear and logistic regression for genome-wide association studies, *BMC Bioinformatics* 14 (2013) 166.
18. A. Voorman, K. Rice, T. Lumley, Fast computation for genome-wide association studies using boosted one-step statistics, *Bioinformatics* 28 (14) (2012) 1818–1822.
19. G. H. Golub, C. F. Van Loan, *Matrix computations* (3rd ed.), Johns Hopkins University Press, Baltimore, MD, USA, 1996.
20. D. Fabregat-Traver, P. Bientinesi, Knowledge-based automatic generation of Partitioned Matrix Expressions, in: *Computer Algebra in Scientific Computing*, Vol. 6885 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2011, pp. 144–157.
21. S. Toledo, *External memory algorithms*, American Mathematical Society, Boston, MA, USA, 1999, Ch. A survey of out-of-core algorithms in numerical linear algebra, pp. 161–179.