


METHOD

Open Access

# MegaLMM: Mega-scale linear mixed models for genomic predictions with thousands of traits



Daniel E. Runcie<sup>1\*</sup> , Jiayi Qu<sup>1</sup>, Hao Cheng<sup>1</sup> and Lorin Crawford<sup>2</sup>

\*Correspondence:

[deruncie@ucdavis.edu](mailto:deruncie@ucdavis.edu)

<sup>1</sup>Department of Plant Sciences,  
University of California Davis, Davis,  
CA, USA

Full list of author information is  
available at the end of the article

## Abstract

Large-scale phenotype data can enhance the power of genomic prediction in plant and animal breeding, as well as human genetics. However, the statistical foundation of multi-trait genomic prediction is based on the multivariate linear mixed effect model, a tool notorious for its fragility when applied to more than a handful of traits. We present **MegaLMM**, a statistical framework and associated software package for mixed model analyses of a virtually unlimited number of traits. Using three examples with real plant data, we show that **MegaLMM** can leverage thousands of traits at once to significantly improve genetic value prediction accuracy.

**Keywords:** Multi-trait Linear Mixed Model, Genomic prediction, High-throughput phenotyping, Multi-environment trial

## Background

New high-throughput phenotyping technologies hold promise for a revolution in data-driven decisions in plant and animal breeding programs [1, 2]. For example, drone-based hyperspectral cameras can image fields at high resolution across hundreds of spectral bands [3], wearable sensors can continuously monitor animals health and physiology [4], and RNA sequencing and metabolite profiling can simultaneously assay the concentrations of tens-of-thousands of targets [5]. These high-dimensional traits could allow breeders to rapidly assess many aspects of performance more accurately or earlier in development than was possible using traditional tools. This can increase the rate of gain in target traits by increasing selection accuracy, increasing selection intensity, and reducing breeding cycle durations.

However, efficiently incorporating high-dimensional phenotype data into breeding decisions is challenging. Whenever two traits are genetically correlated, joint analyses can improve the precision of variety evaluation [6]. However, two key problems emerge. First, the number of traits in high-dimensional datasets is often much larger than the number of breeding lines, which means that naive correlation estimates are not robust. Second,



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

phenotypic correlation among traits are often poor approximations to genetic correlation, so not all correlated traits are useful for breeding decisions [7]. For example, plants grown in more productive areas of a field will tend to produce higher yields and be greener (measured by hyperspectral reflectance). Yet, selecting indirectly based on green plants instead of directly on higher yields may be counter-productive because “green-ess” may indicate an over-investment in vegetative tissues at the expense of seed. This contrasts with the problem of predicting genetic values from genotype data (e.g., genomic prediction; [8]), where all correlations between candidate features and performance are useful for selection.

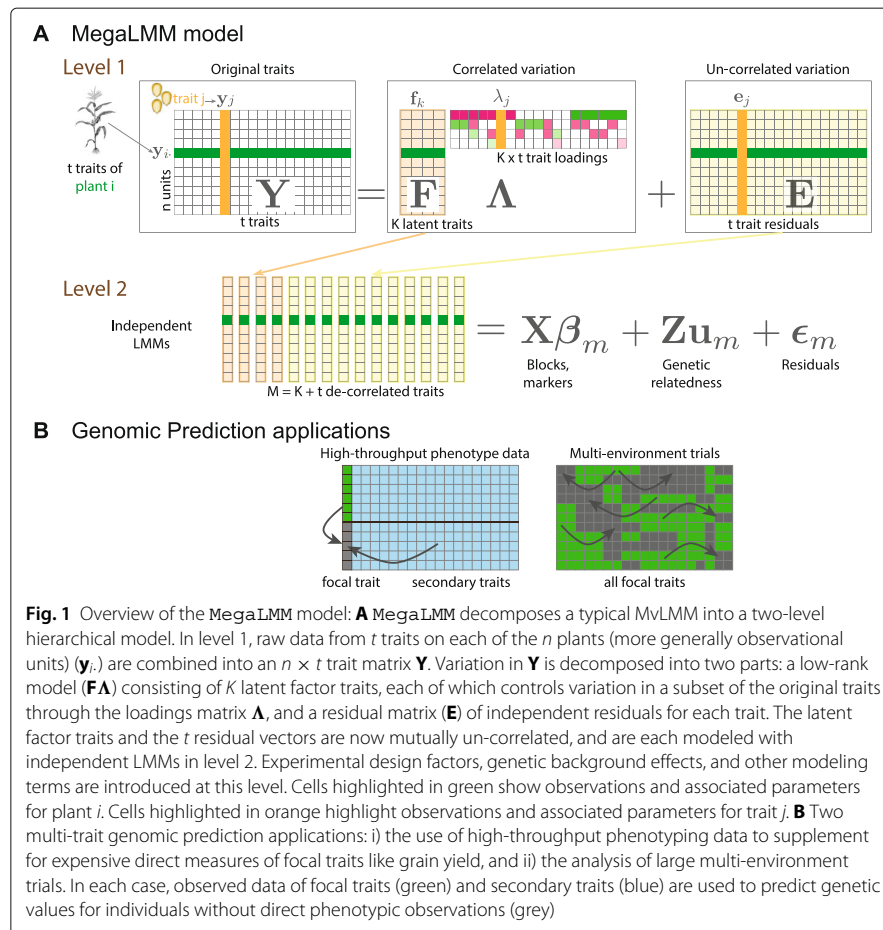
The multivariate linear mixed model (MvLMM) is a widely-used statistical tool for decomposing phenotypic correlations into genetic and non-genetic components. The MvLMM is a multi-outcome generalization of the univariate linear mixed model (LMM) that forms the backbone of the majority of methods in quantitative genetics. The MvLMM was introduced over 40 years ago [9], and has repeatedly been shown to increase selection efficiency [10–12]. Yet, MvLMMs are still rarely used in actual breeding programs because naive implementations of the framework are sensitive to noise, prone to overfitting, and exhibit convergence problems [13]. Furthermore, existing algorithms are extremely computationally demanding. The fragility of naive MvLMMs is due to the number of variance-covariance parameters that must be estimated which increases quadratically with the number of traits. The computational demands increase even more dramatically: from cubically to quintically with the number of traits [14] because most algorithms require repeated inversion of large covariance matrices. These matrix operations dominate the time required to fit a MvLMMs, leading to models that take days, weeks, or even years to converge.

Here, we describe MegaLMM (linear mixed models for millions of observations), a novel statistical method and computational algorithm for fitting massive-scale MvLMMs to large-scale phenotypic datasets. Although we focus on plant breeding applications for concreteness, our method can be broadly applied wherever multi-trait linear mixed models are used (e.g., human genetics, industrial experiments, psychology, linguistics, etc.). MegaLMM dramatically improves upon existing methods that fit low-rank MvLMMs, allowing multiple random effects and un-balanced study designs with large amounts of missing data. We achieve both scalability and statistical robustness by combining strong, but biologically motivated, Bayesian priors for statistical regularization—analogueous to the  $p \gg n$  approach of genomic prediction methods—with algorithmic innovations recently developed for LMMs. In the three examples below, we demonstrate that our algorithm maintains high predictive accuracy for tens-of-thousands of traits, and dramatically improves the prediction of genetic values over existing methods when applied to data from real breeding programs.

## Results

### Methods overview

MegaLMM fits a full multi-trait linear mixed model (MvLMM) to a matrix of phenotypic observations for  $n$  genotypes and  $t$  traits (level 1 of Fig. 1A). We decompose this matrix into fixed, random, and residual components, while modeling the sources of variation and covariation among all pairs of traits. The main statistical and computational challenge of fitting large MvLMMs centers around the need to robustly estimate  $t \times t$



covariance matrices for the residuals and each random effect. Each covariance matrix has  $t(t-1)/2 + t$  free parameters, and any direct estimation approach is computationally demanding because it requires repeatedly inverting these matrices (an  $O(t^3)$  operation).

We solve both of these problems by introducing  $K$  un-observed (latent) traits called factors ( $\mathbf{f}_k$ ) to represent the causes of covariance among the  $t$  observed traits. We treat each latent trait just as we would any directly measured trait and decompose its variation into the same fixed, random and residual components using a set of parallel univariate linear mixed models (level 2 of Fig. 1A). We then model the pairwise correlations between each latent trait and each observed trait through  $K$  loadings vectors  $\lambda_k$ .

Together, the set of parallel univariate LMMs and the set of factor loading vectors result in a novel and very general re-parameterization of the MvLMM framework as a mixed-effect factor model. This parameterization leads to dramatic computational performance gains by avoiding all large matrix inversions. It also serves as a scaffold for eliciting Bayesian priors that are intuitive and provide powerful regularization which is necessary for robust performance with limited data. Our default prior distributions encourage: i) shrinkage on the factor-trait correlations ( $\lambda_{jk}$ ) to avoid over-fitting covariances, and ii) shrinkage on the factor sizes to avoid including too many latent traits. This two-

dimensional regularization helps the model focus only on the strongest, most relevant signals in the data.

While others have used latent factor approaches to reduce dimensionality of MvLMMs (e.g., [15–18]), these methods only use factors for a single random effect (usually the matrix of random genetic values)—with the exception of BSFG which uses factors for the combined effect of a single random effect and the residuals [17]. In MegaLMM, we expand this framework and use factors to model the joint effects of all predictors: fixed, random and residual factors on multiple traits.

We combine this efficient factor model structure with algorithmic innovations that greatly enhance computational efficiency, drawing upon recent work in LMMs [19–22]. While Gibbs samplers for MvLMMs are notoriously slow, we discovered extensive opportunities for collapsing sampling steps, marginalizing over missing data, and discretizing variance components so that intermediate results can be cached (Additional file 1: Supplemental Methods).

Genomic prediction using MegaLMM works by fitting the model to a partially observed trait matrix, with the traits to be predicted imputed as missing data. MegaLMM then estimates genetic values for all traits (both observed and missing) in a single step (Fig. 1B).

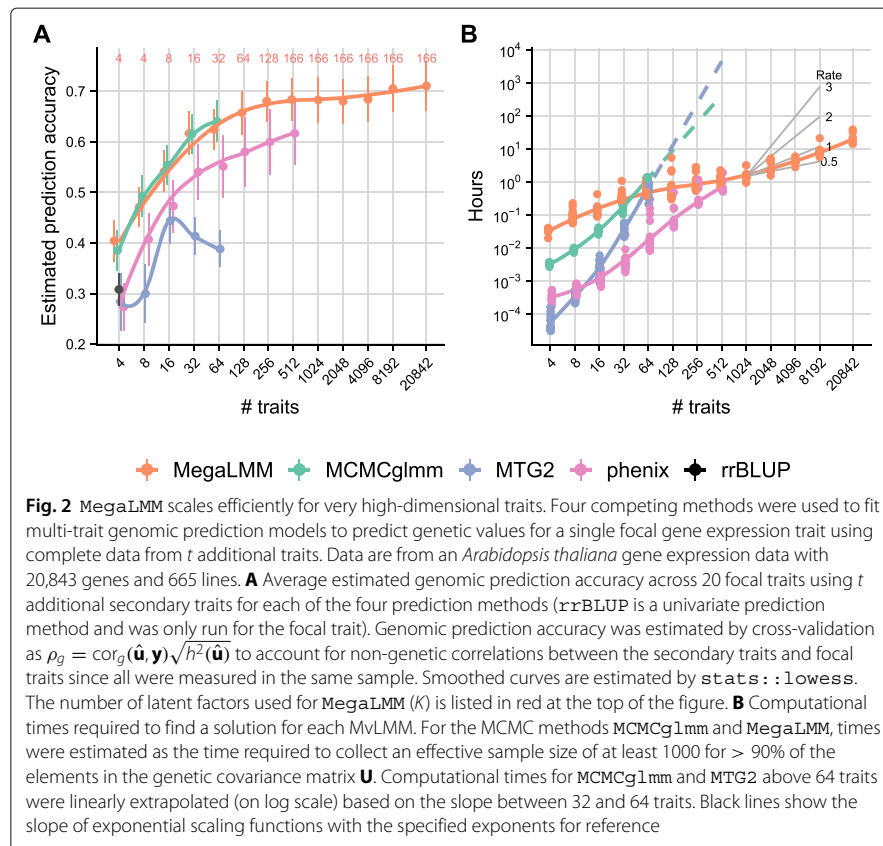
#### MegaLMM is efficient and effective for large datasets

We used a gene expression matrix with 20,843 genes measured in each of 665 *Arabidopsis thaliana* accessions (a total of nearly 14 million observations), to evaluate the accuracy and time requirements for trait-assisted genomic prediction—a classic example of an applied use of MvLMMs—across a panel of existing software packages. We created datasets with 4 to 20,842 “secondary” traits with complete data, and used these data to predict the genetic values of a single randomly selected “focal” gene with 50% missing data.

