

rMVP: A Memory-efficient, Visualization-enhanced, and Parallel-accelerated tool for Genome-Wide Association Study

Lilin Yin^{1,2,#}, Haohao Zhang^{3,#}, Zhenshuang Tang^{1,2}, Jingya Xu^{1,2}, Dong Yin^{1,2}, Zhiwu Zhang⁴, Xiaohui Yuan³, Mengjin Zhu^{1,2}, Shuhong Zhao^{1,2}, Xinyun Li^{1,2*}, and Xiaolei Liu^{1,2*}

¹Key Laboratory of Agricultural Animal Genetics, Breeding and Reproduction, Ministry of Education & College of Animal Science and Technology, Huazhong Agricultural University, Wuhan, Hubei, 430070, PR China;

²Key Laboratory of Swine Genetics and Breeding, Ministry of Agriculture, Huazhong Agricultural University, Wuhan 430070, Hubei, PR China;

³School of Computer Science and Technology, Wuhan University of Technology, Wuhan, China;

⁴Department of Crop and Soil Sciences, Washington State University, Pullman, Washington, USA.

#Equal contribution.

*Corresponding author(s).

Email: xiaoleiliu@mail.hzau.edu.cn (Liu X), hzaulxy@163.com (Li X).

Abstract

Along with the development of high-throughput sequencing technologies, both sample size and number of SNPs are increasing rapidly in Genome-Wide Association Studies (GWAS) and the associated computation is more challenging than ever. Here we present a Memory-efficient, Visualization-enhanced, and Parallel-accelerated R package called “rMVP” to address the need for improved GWAS computation. rMVP can: (1) effectively process large GWAS data; (2) rapidly evaluate population structure; (3) efficiently estimate variance components by EMMAX, FaST-LMM, and HE regression algorithms; (4) implement parallel-accelerated association tests of markers using GLM, MLM, and FarmCPU methods; (5) compute fast with a globally efficient design in the GWAS processes; and (6) generate various visualizations of GWAS related information. Accelerated by block matrix multiplication strategy and multiple threads, the association test methods embedded in rMVP are approximately 5-20 times faster than PLINK, GEMMA, and FarmCPU_pkg. rMVP is freely available at <https://github.com/xiaolei-lab/rMVP>.

Key words: Memory-efficient, Visualization-enhanced, Parallel-accelerated, rMVP, GWAS

Introduction

The computation burden of GWAS is partially caused by the increasing sample size and marker density applied for these studies. As a result, how to efficiently analyse the big data is a big challenge. Additionally, GWAS have been widely used for detecting candidate genes that control human diseases and agricultural economic traits, where the accuracy of the results is of significant implications. Thus, how to achieve higher statistical power under a reasonable level of type I error is another challenge[1]. To efficiently detect more candidate genes with lower false positive rates is the current working goal for GWAS algorithms and tools[2, 3].

39 Introducing the population structure concept into GWAS has dramatically improved
40 accuracy of detection. For example, incorporating the fractions of individuals belonging to
41 subpopulations, namely Q matrix, reduces both false positive and false negative signals[4].
42 Principal components (PCs) are widely used to represent subpopulations and to enable the
43 incorporation of population structure into GWAS[5]. Implementing the General Linear
44 Model (GLM) to incorporate either the Q matrix or PCs as covariates, PLINK has become the
45 most popular software package for GWAS[6].

46 False positives also stem from individuals that exhibit high variability in pairwise
47 relatedness presumptively classified into different subpopulations. In addition to integrating
48 population structure, statistical power can be substantially improved by the incorporation of
49 hidden relationships in a mixed linear model (MLM) - particularly when population structure
50 is less dominant than the cryptic relatedness[7]. Multiple algorithms have been developed to
51 boost both the computational efficiency and statistical power of MLM methods[8-11].
52 Various software packages have also been developed for the implementation of these
53 algorithms, including TASSEL[12], GAPIT[13, 14], GenABEL[15], EMMAX[16],
54 GEMMA[17], and GCTA[18]. Even though the number of GWAS literature applying MLM-
55 based methods is increasing rapidly, it is still not comparable to that of PLINK, primarily
56 because PLINK operates much faster than MLM-based method software.

57 Besides the difference in computing time, MLM does not provide high statistical
58 power compared to GLM. The difference in statistical power between GLM and MLM is
59 negligible in some scenarios, such as mapping genes under the same false discovery rate in
60 populations with strong population structure[19]. These populations include human
61 populations, as well as animal and plant populations that have been isolated by breeding
62 programs. Our newly developed method, FarmCPU (Fixed and random model Circulating
63 Probability Unification) has higher statistical power than both GLM and MLM for evaluating
64 populations with either weak or strong population structure. FarmCPU splits MLM into a
65 fixed effect model (FEM) and a random effect model (REM), using them iteratively to
66 increase the power for detecting candidate genes associated with population structure.
67 Association tests in FarmCPU are validated by FEM with the same computing efficiency as
68 GLM while the statistical power surpasses that of MLM at the same level of type I error.