Despite the limited number of independent lines in this data set, adding up to  $\approx 200$  secondary traits improved the genomic prediction accuracy of MegaLMM and two other Bayesian methods: MCMCglmm and phenix (Fig. 2A). The maximum likelihood method MTG2 [23], on the other hand, did only marginally better than single-trait prediction, and genomic prediction accuracy declined with 32 traits, likely due to overfitting. We note that the results here are averages over 20 randomly selected focal genes. The prediction accuracy and benefits of multi-trait prediction varied considerably among genes (Additional file 1: Figure S1 and Figure S2), but comparisons among methods were largely correlated. Using simulated datasets where we knew the true genetic and residual covariances among traits, we also found that MegaLMM was at least as accurate in estimating covariance parameters as the competing methods (Additional file 1: Figure S3).

Beyond 32 secondary traits, computational times for MCMCglmm and MTG2 became prohibitive (Fig. 2B). Using extrapolation, we estimated that fitting these methods for 512 traits would take 20 days and 217 days, respectively, without considering issues of model convergence. In contrast, phenix and MegaLMM were both able to converge on good model fits for 512 traits in approximately one hour.

Beyond 512 traits, MegaLMM was the only viable method as phenix cannot be applied to datasets with  $t > n$  phenotypes. Although the genomic prediction accuracy of MegaLMM did not increase further after  $\approx 256$  traits, performance did not suffer even with the full dataset of  $> 20,000$  traits and the analysis was completed in less than a day.



This shows that MegaLMM is feasible to apply to very high-dimensional studies and, in most cases, does not require pre-filtering of traits—something that requires great care in genomic prediction applications to avoid misleading results [24].

An important feature of MegaLMM is that the choice of the number of latent factors  $K$  is less critical than in most factor models. Since factors are ordered from most-to-least important by the prior (See Methods), as long as enough factors are specified to capture the majority of the covariance among traits, adding additional latent factors does not lead to over-fitting (Additional file 1: Figure S4A). Additional factors do increase the run-time of the algorithm, though (Additional file 1: Figure S4B), so some optimization of  $K$  during the burn-in period can reduce computational demands during posterior sampling.

## Applications to real breeding programs

To demonstrate the utility of MegaLMM, we developed two classes of genomic prediction models for high-dimensional phenotype data in real plant breeding programs.

### Genomic prediction using hyperspectral reflectance data

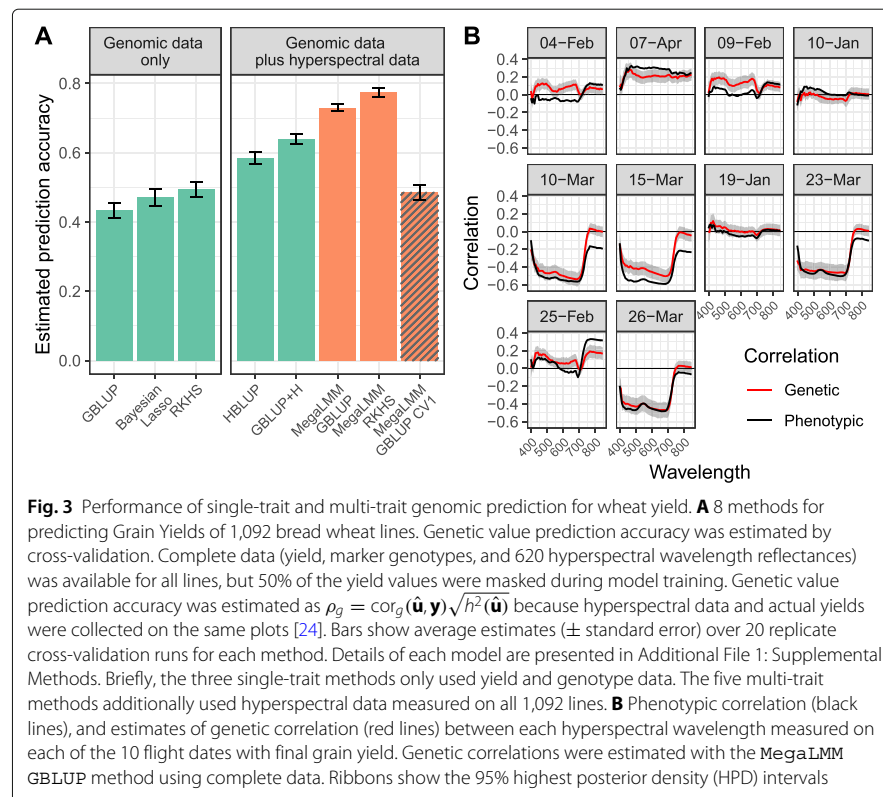
When the final performance of a variety is difficult or costly to obtain, breeding programs can supplement direct measures of performance with data from surrogate traits that can be measured cheaply, earlier in the breeding cycle, and on more varieties. For example, in the bread wheat breeding program at CIMMYT, hyperspectral reflectance data can be

collected rapidly and repeatedly by aerial drones on thousands of plots [25]. We developed a multi-trait genomic prediction model to incorporate 62-band hyperspectral reflectance data from 10 different drone flights over the course of one growing season, and compared the accuracy of these genetic value predictions against more traditional approaches.

We first compared three standard univariate methods: GBLUP [26], Bayesian LASSO (BL) [27], and Reproducing kernel Hilbert space (RKHS) regression [28]. GBLUP achieved a prediction accuracy of  $\rho_g = 0.43$  for yield (Fig. 3A). Both the BL and RKHS methods showed modest improvements, with  $\rho_g = 0.47$  and  $\rho_g = 0.49$ , respectively in these data. The RKHS model often out-performs GBLUP in plant breeding datasets, but improvements are generally slight and inconsistent depending on the genetic architecture of the targeted trait.

In the original analysis of this dataset, [25] achieved increased performance by replacing the genomic kernel (**K** in our notation) with a kernel based on the cross-product of hyperspectral reflectances across all wavelengths and time points (termed the **H** matrix). We replicated these results, achieving a prediction accuracy of  $\rho_g = 0.58$  (HBLUP method). These authors also proposed a multi-kernel model combining the **K** and **H** kernel matrices, although they only applied this to cross-treatment genotype-by-environment predictions. We found that applying this multi-kernel method to the within-environment data resulted in additional accuracy gains ( $\rho_g = 0.64$ ) (GBLUP+H method; Fig. 3A).

While more effective than univariate methods, predictions based on the **H** kernel matrix are biased by non-genetic correlations between surrogate traits and yield because



they do not directly model the genetic component of these correlations. MegaLMM implements a full multi-trait mixed model and thus can separate these sources of correlation. We fit three different multi-trait prediction models with MegaLMM. The first was a standard multi-trait mixed model with a single random effect based on the genomic relationship matrix  $\mathbf{K}$ . This method achieved a dramatically higher prediction accuracy than any of the previous approaches ( $\rho_g = 0.73$ ). Second, because the RKHS model had the highest performance among univariate predictions, we implemented an approximate RKHS method in MegaLMM based on averaging over three kernel matrices [28]. We are not aware of any other high-dimensional MvLMM implementations that allow models with multiple random effects. This model achieved the highest predictive accuracy ( $\rho_g = 0.77$ ). Finally, we repeated the MegaLMM-GBLUP analysis but this time masking all phenotype data (both grain yield and hyperspectral data) from the testing set. We called this approach MegaLMM-GBLUP-CV1 following the nomenclature from [29]. Genetic prediction accuracy in the CV1 setting was similar to the univariate methods ( $\rho_g = 0.49$ ), showing that nearly all benefit of MegaLMM in this dataset came through the optimal use of secondary trait phenotypes on the lines of interest.

In summary, by directly modeling the genetic covariance between the surrogate traits (hyperspectral reflectance measures), we achieved performance increases of 56%–79%, and up to 36% over the HBLUP method. To show that these conclusions were robust in other datasets, we repeated the same analyses in the other 19 trials reported by [25] and results were highly similar in all trials (Additional file 1: Figure S5).

To explore *why* directly modeling the genetic correlation is important, we compared the estimated genetic correlations between each hyperspectral band and grain yield to the corresponding phenotypic correlations (Fig. 3B). Most genetic correlation estimates closely paralleled the phenotypic correlations, with the largest values for low-to-intermediate wavelengths occurring on dates towards the end of the growing season while plants were in the grain filling stage [25]. However, there were notable differences. For example, genomic correlations were moderate ( $\rho_g \approx 0.2$ ) for most wavelengths during early February sampling dates while phenotypic correlations were near zero; yet, during early March time points, phenotypic correlations between yield and bands around 800 nanometers were moderate ( $\rho_y \approx -0.2$ ) but genomic correlations were near-zero. MegaLMM is able to model the discrepancy between genomic and phenotypic correlations, but methods based on the  $\mathbf{H}$  matrix (e.g., HBLUP) are not.

### Genomic prediction of agronomic traits across multi-environment trials

Multi-trait mixed models are also used to analyze data from multi-environment trials to account for genotype-environment interactions and select the best genotypes in each environment. The Genomes2Field initiative (<https://www.genomes2fields.org/>) is an ongoing multi-environment field experiment of maize hybrid genotypes across 20 American states and Canadian provinces. Data from the years 2014–2017 included 119 trials with a total of 2102 hybrids. As in many large-scale multi-environment trials, only a small proportion of the available genotypes were grown in each trial. Therefore, the majority of trial-genotype combinations were un-observed.

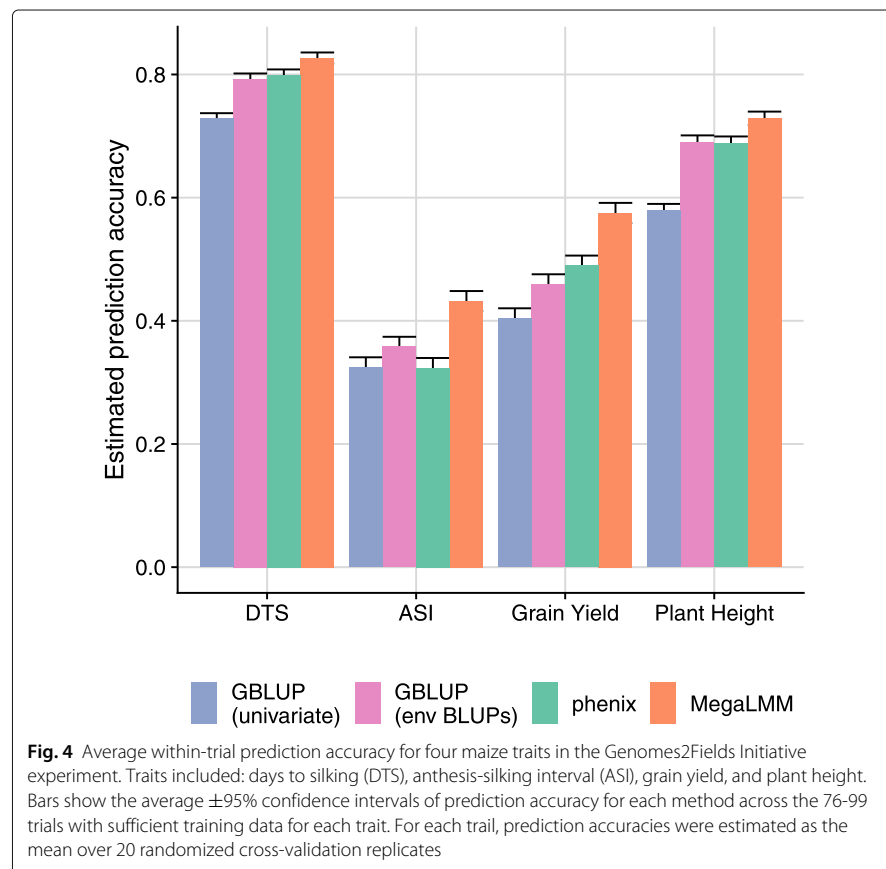
We selected four representative agronomically important traits and compared the ability of four modeling approaches to impute the missing measurements. Including across-trial information was beneficial for each of the four traits, suggesting generally



positive genetic correlations across trials. However, applying MegaLMM to each of the four trait datasets improved predictions dramatically, with average benefits across trials ranging from  $\rho_y = 0.10$  to  $\rho_y = 0.17$  (Fig. 4). The performance of phenix was inconsistent across traits and trials, likely because its model for the non-additive genetic covariance (i.e., the residual) is less flexible than MegaLMM.

To explore *why* jointly modeling all genetic and non-genetic covariances for each pair of trials improved prediction accuracy for each trait, we assessed the per-trial differences in performance between MegaLMM and the corresponding within-trial genomic prediction model. Trials varied considerably in how much MegaLMM improved genomic prediction accuracy, with several trials seeing improvements of  $\rho > 0.4$ . The magnitude of the MegaLMM effect on genomic prediction accuracy was largely explained by the maximum genetic covariance between that trial and any other trial in the dataset (Additional file 1: Figure S6). This is expected because the benefit of a MvLMM is largely dependent on the magnitude of genetic covariances between traits.

A common approach in multi-environment trials is to combine similar trials (based on geographic location or similar environments) into clusters and make genetic value predictions separately for each cluster [30]. However, this will not be successful if clusters cannot be selected a priori because using the trial data itself to identify clusters can lead to overfitting if not performed carefully [24]. In these data, the distribution of genomic correlations between trials differed among traits, so it is not straightforward to identify which

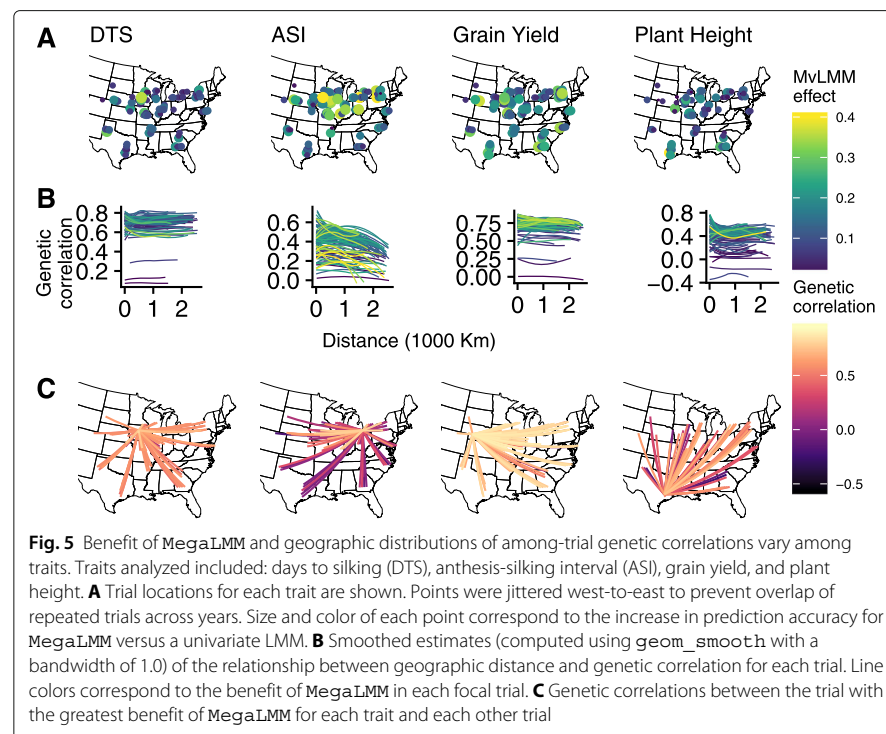




pairs or subsets of trials could be combined. The most obvious predictor of trial similarity is geographic distance, but we did not see consistent spatial patterns in the among-trial covariances across the four traits. The trials with the greatest benefit from our MvLMM showed geographic clustering in the central mid-west for the anthesis-silking interval (ASI) but not for the other three traits (Fig. 5A). Genetic correlations tended to decrease over long distances for ASI and over short distances for plant height, but not for the other two traits (Fig. 5B), resulting in obvious geographic clustering of genetic correlations for ASI but not the other traits (Fig. 5C). This suggests that including all trials together in one model is necessary to maximize the benefit of the MvLMM approach to multi-environment plant breeding.

## Discussion

Novel statistical methods can help optimize plant and animal breeding programs to meet future food security needs. In the above examples, we highlighted two areas where large-scale phenotype data can improve the accuracy of genomic prediction in realistic plant breeding scenarios: by incorporating high-throughput phenotyping data from remote sensors, and by synthesizing data on gene-environment interactions across large-scale multi-environment trials. In both examples, we apply high-dimensional multivariate linear mixed models to efficiently integrate all available genotype and phenotype data into genetic value predictions. MegaLMM is a scalable tool that extends the feasible range of input data for multivariate linear mixed models by at least two orders of magnitude over existing methods, while providing the flexibility to plug directly into existing breeding programs.



### Computational and statistical efficiency

Computational issues in single-trait LMMs have been studied extensively, allowing implementations for large datasets [14, 21, 22, 31]. Most of these algorithms diagonalize the genomic relationship matrices to improve computational efficiency. This technique dramatically improves the performance of simple, low-dimensional MvLMMs as well (e.g., [14, 23]). However, diagonalization does not address the computational challenge imposed by large trait-covariance matrices, and can only be applied to models with a single random effect and no missing data. Therefore, these tools cannot be applied to the datasets studied here or, more generally, to most large-scale studies of gene-environment interactions that frequently have large proportions of missing data [10] (Fig. 1) and to studies that have experimental designs with multiple sources of covariance (e.g., spatial environmental variation or non-additive genetics).

Our work builds on the factor-analytic approach to regularizing MvLMMs [15–18] and is most similar to BSFG [17] and *phenix* [18], which improve upon traditional quantitative genetic factor models by specifying sparse or low-rank factor matrices to improve robustness in high dimensions. Importantly, however, these models are limited to a single random effect and are not tractable for datasets with large numbers of traits because of computational inefficiencies (BSFG), or a lack of strong regularization on the residual covariance matrix (*phenix*). *MegaLMM* generalizes both methods and dramatically improves their weaknesses, allowing analyses with >20,000 traits to be completed in less than one day. Since *MegaLMM* scales approximately linearly with the number of traits (Fig. 2), applying it to datasets with many more traits may be feasible. While we have designed many of our routines to take advantage of multi-core CPUs, graphical processing units may offer additional performance gains.

Two key advantages of *MegaLMM* are its flexibility and generality. We have designed the *MegaLMM* R package to be as general as possible so that it can be applied to a wide array of problems in quantitative genetics. *MegaLMM* tolerates unbalanced designs with incomplete observations (something that makes *MCMCglmm* and *MTG2* very slow), arbitrarily complex fixed effect specifications to model experimental blocks, covariates, or other sources of variation among samples (unlike *phenix*), and most importantly, multiple random effects (unlike *phenix*, *GEMMA*, or *MTG2*). Multiple random effect terms can be used to account for spatially correlated variation across fields, non-additive genetic variation that is not useful for breeding, or to more flexibly model non-linear genetic architectures as we demonstrated with the approximate RKHS regression approach in the wheat application (Fig. 3). To make multiple-random-effect models computationally efficient, we take our earlier work with LMMs [22] and extend the same discrete estimation procedure to MvLMMs where the impact on computational efficiency is exponentially greater. Other commonly used tools for fitting MvLMMs such as *ASREML* [32] allow more flexibility in the specification of multiple variance-component models with correlated random effects that are not currently possible in *MegaLMM*. However, these tools do not scale well beyond  $\approx 10$  traits, so are not feasible to apply directly to large-scale datasets in plant breeding.

### Applicability to high-throughput phenotypic data

Large-scale phenotype data collection is rapidly emerging as a standard tool in plant breeding and other fields that use quantitative genetics [1, 33, 34]. These deep pheno-

typing datasets can be used as high-dimensional features to predict genetic values in agronomically important traits and serve as substitutes for direct assays where these are more time-consuming or expensive to collect.

Breeding objectives differ from the goals of polygenic risk score predictions for human diseases because the prediction target is not the phenotype of an individual, but its genetic value [24]. Genetic values quantify the expected phenotype of a plant's offspring, and so exclude impacts of the plant's own microenvironment on its phenotype [7]. Therefore, accurate genetic value prediction requires models that can distinguish between genetic and non-genetic sources of covariation among traits.

The MvLMM is considered the gold-standard method for isolating genetic correlations from non-genetic correlations in genetic value prediction [10]. However, it has rarely been applied in breeding programs because of the computational challenges associated with estimating multiple large covariance matrices. With high-throughput phenotype (HTP) data, MvLMMs have only been applied directly to sets of  $\approx 2 - 5$  traits. Instead, several authors have used a prior round of feature selection or calculated summary statistics of the HTP to generate model inputs rather than using the raw high-dimensional data itself (e.g., [3, 12, 35–37]). Other authors have replaced the MvLMM with a direct regression on the HTP data, using techniques such as factorial regression [38], functional regression [39], kernel regression [25], and deep learning [40]. While straightforward to implement, this conditioning on the HTP traits creates a form of collider bias which can induce genotype-phenotype associations that do not actually exist and impede genetic value predictions [24]. Alternative methods including IBCF [41]) and regularized selection indexes [42] avoid computational complexities of the full MvLMMs, but do not make full use of the trait correlations in the data.