69 Although recently developed methods have improved statistical power under certain
70 assumptions, determining the most appropriate method for a given dataset is still convoluted.
71 Human genetic studies often use large datasets with simple models, while plant and animal
72 genetic studies prefer complex models with limited sample sizes. For a specific trait, it is

73 usually difficult to identify the real genetic architecture and the most appropriate method to
74 be used. Researchers have to try out multiple methods and identify candidate genes based on
75 both statistical and biological evidence. Additionally, existing GWAS software rarely focuses
76 on providing a flexible plotting function to display GWAS related information in a way that
77 satisfies the personal aesthetic requirements of the researchers. Furthermore, with the
78 development of multi-traits methods, such as GSA-SNP2[20], MTMM[21], mvLMMs[22],
79 and mtSet[23], results from multiple-group GWAS need to be displayed simultaneously for
80 easier comparisons. Therefore, there appears a need for analysing big data with limited
81 computing memory, reasonable time, and reduced false positive rates, while displaying the
82 results in high-quality figures. To address all of the above requirements, we developed the
83 Memory-efficient, Visualization-enhanced, and Parallel-accelerated package (rMVP) in R.

84 **Methods and materials**

85 We split the entire GWAS procedure into six sections: data preparation, evaluation of
86 population structure, estimation of variance components, association tests, globally efficient
87 design on GWAS process computing, and result visualization. Abundant functions have been
88 implemented in rMVP for each section:

89 **(1) Data preparation.** rMVP accepts multiple popular formats for genotype files, such as
90 PLINK binary, Hapmap, VCF, and Numeric (e.g., genotype data can be coded as integer (0, 1,
91 2) or dosage/probability (0.1, 0.3, 0.6)). All above formats will be converted to the
92 ‘big.matrix’ format. The advantage of converting genotype files into ‘big.matrix’ is that the
93 size of the file is only limited by the storage capacity of the hard disk but not the processing
94 capacity of Random Access Memory (RAM, and ‘memory’ is referred to RAM in this
95 manuscript)[24].

96 **(2) Evaluation of population structure and individual relationship.** For population
97 structure analysis, PCs can be calculated using all available markers. An ideal population for
98 GWAS assumed that the individuals were randomly selected from a big population, but the
99 population could always be classified to multiple subpopulations in fact. The alleles with
100 different frequencies in different subpopulations would generate false positives, we
101 recommend to integrate the 3-5 top PCs as covariates into model to control false positives
102 caused by population structure following previous studies[5, 19]. VanRaden Method is
103 implemented in rMVP for the efficient construction of genomic relationship matrix
104 (GRM)[25].

105 **(3) Estimation of variance components.** Four algorithms are implemented for estimating
106 variance components in rMVP: Brent (default method in rMVP)[26], EMMAX (Efficient

107 Mixed-Model Association eXpedited) / P3D (Population Parameters Previously
108 Determined)[8, 16], FaST-LMM (Factored Spectrally Transformed Linear Mixed Model)
109 algorithms[9] and HE regression (Haseman-Elston regression)[27]. Different algorithms have
110 different performances in terms of accuracy and efficiency. For instance, the Brent and
111 EMMAX use eigen decomposition on genomic relationship matrix to avoid computing the
112 inverse of big matrix, the FaST-LMM uses singular value decomposition on genotype matrix,
113 which can be more efficient when the number of markers is far less than the number of
114 individuals, the HE regression, which uses linear regression model to fit the similarity of
115 phenotype and genomic relationship matrix among individuals, is less accurate but can be
116 much more memory-efficient and time-saving, making it more promising in very large
117 datasets.

118 **(4) Association tests.** General Linear Model, Mixed Linear Model, and FarmCPU methods
119 are implemented in rMVP for association tests. When there is more than one covariate (e.g.
120 PCs) added to association test models, the inverse of the design matrix corresponding to the
121 covariates will be calculated n times, where n is marker size. Block matrix multiplication
122 strategy can be used to speed up the processes including inverse of the design matrix
123 corresponding to the covariates and the testing markers. This strategy is used in all available
124 association test methods in rMVP. Take GLM as an example, the fixed effect model function
125 can be written as:

126
$$y = Xb + e \dots\dots\dots (1)$$

127 where y is a vector of phenotype, X is a matrix of fixed effects and test SNP, b is
128 an incidence matrix for X , and e is a vector of residual that followed a normal distribution
129 with mean of zero and $I\sigma_e^2$ covariance, where I is the identity matrix and σ_e^2 is the
130 unknown residual variance. Equation (1) can be reformulated by following steps:

131
$$X'y = X'Xb$$

132
$$b = (X'X)^{-1} X'y \dots\dots\dots (2)$$

133 Where X' is the transpose matrix of X . If there are k fixed effect vectors added as
134 covariates in the model, X and b can be written as:

135
$$X = [C'_1, C'_2, C'_3, \dots, C'_k, SNP']$$

136
$$b = [b_1, b_2, b_3, \dots, b_k, c]$$

137 where $C_1, C_2, C_3, \dots, C_k$ represent k fixed effect vectors and SNP represents the test
138 SNP vector. Equation (2) can be written as

$$139 \quad \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \dots \\ b_k \\ c \end{bmatrix} = \left(\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ \dots \\ C_k \\ SNP \end{bmatrix} [C'_1, C'_2, C'_3, \dots, C'_k, SNP'] \right)^{-1} \begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ \dots \\ C_k \\ SNP \end{bmatrix} [y'] \dots \dots \dots (3)$$

140 The most time-consuming part in equation (3) is the inverse of M matrix. And M is
141 defined as:

$$142 \quad M = \begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ \dots \\ C_k \\ SNP \end{bmatrix} [C'_1, C'_2, C'_3, \dots, C'_k, SNP']$$