MegaLMM, on the other hand, fits a full MvLMM to an arbitrary number of HTP traits and should be more efficient at leveraging high-dimensional genetic correlations while accounting for non-genetic sources of covariance, particularly for datasets when HTP traits and focal performance traits are measured on the same plants. Non-genetic correlations will be less important on datasets where these sets of traits are measured on different plots. At least in the wheat breeding trial datasets we examined, the benefit of multi-trait modeling was much greater when traits were partially observed on each individual than when secondary traits were only observed in the training partition. This is expected theoretically and has been demonstrated previously in simulations [24], but the magnitude of the benefit was particularly dramatic here. This suggests that breeding programs should focus on developing HTP technologies that can measure secondary traits on the target individuals; HTP measurements on training individuals may be less useful for prediction applications. Unlike other methods, including too many traits, or including redundant traits that are highly correlated is unlikely to significantly impact prediction accuracy, reducing the need to carefully choose which traits to include and which to exclude a priori; MegaLMM allows users to simply include all traits they have at once.

#### **Applicability to multi-environment trial data**

The analysis of multi-environment trials provides a separate set of computational and statistical challenges for plant breeders. Multi-environment trials (METs) are necessary because gene-environment interactions (GEIs) often prevent the same variety from performing best in all locations where a crop is grown [10]. However, METs are expensive and

logistically difficult. Genomic predictions in METs could reduce the need to test every variety in every environment, allowing smaller individual trials [43].

GEIs can be modeled in two ways: (i) as changes in variety effects on the same trait across environments (i.e., variety-by-environment interactions), or (ii) as a set of genetically correlated traits, with each trait-environment combination considered as a different phenotype [10]. When formulated with linear mixed models and random genetic effects, these two approaches are mathematically equivalent. Traditionally, the most common model for analyzing METs has been the AMMI model in which the genetic effects of each variety in each environment are modeled using a set of products between genetic and environmental vectors [44]. AMMI models are used to rank genotypes in different environments and to identify environmental clusters with similar rankings of varieties. However, AMMI models cannot easily incorporate marker data. When genetic values are treated as random effects, AMMI models becomes factor models (generally called factor analytic models in this literature) (e.g. [45, 46]), and can incorporate genetic marker data (e.g. [47]). MegaLMM extends this factor-analytic method for analyzing METs, making the methods robust for METs with hundreds or more individual trials.

A limitation of the AMMI factor-analytic approach to analyzing METs is that there is no mechanism for extending predictions to new environments outside of those already tested. Even large-scale commercial trials cannot test every field a farmer might use. Several authors have proposed using environmental covariates (ECs) to model environmental similarity in METs and predict GEIs for novel environments (e.g., [47–49]). These approaches all involve regressions of genetic variation on the ECs, and so, if relevant ECs are missing or the relationship between variety plasticity and ECs is non-linear, these models will under-fit the GEIs. Nevertheless, these approaches are promising and have been successfully applied to large METs (e.g. [47]). MegaLMM cannot currently incorporate ECs to predict novel environments. However, a possible extension could involve replacing the *iid* prior on the elements of the factor loadings matrix with a regression on the ECs. This hybrid of ECs and a full MvLMM could leverage the strengths of both approaches.

### Model limitations

While MegaLMM works well across a wide range of applications in breeding programs, our approach does have some limitations.

First, since MegaLMM is built on the Grid-LMM framework for efficient likelihood calculations [22], it does not scale well to large numbers of observations (in contrast to large numbers of traits), or large numbers of random effects. As the number of observational units increases, MegaLMM's memory requirements increase quadratically because of the requirement to store sets of pre-calculated inverse-variance matrices. Similarly, for each additional random effect term included in the model, memory requirements increase exponentially. Therefore, we generally limit models to fewer than 10,000 observations and only 1-to-4 random effect terms per trait. There may be opportunities to reduce this memory burden if some of the random effects are low-rank; then these random effects could be updated *on the fly* using efficient routines for low-rank Cholesky updates. We also do not currently suggest including regressions directly on markers and have used marker-based kinship matrices here instead for computational efficiency. Therefore as a stand-alone prediction method, MegaLMM requires calculations involving the Schur com-

plement of the joint kinship matrix of the testing and training individuals which can be computationally costly.

Second, MegaLMM is inherently a linear model and cannot effectively model trait relationships that are non-linear. Some non-linear relationships between predictor variables (like genotypes) and traits can be modeled through non-linear kernel matrices, as we demonstrated with the RKHS application to the Bread Wheat data. However, allowing non-linear relationships among traits is currently beyond the capacity of our software and modeling approach. Extending our mixed effect model on the low-dimensional latent factor space to a non-linear modeling structure like a neural network may be an exciting area for future research. Also, some sets of traits may not have low-rank correlation structures that are well-approximated by a factor model. For example, certain auto-regressive dependence structures are low-rank but cannot efficiently be decomposed into a discrete set of factors.

Nevertheless, we believe that in its current form, MegaLMM will be useful to a wide range of researchers in quantitative genetics and plant breeding.

#### Potential extensions

Beyond the examples we show in this work, the scalability and statistical power of MegaLMM can open up new avenues for innovation in genomic prediction applications across the fields of quantitative genetics—both in breeding programs as we have described here and, potentially, in human genetics. Genomic prediction is also used for the calculation of polygenic risk scores for complex human traits and diseases [50]. MegaLMM may help leverage past case histories, survey responses, molecular tests, and the genetic architecture of other correlated traits to provide a more comprehensive multi-trait polygenic risk score (e.g. [51]).

We have focused here on simple scalar phenotypes: the expression of a single gene, the total grain yield, and individual measures of agronomic performance. However, many important traits in plants, animals, and humans cannot easily be reduced to a scalar value. Examples include time-series traits such as growth curves [52], metabolic traits such as the relative concentrations of different families of metabolites [53], and morphological traits such as shape or color [54]. Each of these traits can be decomposed into vectors of interrelated components, but treating these components as independent prediction targets using existing univariate LMM or low-dimensional MvLMM genomic prediction tools is inefficient because of their statistical dependence. MegaLMM can be adapted to make joint predictions on vectors of hundreds or thousands of correlated trait components, which could be fed into high-dimensional selection indices for efficient selection of these important plant characteristics. In human genetics, MegaLMM may provide a way to derive multi-ethnic polygenic risk scores [55] by treating outcomes within each ethnic, geographic, or other stratified population group as correlated traits, similar to the analysis of the multi-environment trials above.

MegaLMM should be straightforward to extend to more flexible genetic models including the Bayesian Alphabet family of mixture priors on marker effect sizes. These effects can be incorporated into the parameters  $\mathbf{B}_{2R}$  and  $\mathbf{B}_{2F}$  by adapting the prior structure. This will be further explored in future manuscripts.

Lastly, we have only focused on Gaussian MvLMMs, in which observations are assumed to marginally follow a Gaussian distribution. However, many other types of data require

more flexible models. It should be possible to extend MegaLMM to the broader family of generalized LMMs. These approaches model the relationships among predictor variables in a latent space, which is then related to the observed data through a link function and an exponential family error distribution. More generally, link-functions could be any non-linear function of multiple parameters such as a polynomial or spline basis, or a mechanistic model. In this case, we would model the correlations among model parameters on this link-scale and then use the link-function to relate the latent scale variables to the observed data. Extending MegaLMM to accommodate such generalized LMM structures would require new sampling steps in our MCMC algorithm (see Methods), but we do not see any conceptual challenges with this approach.

## Conclusions

MegaLMM is a flexible and powerful framework for the analysis of very high-dimensional datasets in genetics. Multivariate linear mixed models are widely used for analyzing correlated traits, but have been limited to a maximum of a dozen or so traits at a time by the curse of dimensionality. We developed a novel re-parameterization of the MvLMM that allows powerful statistical regularization and efficient computation with thousands of traits. When applied to real plant breeding objectives, MegaLMM efficiently leverages information across traits to improve genetic value predictions. Our open-source software package will enable users to apply and extend this method in many directions, opening up new areas of research and development in breeding programs.

## Methods

### Multivariate linear mixed models

Multivariate linear mixed models (MvLMMs) are widely used to model multiple sources of covariance among related observations. Let the  $n \times t$  matrix  $\mathbf{Y}$  represent observations on  $t$  traits for  $n$  observational units (i.e., individual plants, plots, or replicates). A general MvLMM takes on the following form

$$\mathbf{Y} = \mathbf{XB} + \mathbf{ZU} + \mathbf{E}, \quad (1)$$

where  $\mathbf{X}$  is a  $n \times b$  matrix of “fixed” effect covariates with effect sizes matrix  $\mathbf{B}$ ,  $\mathbf{U}$  is an  $r \times t$  matrix of random effects for each of the  $t$  traits, with corresponding random effect design matrix  $\mathbf{Z}$ , and  $\mathbf{E}$  is a  $n \times t$  matrix of residuals for each of the  $t$  traits.

MegaLMM uses this formulation to accommodate a large number of designs through different specifications of  $\mathbf{X}$  and  $\mathbf{Z}$ , and different priors on  $\mathbf{B}$ ,  $\mathbf{U}$  and  $\mathbf{E}$ . The distinction between “fixed” and “random” effects in Bayesian mixed models is not well-defined because every parameter requires a prior. However, we use the following distinction here: “fixed” effects are covariates assigned flat (i.e., infinite variance) priors or priors with independent variances on each coefficient; “random” effects, in contrast, are grouped in sets that can be thought of as (possibly correlated) samples from a common population distribution. Generally, “fixed” effects are used to model experimental design terms such as blocks, time, sex, etc, genetic principal components, or specific genetic markers; while “random” effects are used to model genetic values, spatial variation, or related effects.

An important feature of MegaLMM is that multiple random effect terms can be included in the model. We specify this as

$$\mathbf{Z}\mathbf{U} = \sum_{m=1}^M \mathbf{Z}_m \mathbf{U}_m = [\mathbf{Z}_1, \dots, \mathbf{Z}_M] [\mathbf{U}_1^T, \dots, \mathbf{U}_M^T]^T,$$

where each  $\mathbf{Z}_m$  is an  $n \times r_m$  design matrix for a set of related parameters with corresponding coefficient matrix  $\mathbf{U}_m$ . For example,  $\mathbf{U}_1$  may model additive genetic values for each individual, while  $\mathbf{U}_2$  may model spatial environmental effects for each individual. The distribution of each random effect coefficient matrix is  $\mathbf{U}_m \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_m, \mathbf{G}_m)$ , where  $\mathcal{N}(\mathbf{M}, \mathbf{\Sigma}, \mathbf{\Psi})$  is the matrix normal distribution with mean matrix  $\mathbf{M}$ , among-row covariance  $\mathbf{K}_m$  and among-column (i.e., among-trait) covariance  $\mathbf{G}_m$ . We assume that both  $\mathbf{Z}_m$  and  $\mathbf{K}_m$  are known, while  $\mathbf{G}_m$  is unknown and must be learned from the data. Note that  $\mathbf{K}_m$  must be positive semi-definite, while  $\mathbf{G}_m$  is positive-definite. The covariance among different coefficient matrices is assumed to be zero.

To complete the specification of the MvLMM, we assign the residual matrix the distribution  $\mathbf{E} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n, \mathbf{R})$  where  $\mathbf{I}_n$  is the  $n \times n$  identity matrix and  $\mathbf{R}$  is an unknown  $t \times t$  positive-definite covariance matrix.

#### Computational challenges with large multi-trait mixed models

Fitting Eq. (1) is challenging because the columns of  $\mathbf{U}$  and  $\mathbf{E}$  are correlated. This means that data from individual traits (columns of  $\mathbf{Y}$ ) cannot be treated independently. Maximum-likelihood approaches for fitting MvLMMs (e.g., MTG2) compute the full (or restricted) likelihood of  $\mathbf{Y}$ , which involves calculating the inverse of an  $nt \times nt$  matrix many times during model optimization. This is computationally prohibitive when  $n$  and/or  $t$  are large (Fig. 2A). Gibbs samplers (e.g., MCMCglmm) avoid forming and computing the inverse of this extremely large matrix, but still require inverting each of the  $\mathbf{G}_m$  and  $\mathbf{R}$  matrices repeatedly, which is still prohibitive when  $t$  is large. Furthermore, the number of parameters in each  $\mathbf{G}_m$  and  $\mathbf{R}$  grow with the square of  $t$  and quickly get larger than the total number of observations ( $nt$ ) when  $t$  is large. This means that  $\mathbf{G}_m$  and  $\mathbf{R}$  are not identifiable in many datasets and estimates require strong regularization.

#### Mixed effect factor model

If both  $\mathbf{G}_m$  and  $\mathbf{R}$  were diagonal matrices, the  $t$  traits would be uncorrelated. Fitting Eq. (1) then could be done in parallel across traits, greatly reducing the computational burden. While we cannot directly de-correlate traits, if we can identify the sources of variation that cause trait correlations, the residuals of each trait on these causal factors will be un-correlated. We circumvent this issue by re-parameterizing Eq. (1) as a factor model, where we introduce a set of un-observed (or latent) factors that account for all sources of correlation among the traits. Conditional on the values of these factors, the model reduces to a set of independent linear mixed models. Our re-parameterized multi-trait mixed effect factor model is

$$\begin{aligned} \mathbf{Y} &= \mathbf{F}\mathbf{A} + \mathbf{X}_1\mathbf{B}_1 + \mathbf{X}_2\mathbf{B}_{2R} + \mathbf{Z}\mathbf{U}_R + \mathbf{E}_R \\ \mathbf{F} &= \mathbf{X}_2\mathbf{B}_{2F} + \mathbf{Z}\mathbf{U}_F + \mathbf{E}_F \end{aligned} \quad (2)$$

where  $\mathbf{F}$  is an  $n \times K$  matrix of latent factors,  $\mathbf{A}$  is a  $K \times t$  factor loadings matrix,  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$  is a partition of the  $n \times b$  fixed effect covariate matrix between the  $b_1$  covariates with improper priors and the  $b_2 = b - b_1$  covariates with proper priors, and  $\mathbf{U}_R$  and  $\mathbf{U}_F$  coefficients matrices are specified as:



$$\mathbf{U}_R = [\mathbf{U}_{R1}^T, \dots, \mathbf{U}_{RM}^T]^T$$

$$\mathbf{U}_F = [\mathbf{U}_{F1}^T, \dots, \mathbf{U}_{FM}^T]^T.$$

The distributions of the random effects are specified as:

$$\mathbf{U}_{Rm} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_m, \Psi_{Rm}), \quad \mathbf{U}_{Fm} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_m, \Psi_{Fm})$$

$$\mathbf{E}_R \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n, \Psi_{RE}), \quad \mathbf{E}_F \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n, \Psi_{FE})$$

where  $\Psi_{Rm}$ ,  $\Psi_{Fm}$ ,  $\Psi_{RE}$ , and  $\Psi_{FE}$  are all diagonal matrices. Diagonal elements of  $\Psi_{Fm}$  and  $\Psi_{FE}$  are non-negative, while diagonal elements of  $\Psi_{Rm}$  and  $\Psi_{RE}$  are strictly positive.

Conditional on  $\mathbf{F}$  and  $\mathbf{A}$ , the variation in each of the  $t$  columns of  $\mathbf{Y}$  are uncorrelated and can be fitted to the remaining terms independently. Similarly, the  $K$  columns of  $\mathbf{F}$  are also uncorrelated and can be modeled independently as well. Therefore, we can fit Eq. (2) without requiring calculating the inverses of any  $t \times t$  matrices, and many calculations can be done in parallel across different CPU cores.

As long as  $K$  is sufficiently large, Eq. (2) is simply a re-parameterization of Eq. (1). To see how Eq. (2) can represent the terms of Eq. (1), let:

$$\mathbf{B} = [\mathbf{B}_1^T, (\mathbf{B}_{2R} + \mathbf{B}_{2F}\mathbf{A})^T]^T$$

$$\mathbf{U} = \mathbf{U}_R + \mathbf{U}_F\mathbf{A}$$

$$\mathbf{E} = \mathbf{E}_R + \mathbf{E}_F\mathbf{A}$$

Based on the properties of matrix normal random variables, we can integrate over  $\mathbf{U}_R$ ,  $\mathbf{U}_F$ ,  $\mathbf{E}_R$  and  $\mathbf{E}_F$  to calculate the distributions of each  $\mathbf{U}_m$  and  $\mathbf{E}$  as:

$$\mathbf{U}_m \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_m, \Psi_{Rm} + \mathbf{A}^T \Psi_{Fm} \mathbf{A})$$

$$\mathbf{E} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n, \Psi_{RE} + \mathbf{A}^T \Psi_{FE} \mathbf{A})$$

Therefore, each  $\mathbf{G}_m$  is re-parameterized as  $\Psi_{Rm} + \mathbf{A}^T \Psi_{Fm} \mathbf{A}$  and  $\mathbf{R}$  is re-parameterized as  $\Psi_{RE} + \mathbf{A}^T \Psi_{FE} \mathbf{A}$ , such that all off-diagonal elements of each matrix are controlled by  $\mathbf{A}$ .

Although these equations appear to imply that our mixed effect factor model constrains  $\mathbf{B}$ ,  $\mathbf{U}$  and  $\mathbf{E}$  (and thus each  $\mathbf{G}_m$  and  $\mathbf{R}$ ) to be correlated due to the shared dependence on  $\mathbf{A}$ , this is not necessarily the case. When any diagonal element of any  $\Psi_{Fx}$  matrix is set to zero, the corresponding row of  $\mathbf{A}$  does not contribute to that term. If at least  $t$  linearly independent rows of  $\mathbf{A}$  contribute to each matrix, then any set of positive-definite matrices can be represented as above. Therefore, we can represent any set of positive-definite matrices  $\mathbf{G}_m$  and  $\mathbf{R}$  with our model as long as  $K \geq t(M + 1)$ .

Of course, the reason that we parameterize our model in this way is that we do expect some correlation among the genetic and residual covariance matrices. From a statistical perspective, when it is reasonable (given the data) to use the same row of  $\mathbf{A}$  for multiple covariance matrices, we can save parameters in the model. From a biological perspective, shared factors provide a biologically realistic explanation for correlations among traits. If we consider the columns of  $\mathbf{F}$  to be  $K$  traits that simply have not been observed, then it is reasonable to propose that each of these traits is regulated by the same sources of genetic and environmental variation as any of the observed traits.

In Eq. (2), the  $K$  latent traits ( $\mathbf{F}$ ) are the key drivers of all phenotypic co-variation among the  $t$  observed traits ( $\mathbf{Y}$ ). These latent traits may not account for all variation in the observed traits. But, by definition, this residual variation (e.g., measurement errors

in each trait) is unique to each trait and uncorrelated with the residual variation in other traits.

#### **Prior parameterization**

The intuitive structure of the mixed effect factor model (Eq. (2) and Fig. 1) makes prior specification and elicitation easier than for Eq. (1) because we do not need to define prior distributions for very large covariance matrices directly. Instead, priors on the random effect variance components and fixed effect regression coefficients are separable and can be described independently, while priors on trait correlations are specified indirectly as priors on the factor loading matrix  $\Lambda$ .

In MegaLMM, we have chosen functional forms for each prior parameter that balance between interpretability (for accurate prior elicitation), and compatibility with efficient computational approaches. For the variance components, we use the non-parametric discrete prior on variance proportions we previously introduced in GridLMM [22] that approximates nearly any joint distribution for multiple random effects. For the factor loadings matrix and matrices of regression coefficients, we use a two-dimensional global-local prior based on the horseshoe prior [56], parameterized in terms of the effective number of non-zero coefficients. For the factor loadings matrix specifically, our prior achieves both regularization and interpretability of the factor traits without having to carefully specify  $K$  itself. Full details of each prior distribution are provided in Additional file 1: Table S1 lists the default hyperparameters for each prior used in the analyses reported here and provided as defaults in the MegaLMM R package.

#### **Computational details and posterior inference**

We use a carefully constructed MCMC algorithm to draw samples from the posterior distribution of each model parameter. We gain efficiency in both per-iteration computational time and in effective samples per iteration through careful uses of diagonalization, sparse matrix algebra, parallelization, and integration (or partial collapsing). In particular, our algorithm synthesizes and extends several recent innovations in computational approaches to linear mixed models [17, 20, 22, 57]. Full details of the computational algorithm are provided in Additional file 1: Supplemental Methods.

#### **Data Analyses**

We demonstrate MegaLMM using three example datasets.

##### **Scaling performance with gene expression data**

To compare the scalability of MegaLMM to other multi-trait mixed model programs, we used a large gene expression dataset of 24,175 genes across 728 *Arabidopsis thaliana* accessions. We downloaded the data from NCBI GEO [58] [59] and removed genes with average counts  $< 10$ . We then normalized and variance stabilized the counts using the `varianceStabilizingTransformation` function from DESeq2 [60]. We downloaded a corresponding genomic relationship matrix  $\mathbf{K}$  from the 1001 genomes project [61] and subsetting to the 665 individuals present in both datasets.

We generated datasets of varying sizes from  $t = 2$  to  $t = 24,175$  genes by randomly sampling. We selected one gene as the “focal” trait in each dataset, masked 50% of its values, fit the model in Eq. (1) using four different representative MvLMM programs to the remaining data, and used the results to predict the genetic values of each masked indi-

dual for this “focal” gene. Prediction accuracies were estimated as  $\rho_g = \text{cor}_g(\hat{\mathbf{u}}, \mathbf{y})\sqrt{h^2(\hat{\mathbf{u}})}$ , where  $\text{cor}_g$  is the estimated genetic correlation evaluated in the testing lines only, and  $h^2(\hat{\mathbf{u}})$  is the heritability of the predictor  $\hat{\mathbf{u}}$  estimated using a univariate LMM [6, 42]. The simpler Pearson’s correlation estimate of prediction accuracy is not valid in these data because all genes were measured together in the same sample, and therefore some correlation among genes is caused by non-genetic factors [24]. The four MvLMM prediction methods were:

1. MTG2 [23]: a restricted maximum-likelihood method written in fortran. We pre-calculated the eigenvalue decomposition for  $\mathbf{K}$ , thus this additional time is not included in the results. MTG2 does not work well with a high percentage of missing data, so genetic value predictions were made with the two-step approach from [24] which involves estimating model parameters only from the individuals with complete observations, and then incorporating secondary trait values of the new individuals in the second step.
2. MCMCg1mm [62]: a Bayesian MCMC algorithm largely written in C++. We used “default” priors for  $\mathbf{R}$  and  $\mathbf{G}$  with diagonal means and  $\nu = p$ , and ran a single MCMC chain for 7000 iterations, discarding the first 5000 samples as burnin. To speed up calculations (and make the timing results more comparable with the MegaLMM algorithm), we rotated the input data by pre-multiplying by the eigenvectors of  $\mathbf{K}$  so that the input relationship matrix was diagonal. Since this matrix rotation is only possible with complete data, we again used the two-step multi-trait prediction approach [24].
3. phenix [18]: a variational Bayes algorithm written in R that uses a low-rank representation of  $\mathbf{G}$  but a full-rank prior for  $\mathbf{R}$ . We set the maximum number of factors to  $p/4$  and used the eigendecomposition of  $\mathbf{K}$  as the input. Again, we excluded this calculation from the time estimates.
4. MegaLMM: we ran MegaLMM using “default priors” with  $K = \min(n/4, p/2)$  and collected 6000 MCMC samples, discarding the first 5000 as burnin. We excluded the preparatory calculations, only including the MCMC iterations in the time calculations. For small datasets, these calculations were significant, but were a miniscule part of the analyses of larger datasets.