143 If we use w and x represent $C_1, C_2, C_3, \dots, C_k$ and SNP , respectively, the inverse of
144 M matrix can be written as:

$$145 \quad M^{-1} = \left(\begin{bmatrix} w' \\ x' \end{bmatrix} [w, x] \right)^{-1} = \begin{bmatrix} w'w & w'x \\ x'w & x'x \end{bmatrix}^{-1} = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix}$$

146 where

$$147 \quad M_{11} = (w'w)^{-1} + (w'w)^{-1}w'x(x'x - x'w(w'w)^{-1}w'x)^{-1}x'w(w'w)^{-1}$$

$$148 \quad M_{12} = -(w'w)^{-1}w'x(x'x - x'w(w'w)^{-1}w'x)^{-1}$$

$$149 \quad M_{21} = -(x'x - x'w(w'w)^{-1}w'x)^{-1}x'w(w'w)^{-1}$$

$$150 \quad M_{22} = (x'x - x'w(w'w)^{-1}w'x)^{-1}$$

151 The inversion of $w'w$ matrix is recomputed n times when constructing $M_{11}, M_{12}, M_{21}, M_{22}$
152 matrix for each test marker. For the matrix operations in GLM, MLM, and each iteration of
153 FarmCPU, the w matrix is fixed, and the inversion of $w'w$ can be calculated only once using
154 block matrix multiplication strategy. As it is repeated n times when testing the SNPs, more
155 time will be saved when there are more covariates in the model or more SNPs to be tested.

156 **(5) Globally efficient design of GWAS calculations.** A standard GWAS pipeline generally
157 includes PC derivation, GRM construction, variance components estimation, and association
158 tests. There are three commonly used strategies for deriving the PCs: (a) the Eigen
159 decomposition results of the matrix that represents the correlation coefficients between pairs
160 of markers could be derived by $(M^T M)_v = \lambda_v$, where M is a n by m scaled genotype matrix

161 (n is the number of individuals, m is the number of SNPs), the Eigen decomposition analysis
162 is conducted on the correlation matrix $M^T M$, the dimension of which is m by m , and this
163 would lead to high requirements of both memory and computing time with the increasing
164 number of SNPs; (b) The Singular Value Decomposition analysis could be conducted on the
165 M matrix by $M = U \Sigma V^*$, its computational complexity is relative smaller than the method
166 that described in (1), as it only needs to decompose a n by m matrix; (c) The PCs could be
167 also derived by performing the Eigen decomposition of the GRM, which could be calculated
168 by $GRM = M^T M / m$, and its dimension is n by n . In the majority of cases, the number of
169 markers (m) is far more than the number of individuals (n), this method has the smallest
170 computational complexity compared with the other two. Moreover, the construction of GRM
171 is always a key part in commonly used GWAS procedure, which has been precomputed
172 already. Not only that, as shown in Figure S3, the Eigen decomposition results of GRM could
173 be easily applied to processes of variance components estimation and association tests. By the
174 default sets in rMVP, the Eigen decomposition analysis was conducted on GRM, which was
175 constructed by VanRaden method[25], the methodological formula of VanRaden method can
176 be defined as:

177
$$G = Z^T Z / \sum_{i=1}^n p_i (1 - p_i) \dots\dots\dots (4)$$

178 Where Z is a dimension of m by n matrix, m is the number of markers and n is the number of
179 individuals, it can be derived from centering the additive genotype matrix which was coded
180 as 0, 1, 2 for genotype AA, AB, BB respectively, p is the minor allele frequency. After the
181 eigen decomposition was finished, the eigen values and eigen vectors could be applied to the
182 estimation of variance components using Brent method[26], which has fast convergence
183 determined via the absolute tolerance of heritability rather than all variance components; and
184 the results of eigen decomposition could be also used for solving the mixed model equation
185 when MLM is selected for the association tests. The globally efficient calculation design of
186 GWAS process makes rMVP only need to do the eigen decomposition once instead of doing
187 it multiple times, its results could be directly used in calculations of PC derivation, variance
188 components estimation, and association tests, and the computing time is greatly decreased.

189 **(6) Visualization of results.** High-quality figures are generated to display data information,
190 population structure and GWAS results, including marker density plot, phenotype distribution
191 plot, PCA plot, Manhattan plot, and Q-Q (Quantile-Quantile) plot.

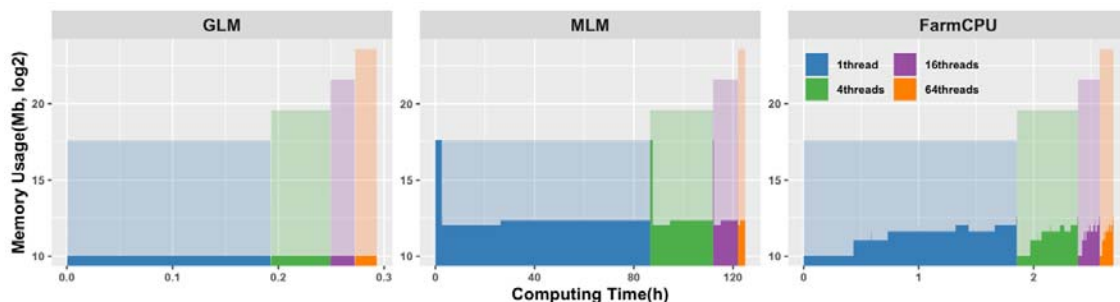
192

193 Results

194 Memory-efficient: Efficient memory usage in data loading and parallel computation

195 Genotype matrices are the biggest datasets for GWAS. In rMVP, genotype data in multiple
196 formats are converted to 'big.matrix', which can minimize RAM usage through generating a
197 bridge that facilitates RAM accessing the data on the hard disk instead of loading it to RAM
198 directly as the most software tools do. rMVP achieves this goal by using the 'bigmemory'
199 package to build data mirrors that are accessible to RAM, while the actual data remain on the
200 hard drive. In this way, very little RAM capacity is needed for the temporary storage of the
201 data. Once the data mirrors are built, users will never need to re-build them again and the
202 time of loading input data is negligible. When multiple threads are used to accelerate the
203 association tests, no additional data mirrors will be copied for each thread as all threads will
204 share the same data mirrors.