Each method was run 20 times on different randomly sampled datasets. For the two MCMC methods, we estimated the effective sample size of each element of  $\mathbf{U}$  using the `ess_bulk` function of the `rstan` package [63], and used this to estimate the time necessary for the effective sample size to be at least 1000 for 90% of the  $u_{ij}$ . We ran MTG2 and MCMCg1mm for datasets up to  $t = 64$  because computational times were prohibitively long for larger datasets. We linearly extrapolated the (log) computational times for these methods out to  $t = 512$  for comparisons. phenix fails when  $t \geq n$ , so this method is limited to  $t < 665$  in this dataset.

To assess the accuracy of each method for estimating genetic and non-genetic covariances, we generated new datasets with 128 genes by calculating empirical correlation matrices for  $\mathbf{G}$  and  $\mathbf{R}$  from two separate samples of 128 genes from the full expression dataset, and then generating genetic and residual values for 128 traits from multivariate normal distributions based on these correlation matrices. For each trait, we converted the correlation matrices into covariance matrices by sampling an independent heritability

lity value for each trait between 0.1 and 0.8. We then estimated the genetic and residual covariance matrices for subsets of these simulated datasets using each of the four above methods. In this example, we found that setting  $K$  larger ( $2p$ ) gave better results, probably because the  $\mathbf{G}$  and  $\mathbf{R}$  matrices were largely uncorrelated and so independent factors were needed to model the two sets of covariances. Accuracy was reported as the Pearson correlation between the estimated covariance parameters and the true covariance parameters (excluding the variance parameters on the diagonal).

#### ***Wheat yield prediction using hyperspectral data***

We used data from a bread wheat breeding trial to demonstrate how MegaLMM can leverage “secondary” traits from high-throughput phenotyping technologies to better predict genetic values of a single target trait. We downloaded grain yield and hyperspectral reflectance data from the bread wheat trials at the Campo Experimental Norman E. Borlaug in Ciudad Obregón, México reported in [25] [64]. We selected the 2014–2015 Optimal Flat site-year for our main analysis because it had among the greatest number of hyperspectral reflectance data points, and [25] reported relatively low predictive accuracy for grain yield in this site-year. Best linear unbiased estimates (BLUEs) and best linear unbiased predictors (BLUPs) of the line means for grain yield (GY) and 62 hyperspectral bands collected at each of 10 time-points during the growing season, and genotype data from 8519 markers were provided for 1,092 lines in this trial. All other trials were analyzed in the analysis presented in Additional file 1: Figure S5.

We compared eight methods for predicting the GY trait based on the genetic marker and hyperspectral data. The first five were “standard” methods using state-of-the-art models for genomic prediction. The final three were new models implemented within the MegaLMM framework.

1. GBLUP: implemented using the R package `rrBLUP` [65], with the genomic relationship matrix  $\mathbf{K}$  calculated with the `A.mat` function of `rrBLUP` as in [66].
2. *Bayesian Lasso* (BL): implemented using the R package `BGLR` [67]. We first removed markers with > 50% missing data, and imputed the remaining missing genotypes with the population mean allele frequency. We used the default prior parameters for the *Bayesian Lasso* in `BGLR`, and collected 9,000 posterior samples with a thinning rate of 5 after a 5,000 iteration burnin.
3. RKHS: implemented using `rrBLUP`. We used the same thinned and imputed genotype dataset as for the BL method to calculate a genomic distance matrix ( $\mathbf{D}$ ). We also used the default `theta.seq` parameter to automatically choose the scale parameter of the Gaussian kernel.
4. HBLUP: implemented using the R package `lme4qt1`. This replicates the analysis reported by [25], which uses the GBLUP method but replaces the genomic relationship matrix described above with  $\mathbf{H}$ , a hyperspectral reflectance relationship matrix calculated as  $\mathbf{H} = \mathbf{S}\mathbf{S}^T/620$ , where  $\mathbf{S}$  is a matrix of centered and standardized BLUEs of hyperspectral bands from each timepoint.
5. GBLUP+H: implemented in the R package `lme4qt1` [68]. This is a two-kernel method, where we use two relationship matrices:  $\mathbf{K}$  and  $\mathbf{H}$ . This method is analogous to the methods proposed by [25] for leveraging the hyperspectral data in prediction; however, those authors only used two-kernel methods for  $G \times E$  prediction across site-years. Since `lme4qt1` does not predict random effects for

un-measured observations, we formed predictions as:  $\mathbf{K}_{no}\mathbf{K}_{oo}^{-1}\hat{\mathbf{u}}_{ko} + \mathbf{H}_{no}\mathbf{H}_{oo}^{-1}\hat{\mathbf{u}}_{ho}$  where  $\mathbf{K}_{no}$  is the  $n_n \times n_o$  quadrant of  $\mathbf{K}$  specifying the genomic relationships among the  $n_n$  “new” un-observed lines,  $\mathbf{K}_{oo}$  is the  $n_o \times n_o$  quadrant of  $\mathbf{K}$  specifying the genomic relationships among the “old” observed lines,  $\hat{\mathbf{u}}_{ko}$  is the vector of BLUPs for the genomic random effect, and  $\mathbf{H}_{no}$ ,  $\mathbf{H}_{oo}$  and  $\hat{\mathbf{u}}_{ho}$  are similar quantities for the hyperspectral random effect.

6. **MegaLMM-GBLUP**: we modeled the combined trait data  $\mathbf{Y} = [\mathbf{y}, \mathbf{S}]$  with the model specified in Eq. (2) using a single random effect with relationship matrix  $\mathbf{K}$  as above, no fixed effects besides an intercept ( $\mathbf{X}$  was a column of ones and  $\mathbf{X}_2$  had zero columns). We ran MegaLMM with  $K = 100$  factors, “default” priors (see Additional file 1: Table S1), and two partitions of the trait data (the first containing grain yield with the masked training set as described below, and the second containing all 620 hyperspectral bands with complete data). We collected 500 posterior samples of the quantity:  $\mathbf{u}_1 = \mathbf{u}_{R1} + (\mathbf{U}_F\boldsymbol{\lambda}_1)$  at a thinning rate of 2, discarding the first 1,000 iterations as burn-in.
7. **MegaLMM-RKHS**: we implemented multi-trait RKHS regression model using the “kernel-averaging” method proposed by [28]. We standardized  $\mathbf{D}$  based on its mean (squared) value, and placed a uniform prior on the set of scaling factors  $h = \{1/5, 1, 5\}$ , which we implemented by calculating three corresponding relationship matrices  $\mathbf{K}_1, \dots, \mathbf{K}_3$  and by specifying three random effects in Eq. (2). We again used “default” priors,  $K = 100$  factors, and treated only the global intercept per-trait as fixed effects. We collected 500 posterior samples of the quantity:  $\mathbf{Zu}_1 = \mathbf{Zu}_{R1} + \mathbf{Z}(\mathbf{U}_F\boldsymbol{\lambda}_1)$  at a thinning rate of 2, discarding the first 1000 iterations as burn-in.
8. **MegaLMM-GBLUP-CV1**: we repeated the MegaLMM-GBLUP method above, but this time without partitioning the trait data. Instead, we masked both the grain yield and the 620 hyperspectral band data from the testing set so all lines in the training data had complete data. Predictions of the genetic values were calculated identically to above.

We used cross-validation to evaluate the prediction accuracy of each method. We randomly selected 50% of the lines for model training, 50% for testing, and masked the GY observations for the testing lines. We fit each model to the partially-masked dataset and collected the predictions of GY for the testing lines. We estimated prediction accuracy as  $\rho_g = \text{cor}_g(\hat{\mathbf{u}}, \mathbf{y})\sqrt{h^2(\hat{\mathbf{u}})}$  because the hyperspectral reflectance data were collected on the same plots as the GY data and therefore non-genetic (i.e., microenvironmental) factors that affect both reflectance and yield may induce non-genetic correlations among traits [24]. BLUPs were used as the predictand except in the 2016–17 year when the BLUPs were poorly correlated with the BLUEs suggesting data quality issues. We used a 50–50 training-testing split of the data to ensure that  $\text{cor}_g$  could be estimated accurately in the testing partition. This cross-validation algorithm was repeated 20 times with different random partitions.

#### **Maize trait imputation in multi-environment trials**

We used data on maize hybrids from the Genomes-To-Fields Initiative experiments to demonstrate how MegaLMM can leverage genetic correlations across locations in multi-environment trials. We downloaded the agronomic data from the 2014–2017 field seasons

from the CyVerse data repository [69] and corresponding genomic data. We used TASSEL5 [70] to build a kinship matrix for each hybrid genotype using the CenteredIBS routine.

A total of 2012 non-check hybrids with phenotype and genotype data from 108 trials (i.e., site-years) were available. We selected four representative agronomic traits: plant height (cm), grain yield (bushels/acre), days-to-silking (days), and the anthesis-silking interval (ASI, days). For each trait in each site-year, we calculated BLUPs for all observed genotypes using the R package lme4 [71] with Rep and Block:Rep as fixed effects to account for the experimental design in each field, and formed them into  $2012 \times 108$  BLUP matrices for each trait. We then dropped site-years where the BLUP variance was zero, or which had fewer than 50 tested lines. On average  $\approx 12\%$  of hybrid-site-year combinations were observed across each of the four BLUP matrices. We then used four methods to predict the BLUPs of hybrids that were not grown in each trial:

1. `GBLUP (univariate)`: missing values were imputed separately for each site:year using the `mixed.solve` function of the `rrBLUP` package.
2. `GBLUP (env BLUPs)`: genetic values for each hybrid were assumed to be constant across all site-years. We estimated these in two steps. In the first step, we estimated hybrid main effects treating lines as independent random effects using lme4, with site:year included as a fixed effect. In the second step, we estimated genetic values using the `mixed.solve` function of the `rrBLUP` package.
3. `phenix`: we used phenix to impute missing observations in  $\mathbf{Y}$  using  $\mathbf{K}$  as a relationship matrix.
4. `MegaLMM`: we fit the model specified in Eq. (2) to the full matrix  $\mathbf{Y}$ , with  $K = 50$  factors and “default”. Here, we partitioned  $\mathbf{Y}$  into 4 sets based on year to minimize the number of missing observations to condition on during the MCMC. We collected 1000 posterior samples of imputed values  $\tilde{\mathbf{Y}} = \mathbf{X}_1\mathbf{B}_1 + \mathbf{F}\mathbf{A} + \mathbf{Z}\mathbf{U}_R$  with a thinning rate of 2, after discarding the first 5000 iterations as burnin.

We estimated prediction accuracy of each method using cross-validation. For each of 20 replicate cross-validation runs per model, we randomly masked 20% of the non-missing BLUPs, and then calculated the Pearson’s correlation between these “observed” values and the imputed values of each method. Pearson’s correlation is appropriate as an estimate of genomic prediction accuracy in this case because different plants were used in each trial, so there is no non-genetic source of correlation among site-years that may bias accuracy estimates.

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-021-02416-w>.

**Additional file 1: Supplementary Methods.** Expanded description of the prior distributions and derivation of the MCMC algorithm used in MegaLMM [72–76]. **Figure S1.** Prediction accuracies for each target gene expression trait relative to MegaLMM. **Figure S2.** Prediction accuracies for each target gene expression trait as a function of dataset size. **Figure S3.** Accuracy of covariance estimates in simulated data. **Figure S4.** Sensitivity of MegaLMM to changes in  $K$ . **Figure S5.** Performance of single-trait and multi-trait genomic prediction for wheat yield across 20 experiments. **Figure S6.** Relationship between genetic correlation and benefit of MvLMM across Genomes2Fields site-years. **Table S1.** Default hyperparameters for user-customizable prior distributions.

**Additional file 2:** Review history.

### Acknowledgments

Not applicable

**Review history**

The review history is available as Additional file 2.

**Authors' contributions**

DER developed the method, wrote the R package, developed and ran the analyses, and wrote the paper. JQ edited the manuscript. HC helped develop the method, design the analysis, and edited the paper. LC helped develop the method, design the analysis, and wrote the paper. The authors read and approved the final manuscript.

**Funding**

This work is supported by Agriculture and Food Research Initiative grants no. 2020-67013-30904 and 2018-67015-27957 from the USDA National Institute of Food and Agriculture to DER and HC. DER is also supported by United States Department of Agriculture (USDA) National Institute of Food and Agriculture (NIFA), Hatch project 1010469. LC is supported by grants P20GM109035 (COBRE Center for Computational Biology of Human Disease; PI Rand) and P20GM103645 (COBRE Center for Central Nervous; PI Sanes) from the NIH NIGMS, 2U10CA180794-06 from the NIH NCI and the Dana Farber Cancer Institute (PIs Gray and Gatsonis), as well as by an Alfred P. Sloan Research Fellowship and a David & Lucile Packard Fellowship for Science and Engineering.

**Availability of data and materials**

All data used in these analyses were downloaded from the publicly accessible repositories described above. Arabidopsis gene expression data was downloaded from the NCBI GEO accession GSE80744 available at <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE80744> [59]. The Arabidopsis kinship matrix was downloaded from [https://1001genomes.org/data/GMI-MPI/releases/v3.1/SNP\\_matrix\\_imputed\\_hdf5/1001\\_SNP\\_MATRIX.tar.gz](https://1001genomes.org/data/GMI-MPI/releases/v3.1/SNP_matrix_imputed_hdf5/1001_SNP_MATRIX.tar.gz). The wheat dataset was downloaded from the CIMMYT Research Data & Software Repository Network available at <http://hdl.handle.net/11529/10548109> [77]. The maize phenotype data were downloaded from the CyVerse data repository based on the links described in [69]. Genomic data were downloaded from ([http://datacommons.cyverse.org/browse/iplant/home/shared/commons\\_repo/curated/Carolyn\\_Lawrence\\_Dill\\_G2F\\_Nov\\_2016\\_V3/b\\_2014\\_gbs\\_data](http://datacommons.cyverse.org/browse/iplant/home/shared/commons_repo/curated/Carolyn_Lawrence_Dill_G2F_Nov_2016_V3/b_2014_gbs_data)) [78]. Scripts for running analyses are available in the GitHub repository: [https://github.com/deruncie/MegaLMM\\_analyses/tree/v0.9.2](https://github.com/deruncie/MegaLMM_analyses/tree/v0.9.2). The R package for MegaLMM is available here: <https://github.com/deruncie/MegaLMM/tree/v0.9.3> and is licensed with the MIT license. The specific versions of the scripts and package codes are archived at zenodo [79, 80].

**Declarations****Ethics approval and consent to participate**

Not applicable

**Consent for publication**

Not applicable

**Competing interests**

LC was employed by Microsoft Research while performing this work.

**Author details**

<sup>1</sup>Department of Plant Sciences, University of California Davis, Davis, CA, USA. <sup>2</sup>Microsoft Research New England, Cambridge, MA, USA.

Received: 2 July 2020 Accepted: 23 June 2021

Published online: 23 July 2021

**References**

1. Araus JL, Kefauver SC, Zaman-Allah M, Olsen MS, Cairns JE. Translating High-Throughput Phenotyping into Genetic Gain. *Trends Plant Sci.* 2018;23(5):451–66.
2. Koltes JE, Cole JB, Clemmens R, Dilger RN, Kramer LM, Lunney JK, McCue ME, McKay SD, Mateescu RG, Murdoch BM, Reuter R, Rexroad CE, Rosa GJM, Serão NVL, White SN, Woodward-Greene MJ, Worku M, Zhang H, Reecy JM. A vision for development and utilization of high-throughput phenotyping and big data analytics in livestock. *Front Genet.* 2019;10:1197. <https://doi.org/10.3389/fgene.2019.01197>.
3. Rutkoski J, Poland J, Mondal S, Autrique E, Pérez LG, Crossa J, Reynolds M, Singh R. Canopy Temperature and Vegetation Indices from High-Throughput Phenotyping Improve Accuracy of Pedigree and Genomic Selection for Grain Yield in Wheat. *G3 Genes Genomes Genetics.* 2016;6(9):2799–808.
4. Neethirajan S. Recent advances in wearable sensors for animal health management. *Sens and Bio-Sens Res.* 2017;12: 15–29.
5. Schrag TA, Westhues M, Schipprack W, Seifert F, Thiemann A, Scholten S, Melchinger AE. Beyond genomic prediction: combining different types of omics data can improve prediction of hybrid performance in maize. *Genetics.* 2018;208(4):1373–85.
6. Thompson R, Meyer K. A review of theoretical aspects in the estimation of breeding values for multi-trait selection. *Livest Prod Sci.* 1986;15(4):299–313.
7. Bernardo R. *Breeding for Quantitative Traits in Plants*, vol 1. 2nd ed. Woodbury: Stemma press; 2010.
8. Meuwissen TH, Hayes BJ, Goddard ME. Prediction of total genetic value using genome-wide dense marker maps. *Genetics.* 2001;157(4):1819–29.
9. Henderson CR, Quaas RL. Multiple Trait Evaluation Using Relatives' Records. *J Anim Sci.* 1976;43(6):1188–97.
10. Piepho HP, Möhring J, Melchinger AE, Büchse A. BLUP for phenotypic selection in plant breeding and variety testing. *Euphytica.* 2007;161(1–2):209–28.



11. Calus MP, Veerkamp RF. Accuracy of multi-trait genomic selection using different methods. *Genet Sel Evol.* 2011;43(1):26.
12. Jia Y, Jannink J-L. Multiple-Trait Genomic Selection Methods Increase Genetic Value Prediction Accuracy. *Genetics.* 2012;192(4):1513–22.
13. Johnstone IM, Titterton DM. Statistical challenges of high-dimensional data. *Phil Trans Ser A Math Phys Eng Sci.* 2009;367(1906):4237–53.
14. Zhou X, Stephens M. Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nat Methods.* 2014;11(4):407–9.
15. de Los Campos G, Gianola D. Factor analysis models for structuring covariance matrices of additive genetic effects: a Bayesian implementation. *Genet Sel Evol.* 2007;39(5):481–94.
16. Meyer K. Multivariate analyses of carcass traits for Angus cattle fitting reduced rank and factor analytic models. *J Anim Breed Genet.* 2007;124(2):50–64.
17. Runcie D, Mukherjee S. Dissecting High-Dimensional Phenotypes with Bayesian Sparse Factor Analysis of Genetic Covariance Matrices. *Genetics.* 2013;194(3):753–67.
18. Dahl A, Iotchkova V, Baud A, Johansson Å, Gyllensten U, Soranzo N, Mott R, Kranis A, Marchini J. A multiple-phenotype imputation method for genetic studies. *Nat Genet.* 2016;48(4):466–72.
19. Kang HM, Zaitlen NA, Wade CM, Kirby A, Heckerman D, Daly MJ, Eskin E. Efficient control of population structure in model organism association mapping. *Genetics.* 2008;178(3):1709–23.
20. Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association studies. *Nat Genet.* 2012;44(7):821–4.
21. Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST linear mixed models for genome-wide association studies. *Nat Methods.* 2011;8(10):833–5.
22. Runcie D, Crawford L. Fast and flexible linear mixed models for genome-wide genetics. *PLOS Genet.* 2019;15(2):1007978.
23. Lee SH, van der Werf JHJ. MTG2: an efficient algorithm for multivariate linear mixed model analysis based on genomic information. *Bioinformatics.* 2016;32(9):1420–2.
24. Runcie D, Cheng H. Pitfalls and remedies for cross validation with multi-trait genomic prediction methods. *G3 Genes Genomes Genet.* 2019;9(11):3727–41. <https://doi.org/10.1534/g3.119.400598>.
25. Krause MR, González-Pérez L, Crossa J, Pérez-Rodríguez P, Montesinos-López O, Singh RP, Dreisigacker S, Poland J, Rutkoski J, Sorrells M, Gore MA, Mondal S. Hyperspectral Reflectance-Derived Relationship Matrices for Genomic Prediction of Grain Yield in Wheat. *G3 Genes Genomes Gene.* 2019;9(4):1231–47.
26. Hayes BJ, Bowman PJ, Chamberlain AC, Verbyla K, Goddard ME. Accuracy of genomic breeding values in multi-breed dairy cattle populations. *Genet Sel Evol.* 2009;41(1):1–9.
27. Park T, Casella G. The Bayesian Lasso. *J Am Stat Assoc.* 2013;103(482):681–6.
28. de Los Campos G, Gianola D, Rosa GJM, Weigel KA, Crossa J. Semi-parametric genomic-enabled prediction of genetic values using reproducing kernel Hilbert spaces methods. *Genet Res.* 2010;92(4):295–308.
29. Burguño J, de los Campos G, Weigel K, Crossa J. Genomic prediction of breeding values when modeling genotype x environment interaction using pedigree and dense molecular markers. *Crop Sci.* 2012;52(2):707–19. <https://doi.org/10.2135/cropsci2011.06.0299>.
30. Piepho HP, Möhring J. Best Linear Unbiased Prediction of Cultivar Effects for Subdivided Target Regions. *Crop Sci.* 2005;45(3):1151–9.
31. Loh P-R, Tucker G, Bulik-Sullivan BK, Vilhjálmsson BJ, Finucane HK, Salem RM, Chasman DI, Ridker PM, Neale BM, Berger B, Patterson N, Price AL. Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet.* 2015;47(3):284–90.
32. Gilmour AR. Mixed model regression mapping for QTL detection in experimental crosses. *Comput Stat Data Anal.* 2007;51(8):3749–64.
33. GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature.* 2017;550(7675):204–13.
34. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D, Delaneau O, O'Connell J, Cortes A, Welsh S, Young A, Effingham M, McVean G, Leslie S, Allen N, Donnelly P, Marchini J. The UK Biobank resource with deep phenotyping and genomic data. *Nature.* 2018;562(7726):203–9.
35. Guo G, Zhao F, Wang Y, Zhang Y, Du L, Su G. Comparison of single-trait and multiple-trait genomic prediction models. *BMC Genet.* 2014;15(1):30.
36. Sun J, Rutkoski JE, Poland JA, Crossa J, Jannink JL, Sorrells ME. Multitrait, Random Regression, or Simple Repeatability Model in High-Throughput Phenotyping Data Improve Genomic Prediction for Wheat Grain Yield. *Plant Genome.* 2017;10(2):0.
37. Crain J, Mondal S, Rutkoski J, Singh RP, Poland J. Combining High-Throughput Phenotyping and Genomic Information to Increase Prediction and Selection Accuracy in Wheat Breeding. - PubMed - NCBI. *Plant Genome.* 2018;11(1):1–14.
38. van Eeuwijk FA, Bustos-Korts D, Millet EJ, Boer MP, Kruijer W, Thompson A, Malosetti M, Iwata H, Quiroz R, Kuppe C, Muller O, Blazakis KN, Yu K, Tardieu F, Chapman SC. Modelling strategies for assessing and increasing the effectiveness of new phenotyping techniques in plant breeding. *Plant Sci.* 2019;282:23–39.
39. Montesinos-López A, Montesinos-López OA, Cuevas J, Mata-López WA, Burguño J, Mondal S, Huerta J, Singh R, Autrique E, González-Pérez L, Crossa J. Genomic Bayesian functional regression models with interactions for predicting wheat grain yield using hyper-spectral image data. *Plant Methods.* 2017;13(1):1.
40. Cuevas J, Montesinos-López O, Juliana P, Guzman C, Pérez-Rodríguez P, González-Bucio J, Burguño J, Montesinos-López A, Crossa J. Deep Kernel for Genomic and Near Infrared Predictions in Multi-environment Breeding Trials. *G3 Genes Genomes Genet.* 2019;9(9):2913–24.
41. Juliana P, Montesinos-López OA, Crossa J, Mondal S, González-Pérez L, Poland J, Huerta-Espino J, Crespo-Herrera L, Govindan V, Dreisigacker S, Shrestha S, Pérez-Rodríguez P, Pinto Espinosa F, Singh RP. Integrating genomic-enabled prediction and high-throughput phenotyping in breeding for climate-resilient bread wheat. *Theor Appl Genet.* 2019;132(1):177–94.