205 Here, we made a rough illustration of 'big.matrix' based memory storage of one and
206 multiple threads for rMVP. The complete GWAS procedure of three methods was recorded
207 for RAM usage test in a Linux server ('RES' – 'SHR'). In this test, the product of genotype
208 data size was measured in standard R matrix format, and 'theoretical RAM cost' for multiple
209 threads in 'fork' mode is defined as $r \times c \times t \times 8$ bytes, where r and c are the number of rows
210 and columns of a matrix respectively, t is the number of threads. From the Figure 1, we
211 concluded that, with more threads, rMVP shares variables in RAM among processers and
212 does not require additional memory compared with single thread by the aid of OpenMP
213 (Open Multi-Processing) parallel acceleration. Moreover, by constructing memory-map file
214 for genotype in disk rather than load it all into RAM, rMVP significantly decrease the
215 memory cost, making rMVP pretty promising in process of big data at a PC with limited
216 computing resources.



217

218 **Figure 1. Comparison of memory usage in response to number of threads used for parallel computation**
219 **under "speed" mode of rMVP.**

220 For each block with a specific colour, the y-axis represents memory usage (Mb) in log2 scale; the x-axis
221 represents computing time. Different colour represents different number of threads used for parallel computation,
222 the height of area in dark colour represents real memory costs while the height of shadow in light colour
223 represents theoretical memory costs which is 1, 4, 16, and 64 times of genotype data size in standard R matrix

224 format under ‘fork’ parallel mode, respectively. Data for speed test was generated by PLINK software and each
 225 data unit represents 1,000 samples and 100,000 SNPs. The data size for testing memory usage was 16 data units
 226 (16,000 samples and 1,600,000 SNPs), 10 PCs are added as covariates in all test methods. All tests were
 227 performed on a Red Hat Enterprise Linux sever with 2.60 GHz Intel(R) Xeon(R) 32CPUs E5-4620 v2, and 512
 228 GB memory.

229 For MLM in Figure 1, a high shoulder peak appears at the beginning of the memory
 230 records, which indicating that the most memory cost part of the MLM is the construction of
 231 GRM. From the computation details of VanRaden method described above (Equation 4), we
 232 can conclude that the calculation of $Z^T Z$ requires gigantic storage space and the requirement
 233 is increasing with both the marker size and the number of individuals. To take care of this
 234 problem, we implement two modes (“speed” and “memory”) in rMVP to handle the big data
 235 with limited computation resources.

236 For the “speed” mode, the genotype matrix is stored in R standard matrix format and
 237 the transpose of Z matrix and the matrix multiplication are carried out by the RcppArmadillo
 238 package, which could be automatically speeded up by the Inter MKL math library based on
 239 Microsoft R Open platform. However, the big genotype data is loaded into RAM and
 240 resulting in a big memory cost as most of the GWAS software tools do, e.g., GEMMA,
 241 GCTA, and GAPIT. For the “memory” mode, all the matrices that required for constructing
 242 the GRM are stored in the ‘big.matrix’ format and the matrix multiplication of ‘big.matrix’ is
 243 implemented by our newly developed C++ function, which could be parallel accelerated by
 244 using the OpenMP (Open Multi-Processing) technology. Although it can significantly
 245 decrease the cost of memory, a little bit more computing time is required (Table 1). Users can
 246 easily adjust the “priority” parameter to get rid of the memory limit or the fastest speed
 247 depending on the data size and computing resources.

248 **Table 1. Comparison of memory and time cost under modes of “speed” and “memory”**

Mem (Gb) /Time (min)	Data Units				
	1	2	4	8	16
Speed mode	0.51/0.05	3.28/0.15	17.80/0.6	73.10/3.2	285.60/34.70
Memory mode	0.06/0.20	0.08/1.61	0.17/9	0.53/42.12	2.06/461.66

249 *Note:* Data for speed test was generated by PLINK software and each data unit represents 1,000 samples and
 250 100,000 SNPs. Parallel computation with 32 CPUs is used to speed up for both two modes. All tests were
 251 performed on a Red Hat Enterprise Linux sever with 2.60 GHz Intel(R) Xeon(R) 32CPUs E5-4620 v2, and 512
 252 GB memory.

253 **Parallel-accelerated: Parallel computation and block matrix multiplication for**
 254 **accelerating association tests**

255 *Speed Up by Block matrix multiplication.* Most GWAS models contain several columns of
 256 covariates, such as PCs and Sex, and the linear model function has to be solved for every

257 single tested marker. This process involves the inverse of the design matrix for covariates and
258 the tested markers. Since the covariates were the same for every tested marker, we partitioned
259 the design matrix into sub-matrices according to the covariates and the testing markers. The
260 inverse of the entire design matrix was calculated from the one-time calculation of the inverse
261 of the sub-matrix of covariates. As the number of covariates and markers increased, sub-
262 matrices partitioning significantly saved computing time (Table 2). Block matrix
263 multiplication strategy has been used in all association tests including GLM, MLM, and
264 FarmCPU.

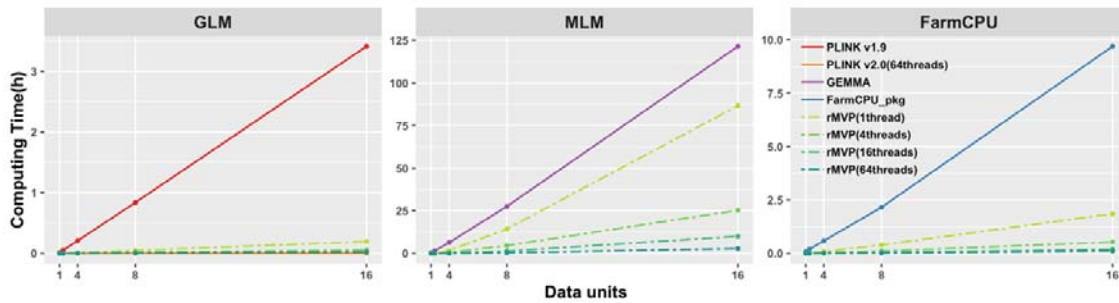
265 **Table 2. Speed performance of general linear model with and without using block matrix multiplication**
266 **strategy**

Time (Seconds)	Methods	
	Without block matrix multiplication strategy	With block matrix multiplication strategy
0 covariates	1,012	597
3 covariates	2,853	614
5 covariates	4,908	623
10 covariates	10,837	681

267 *Note:* 0, 3, 5, and 10 covariates are added in both Plink v1.9 and rMVP for testing speed of general linear model
268 without and with block matrix multiplication strategy, respectively. The advantage of block matrix
269 multiplication is increasing when more covariates added as fixed effects. A dataset includes 16,000 samples
270 with 1,600,000 SNPs was generated by PLINK software and used for test. All tests are performed using single
271 thread.

272 ***Speed Up by Parallel computation.*** There are two levels of parallel computation
273 implemented in rMVP: Data level parallel (DLP) and Thread level parallel (TLP). For DLP,
274 based on (1) Based on the Microsoft R Open platform, multi-threads have been automatically
275 assigned to speed up the mathematical calculation, such as matrix manipulation. For DLP,
276 association tests on millions of markers are allocated to a group of threads and calculated
277 simultaneously. rMVP switches between the two levels of parallel computation to achieve the
278 highest speed based on the biggest computation requirements in different GWAS procedures.
279 Since three association test methods (GLM, MLM, and FarmCPU) nearly generated the
280 consistent results (Figure S1) and same Power/FDR performance (Figure S2) as related
281 methods in PLINK v1.9 (written in C++, <https://www.cog-genomics.org/plink/>) and multiple
282 threads version PLINK v2.0 (written in C++, <https://www.cog-genomics.org/plink/2.0/>),
283 GEMMA (written in C++, <https://github.com/genetics-statistics/GEMMA/>), and
284 FarmCPU_pkg (R package written in pure R, <http://zzlab.net/FarmCPU/>), respectively. The
285 rMVP (written in R and C++) was compared with these software packages for speed
286 performance and the computing time was recorded for each software from loading data to

287 generating results files (Figure 2, Table S1). Detailed software version and scripts for
288 computing speed test are provided in Supplementary Table S2.



289

290 **Figure 2. Comparison of computing speed of PLINK, GEMMA, FarmCPU_pkg, and rMVP (“Speed”**
291 **mode)**