42. Lopez-Cruz M, Olson E, Rovere G, Crossa J, Dreisigacker S, Mondal S, Singh R, de Los Campos G. Regularized selection indices for breeding value prediction using hyper-spectral image data. *bioRxiv*. 2020;125:625251.
43. Heffner EL, Sorrells ME, Jannink J-L. Genomic Selection for Crop Improvement. *Crop Sci*. 2009;49(1):1–12.
44. Gauch HG. Model Selection and Validation for Yield Trials with Interaction. *Biometrics*. 1988;44(3):705–15.
45. Piepho H-P. Empirical best linear unbiased prediction in cultivar trials using factor-analytic variance-covariance structures. *Theor Appl Genet*. 1998;97(1):195–201.
46. Smith A, Cullis B, Thompson R. Analyzing Variety by Environment Data Using Multiplicative Mixed Models and Adjustments for Spatial Field Trend. *Biometrics*. 2001;57(4):1138–47.
47. Jarquín D, Crossa J, Lacaze X, Du Cheyron P, Daucourt J, Lorgeou J, Piraux F, Guerreiro L, Pérez P, Calus M, Burgueño J, de Los Campos G. A reaction norm model for genomic selection using high-dimensional genomic and environmental data. *Theor Appl Genet*. 2014;127(3):595–607.
48. Malosetti M, Bustos-Korts D, Boer MP, van Eeuwijk FA. Predicting Responses in Multiple Environments: Issues in Relation to Genotype × Environment Interactions. *Crop Sci*. 2016;56(5):2210–22.
49. Rincent R, Malosetti M, Ababaei B, Touzy G, Mini A, Bogard M, Martre P, Le Gouis J, van Eeuwijk FA. Using crop growth model stress covariates and AMMI decomposition to better predict genotype-by-environment interactions. *TAG Theor Appl Genet Theor Angew Genet*. 2019;132(12):3399–411.
50. The International Schizophrenia Consortium. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature*. 2009;460(7256):748–52.
51. Turley P, Walters RK, Maghazian O, Okbay A, Lee JJ, Fontana MA, Nguyen-Viet TA, Wedow R, Zacher M, Furlotte NA, et al. Multi-trait analysis of genome-wide association summary statistics using mtag. *Nat Genet*. 2018;50(2):229–37.
52. Campbell M, Walia H, Morota G. Utilizing random regression models for genomic prediction of a longitudinal trait derived from high-throughput phenotyping. *Plant Direct*. 2018;2(9):00080.
53. Chan EKF, Rowe HC, Corwin JA, Joseph B, Kliebenstein DJ. Combining genome-wide association mapping and transcriptional networks to identify novel genes controlling glucosinolates in *Arabidopsis thaliana*. *PLoS Biol*. 2011;9(8):1001125.
54. Demmings EM, Williams BR, Lee C-R, Barba P, Yang S, Hwang C-F, Reisch BI, Chitwood DH, Londo JP. Quantitative Trait Locus Analysis of Leaf Morphology Indicates Conserved Shape Loci in Grapevine. *Front Plant Sci*. 2019;10:36.
55. Márquez-Luna C, Loh P-R, Consortium SATDS, Consortium TSTD, Price AL. Multi-ethnic polygenic risk scores improve risk prediction in diverse populations. *Genet Epidemiol*. 2017;41(8):811–23.
56. Carvalho CM, Polson NG, Scott JG. The horseshoe estimator for sparse signals. *Biometrika*. 2010;97(2):465–80.
57. Makalic E, Schmidt DF. A Simple Sampler for the Horseshoe Estimator. *IEEE Signal Process Lett*. 2016;23(1):179–82.
58. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA, Phillippy KH, Sherman PM, Holko M, et al. Ncbi geo: archive for functional genomics data sets—update. *Nucleic Acids Res*. 2012;41(D1):991–5.
59. Huang S, Kawakatsu T, Jupe F, Schmitz R, Ulrich M, Castanon R, Nery J, Chen H, Ecker J. Epigenomic and genome structural diversity in a worldwide collection of *Arabidopsis thaliana*. *NCBI Gene Expr Omnibus*. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE80744>. Accessed 05 Sept 2018.
60. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome Biol*. 2014;15:550. <https://doi.org/10.1186/s13059-014-0550-8>.
61. Alonso-Blanco C, Andrade J, Becker C, Bemm F, Bergelson J, Borgwardt KM, Cao J, Chae E, Dezwaan TM, Ding W, et al. 1,135 genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana*. *Cell*. 2016;166(2):481–91.
62. Hadfield JD. MCMC methods for multi-response generalized linear mixed models: the MCMCglmm R package. *J Stat Softw*. 2010;33(1):1–22.
63. Stan Development Team. RStan: the R interface to Stan. 2019. R package version 2.19.2 <http://mc-stan.org/>.
64. Mondal S, Krause M, Juliana P, Poland J, Dreisigacker S, Singh R. Use of hyperspectral reflectance-derived relationship matrices for genomic prediction of grain yield in wheat - data for publication. *CIMMYT Res Data Softw Repository Netw*. 2018. <https://hdl.handle.net/11529/10548109>.
65. Endelman JB. Ridge regression and other kernels for genomic selection with r package rrblup. *Plant Genome*. 2011;4:250–55.
66. Endelman JB, Jannink J-L. Shrinkage Estimation of the Realized Relationship Matrix. *G3 Genes Genomes Genet*. 2012;2(11):1405–13.
67. Perez P, de los Campos G. Genome-wide regression and prediction with the bgrr statistical package. *Genetics*. 2014;198(2):483–95.
68. Ziyatdinov A, Vazquez-Santiago M, Brunel H, Martinez-Perez A, Aschard H, Soria JM. lme4qtl: linear mixed models with flexible covariance structure for genetic studies of related individuals. *BMC Bioinformatics*. 2018;19:080. <http://dx.doi.org/10.1186/s12859-018-2057-x>.
69. McFarland BA, AlKhalifah N, Bohn M, Bubert J, Buckler ES, Ciampitti I, Edwards J, Ertl D, Gage JL, Falcon CM, Flint-Garcia S, Gore MA, Graham C, Hirsch CN, Holland JB, Hood E, Hooker D, Jarquín D, Kaeppler SM, Knoll J, Kruger G, Lauter N, Lee EC, Lima DC, Lorenz A, Lynch JP, McKay J, Miller ND, Moose SP, Murray SC, Nelson R, Poudyal C, Rocheford T, Rodriguez O, Romay MC, Schnable JC, Schnable PS, Scully B, Sekhon R, Silverstein K, Singh M, Smith M, Spalding EP, Springer N, Thelen K, Thomison P, Tuinstra M, Wallace J, Walls R, Wills D, Wissler RJ, Xu W, Yeh C-T, de Leon N. Maize genomes to fields (G2F): 2014–2017 field seasons: genotype, phenotype, climatic, soil, and inbred ear image datasets. *BMC Res Notes*. 2020;13(1):1–6.
70. Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. Tassel: software for association mapping of complex traits in diverse samples. *Bioinformatics*. 2007;23(19):2633–5.
71. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *J Stat Softw*. 2015;67(1):1–48. <https://doi.org/10.18637/jss.v067.i01>.
72. Bates D, Edelman D. Fast and elegant numerical linear algebra using the RcppEigen package. *J Stat Softw*. 2013;52(5):1–24. <http://www.jstatsoft.org/v52/i05/>.
73. Anirban B, Antik C, Mallick BK. Fast sampling with Gaussian scale mixture priors in high-dimensional regression. *Biometrika*. 2016;103(4):985–91. <https://doi.org/10.1093/biomet/asw042>. <https://academic.oup.com/biomet/article-pdf/103/4/985/8339159/asw042.pdf>.

74. Bhattacharya A, Dunson DB. Sparse Bayesian infinite factor models. *Biometrika*. 2011;98(2):291–306.
75. Gelman A. Prior distributions for variance parameters in hierarchical models. *Bayesian Anal*. 2006;1(3):515–33.
76. Piironen J, Vehtari A. Sparsity information and regularization in the horseshoe and other shrinkage priors. *Electron J Stat*. 2017;11(2):5018–51.
77. Mondal S, Krause M, Juliana P, Poland J, Dreisigacker S, Singh R. Use of hyperspectral reflectance-derived relationship matrices for genomic prediction of grain yield in wheat - data for publication. 2018. <https://hdl.handle.net/11529/10548109>.
78. Lawrence-Dill C. Genomes To Fields 2014 v.3: CyVerse Data Commons; 2017. [https://datacommons.cyverse.org/browse/iplant/home/shared/commons\\_repo/curated/Carolyn\\_Lawrence\\_Dill\\_G2F\\_Nov\\_2016\\_V3](https://datacommons.cyverse.org/browse/iplant/home/shared/commons_repo/curated/Carolyn_Lawrence_Dill_G2F_Nov_2016_V3).
79. Runcie D. deruncie/MegaLMM: Version for accepted manuscript. Github. 2021. <https://doi.org/10.5281/zenodo.4961220>.
80. Runcie D. deruncie/MegaLMM: Version for accepted manuscript. Github. 2021. <https://doi.org/10.5281/zenodo.4961269>.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