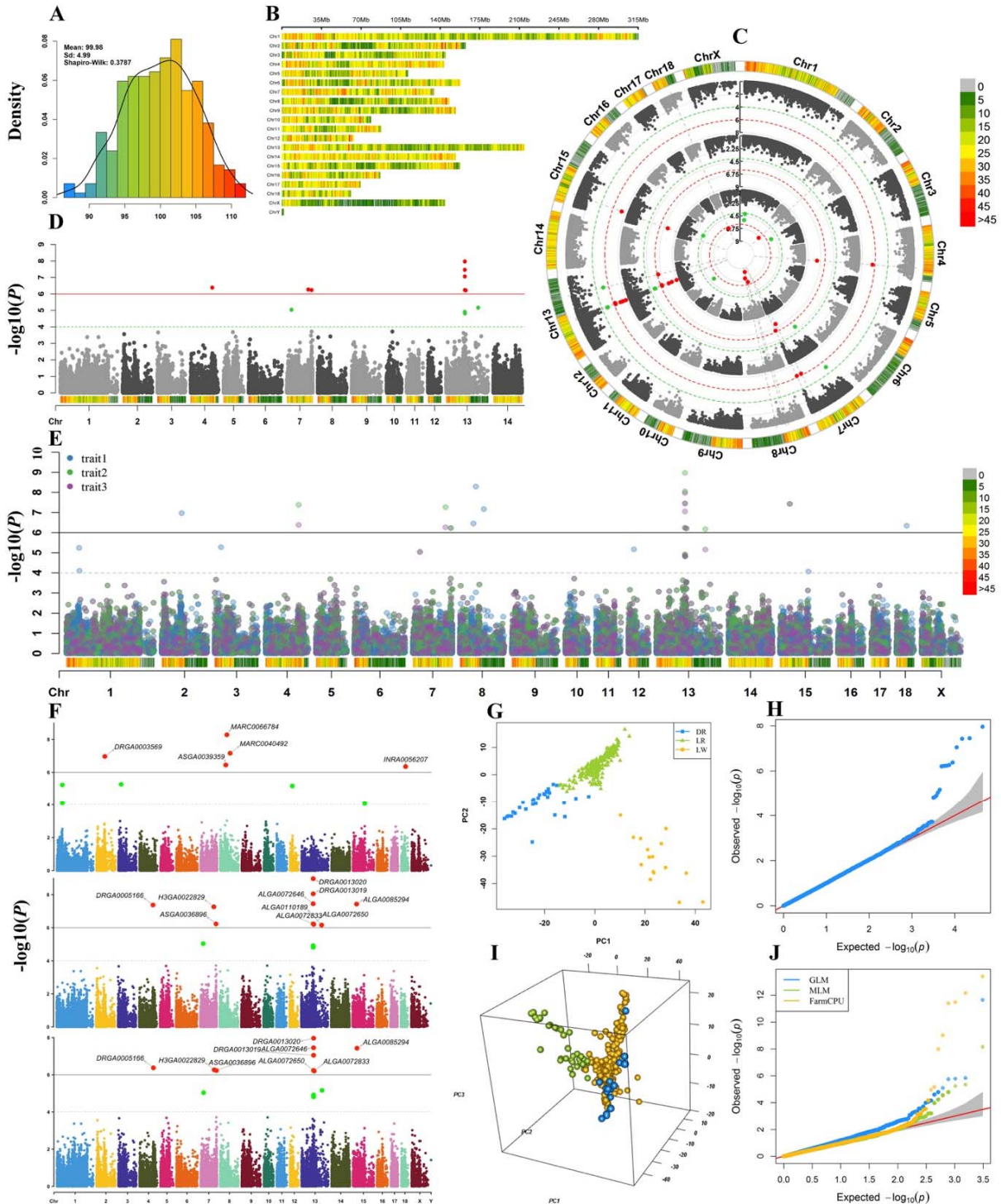
292 Computing time (hours) in response to data units are displayed, 5 PCs are added as covariates in all test methods.
293 Speed performances of association test methods in rMVP were performed using 1, 4, 16, 64 threads and
294 compared with PLINK, GEMMA, and FarmCPU_pkg, respectively. Data for speed test was generated by
295 PLINK software and each data unit represents 1,000 samples and 100,000 SNPs. The biggest data for memory
296 test of all models was 16 data units (16,000 samples and 1,600,000 SNPs). All tests were performed on a Red
297 Hat Enterprise Linux sever with 2.60 GHz Intel(R) Xeon(R) 32CPUs E5-4620 v2, and 512 GB memory.

298 **Visualization enhanced: flexible adjustments for generating high-quality figures**

299 ‘MVP.report’ function provides a pack of high-quality figures for visualizing GWAS related
300 information, including data information, population structure, and GWAS results.

301 Visualization of data information includes phenotype distribution (Figure 3A) and
302 marker density (Figure 3B), which are used to show if the phenotype is normally distributed
303 and the SNPs are evenly covered the entire genome. Skewed phenotype distribution and
304 uneven distributed genotype data would result false positives and biased estimation of
305 population structure and relationship among individuals. Top PCs are visualized in manner of
306 both two and three dimensions to display the population structure (Figure 3G and I).

307 Visualization of GWAS results includes Manhattan plot and Q-Q plot. Marker density
308 information is added to Manhattan plot to show the marker coverage of candidate region
309 (Figure 3D). Multiple-group GWAS results can be visualized on a same Manhattan plot and
310 Q-Q plot for easier comparison, common detected signals can be marked with dotted lines,
311 and users could highlight some SNPs or genes of interest on the Manhattan plot without
312 overlap (Figure 3C, E, F, H, and J). Our ‘MVP.report’ can also easily process GWAS results
313 from other software for visualization, such as PLINK, GEMMA, GCTA, and TASSEL. This
314 function can be further extended to visualizing the results from analyses of multi-omics,
315 correlated traits, and eQTL, and to displaying the commonly detected candidate areas. Users
316 can make a desired output figures using more than 40 parameters, detailed description for all
317 parameters are listed in Supplementary Table S3 and Supplementary File S1.



318
319 **Figure 3. Visualization of GWAS related information**

320 A. Phenotype distribution; B. Marker density, colour lumps with a user-defined window size (e.g. 1 Mb);
321 Manhattan plot for single-group GWAS results with marker density information (D); Manhattan plot for multiple-
322 group GWAS results in both circular manner and rectangular manner (C, E, F); Visualization of population
323 structure in both two dimensions (G) and three dimensions (I); Q-Q plot for single-group GWAS result (H); Q-Q
324 plot for multiple-group GWAS results (J).

325 Discussion

326 A summary of GWAS related functions of rMVP compared with other software tool is listed
 327 in Table 3. At the moment, rMVP does not provide functions of imputation and quality
 328 control, which need to be done before association tests. Instead, rMVP provides functions for
 329 flexible data conversion that can easily accept the data from other software, such as
 330 Beagle[28], which also accepts data in VCF format and provides imputation and quality
 331 control functions.

332 **Table 3.** Summary of GWAS related functions in Plink, GEMMA, FarmCPU_pkg, and rMVP

Functions	Items	Software			
		Plink	GEMMA	FarmCPU_pkg	rMVP
Input	Hapmap	×	×	√	√
	VCF	√	×	×	√
	Binary	√	√	×	√
	Numeric	×	×	√	√
	BIMBAM	×	√	×	×
	Quality control	√	×	×	×
Model	GLM	√	√	√	√
	MLM	×	√	×	√
	FarmCPU	×	×	√	√
Population structure	Principal components	√	×	√	√
	Genomic relationship matrix	×	√	×	√
Variance components estimation	BRENT	×	×	×	√
	EMMA	×	×	√	√
	Fast-LMM	×	×	√	√
	HE regression	×	√	×	√
Output	<i>p-values</i> , SE, Effect	√	√	√	√
	Manhattan plot	×	×	√	√
	QQ-plot	×	×	√	√
	SNP density plot	×	×	×	√
	Phenotype distribution	×	×	×	√
	PCA plot	×	×	√	√

333 rMVP currently only supports DLP and TLP for parallel computation, lacking the
 334 implementation of distributed parallel system (DPS). Compared with TLP that can speed up
 335 the computation using 100 threads on a single node, DPS (e.g. MPI, Hadoop, and Spark) can
 336 distribute the tasks to 1000 threads on multiple nodes. DPS is also better at dealing with
 337 hundreds or thousands of phenotypes and large computing tasks that need to be split, but its
 338 performance is limited by the efficiency of data transfer among multi-nodes through the local

339 network. However, association tests in rMVP can be accomplished within 10 hours for a
340 dataset that includes 500,000 samples and five million markers for each sample using
341 FarmCPU model, suggesting that our rMVP can meet most users' requirements.

342 Future work includes implementing efficient imputation and quality control functions,
343 and supporting DPS to meet the challenge of big datasets with millions of samples. We also
344 plan to incorporate more association test methods, such as logistic regression and multi-trait
345 model, which fits binary and multi-genetically-correlated traits. With the development of
346 GPU technology, we can get thousands of cores and higher memory bandwidth at a low price.
347 Most of the processes in the GWAS analysis have good independence and can give full play
348 to the advantages of GPU parallel computing. But the bottleneck of limited GPU memory
349 makes it difficult to perform GPU-based GWAS analysis on a large population. In future, we
350 plan to extend rMVP to support parallel computing on multiple machines with multiple GPUs
351 for each machine and explore new memory optimization methods. Incorporating the above
352 methods will greatly improve the versatility of rMVP.

353 **Code availability**

354 The rMVP package is available on both CRAN ([https://cran.r-](https://cran.r-project.org/web/packages/rMVP)
355 [project.org/web/packages/rMVP](https://cran.r-project.org/web/packages/rMVP)) and GitHub (<https://github.com/xiaolei-lab/rMVP>).

356 **Authors' contributions**

357 Xiaolei Liu, Xinyun Li, SZ designed the study, LY, HZ, and Xiaolei Liu wrote the software.
358 ZT, JX, and DY tested the stability of our package using various datasets. ZZ, MZ, XY gave
359 professional suggestions on the experimental design. Xiaolei Liu and LY drafted the
360 manuscript, and ALL authors contributed to finalizing the writing.

361 **Competing interests**

362 The authors have declared no competing interests.

363 **Acknowledgements**

364 We thank all rMVP beta version users for giving their valuable feedbacks through GitHub.
365 This work was supported by the National Natural Science Foundation of China [31672391,
366 31730089, 31702087, 31701144]; the National Key Research and Development Program
367 [2016YFD0101900]; the Fundamental Research Funds for the Central Universities
368 [2662020DKPY007, 2662019PY011]; the National Science Foundation [DBI 1661348]; and
369 the National Swine System Industry Technology System [CARS-35].

370

371 **CRCID**

372 0000-0003-4413-7976 (Yin L)
373 0000-0002-7913-5228 (Zhang H)
374 0000-0002-7263-5967 (Tang Z)
375 0000-0002-4951-098X (Xu J)
376 0000-0001-9762-4204 (Yin D)
377 0000-0002-5784-9684 (Zhang Z)
378 0000-0003-0661-5332 (Yuan X)
379 0000-0001-8931-5022 (Zhu M)
380 0000-0002-3997-2320 (Zhao S)
381 0000-0002-2314-4926 (Li X)
382 0000-0002-9954-1426 (Liu X)

383 **References**

384 [1] Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, et al. 10 Years of GWAS
385 Discovery: Biology, Function, and Translation. *Am J Hum Genet* 2017;101:5-22.
386 [2] Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and pitfalls in the
387 application of mixed-model association methods. *Nat Genet* 2014;46:100-6.
388 [3] Zhang Z, Buckler ES, Casstevens TM, Bradbury PJ. Software engineering the mixed model for
389 genome-wide association studies on large samples. *Brief Bioinform* 2009;10:664-75.
390 [4] Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus
391 genotype data. *Genetics* 2000;155:945-59.
392 [5] Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components
393 analysis corrects for stratification in genome-wide association studies. *Nat Genet* 2006;38:904-9.
394 [6] Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a tool set
395 for whole-genome association and population-based linkage analyses. *The American journal of*
396 *human genetics* 2007;81:559-75.
397 [7] Yu J, Pressoir G, Briggs WH, Vroh Bi I, Yamasaki M, Doebley JF, et al. A unified mixed-model
398 method for association mapping that accounts for multiple levels of relatedness. *Nat Genet*
399 2006;38:203-8.
400 [8] Zhang Z, Ersoz E, Lai C-Q, Todhunter RJ, Tiwari HK, Gore MA, et al. Mixed linear model
401 approach adapted for genome-wide association studies. *Nature genetics* 2010;42:355.
402 [9] Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST linear mixed models
403 for genome-wide association studies. *Nature methods* 2011;8:833-5.
404 [10] Segura V, Vilhjálmsson BJ, Platt A, Korte A, Seren Ü, Long Q, et al. An efficient multi-locus
405 mixed-model approach for genome-wide association studies in structured populations. *Nature genetics*
406 2012;44:825.
407 [11] Li M, Liu X, Bradbury P, Yu J, Zhang Y-M, Todhunter RJ, et al. Enrichment of statistical power
408 for genome-wide association studies. *BMC biology* 2014;12:73.
409 [12] Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. TASSEL: software
410 for association mapping of complex traits in diverse samples. *Bioinformatics* 2007;23:2633-5.
411 [13] Lipka AE, Tian F, Wang Q, Peiffer J, Li M, Bradbury PJ, et al. GAPIT: genome association and
412 prediction integrated tool. *Bioinformatics* 2012;28:2397-9.
413 [14] Tang Y, Liu X, Wang J, Li M, Wang Q, Tian F, et al. GAPIT version 2: an enhanced integrated
414 tool for genomic association and prediction. *The plant genome* 2016;9:1-9.
415 [15] Aulchenko YS, Ripke S, Isaacs A, Van Duijn CM. GenABEL: an R library for genome-wide
416 association analysis. *Bioinformatics* 2007;23:1294-6.

- 417 [16] Kang HM, Sul JH, Service SK, Zaitlen NA, Kong S-y, Freimer NB, et al. Variance component
418 model to account for sample structure in genome-wide association studies. *Nature genetics*
419 2010;42:348-54.
- 420 [17] Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association studies. *Nat*
421 *Genet* 2012;44:821-4.
- 422 [18] Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait
423 analysis. *The American Journal of Human Genetics* 2011;88:76-82.
- 424 [19] Liu X, Huang M, Fan B, Buckler ES, Zhang Z. Iterative Usage of Fixed and Random Effect
425 Models for Powerful and Efficient Genome-Wide Association Studies. *PLoS Genet*
426 2016;12:e1005767.
- 427 [20] Yoon S, Nguyen HCT, Yoo YJ, Kim J, Baik B, Kim S, et al. Efficient pathway enrichment and
428 network analysis of GWAS summary data using GSA-SNP2. *Nucleic acids research* 2018;46:e60-e.
- 429 [21] Korte A, Vilhjálmsson BJ, Segura V, Platt A, Long Q, Nordborg M. A mixed-model approach
430 for genome-wide association studies of correlated traits in structured populations. *Nature genetics*
431 2012;44:1066-71.
- 432 [22] Zhou X, Stephens M. Efficient multivariate linear mixed model algorithms for genome-wide
433 association studies. *Nature methods* 2014;11:407-9.
- 434 [23] Casale FP, Rakitsch B, Lippert C, Stegle O. Efficient set tests for the genetic analysis of
435 correlated traits. *Nature methods* 2015;12:755-8.
- 436 [24] Kane MJ, Emerson J, Weston S. Scalable strategies for computing with massive data. *Journal of*
437 *Statistical Software* 2013;55:1-19.
- 438 [25] VanRaden PM. Efficient methods to compute genomic predictions. *Journal of dairy science*
439 2008;91:4414-23.
- 440 [26] Burch BD, Iyer HK. Exact confidence intervals for a variance ratio (or heritability) in a mixed
441 linear model. *Biometrics* 1997:1318-33.
- 442 [27] Zhou X. A unified framework for variance component estimation with summary statistics in
443 genome-wide association studies. *The annals of applied statistics* 2017;11:2027.
- 444 [28] Browning BL, Zhou Y, Browning SR. A one-penny imputed genome from next-generation
445 reference panels. *The American Journal of Human Genetics* 2018;103:338-48.
- 446
- 447

448 **Figure legends**

449 **Table S1.** Computation speed performances of PLINK, GEMMA, and rMVP for five
450 simulated datasets.

451 **Table S2.** Software versions and codes used in performance tests.

452 **Table S3.** Parameter details for flexible visualization of GWAS related information.

453 **Figure S1**

454 Title: Comparisons of association results between rMVP and related software.

455 Legend: x-axis is the computed *p-value* in -log10 format for different GWAS models, y-axis is the computed
456 *p-value* in -log10 format of related software for corresponding GWAS model, the experiment was performed on
457 the simulated 16 data units (16,000 samples and 1,600,000 SNPs).

458 **Figure S2**

459 Title: Comparisons of power and false positive discovery for different GWAS models
460 between rMVP and related software.

461 Legend: The experiment was performed using an Arabidopsis dataset, which includes 1178 individuals and
462 208794 SNPs, the phenotype was simulated by randomly selected 10 QTNs following a normal distribution with
463 mean 0 and variance 0.1, the heritability was 0.5. The final results were the average of 100 replicates.

464 **Figure S3**

465 Title: The road mapping of whole GWAS procedures in rMVP.

466 Legend: K is the Kinship matrix, also known as Genomic relationship matrix (GRM). EigenK represents the
467 eigen decomposition of GRM. PC represents principal components. VC represents variance components.

468 **File S1.** Demo scripts and figures for visualization in rMVP

469

470 **Supplementary material**

471 **Table S1**

472 **Table S2**

473 **Table S3**

474 **Figure S1**

475 **Figure S2**

476 **Figure S3**

477 **File S1**

478