

Examining the relationships between phenotypic plasticity and local environments with genomic structural equation models

Malachy T. Campbell,^{1,2} Haipeng Yu,¹ Mehdi Momen,¹ Gota Morota¹

¹Department of Animal and Poultry Sciences, Virginia Polytechnic Institute and State University, Blacksburg, VA, USA 24061

²Plant Breeding and Genetics Section, School of Integrative Plant Science, Cornell University, Ithaca, NY, USA 14853

ORCID: 0000-0002-8257-3595 (MTC), 0000-0002-8923-9733 (HY), 0000-0002-2562-2741 (MM), 0000-0002-3567-6911 (GM)

1
2
3
4
5
6
7
8
9
10
11

Running Title: SEM for local adaption 12

Keywords: Arabidopsis, structural equation modeling, Bayesian network, adaptation, factor analysis 13

Corresponding Author: 14

Malachy T. Campbell 15

417 Bradfield Hall 16

Cornell University 17

Ithaca, New York 14853 18

campbell.malachy@gmail.com 19

Abstract

Environmental association analyses (EAA) seek to identify genetic variants associated with local adaptation by regressing local environmental conditions at collection sites on genome-wide polymorphisms. The rationale is that environmental conditions impose selective pressure on trait(s), and these traits are regulated in part by variation at a genomic level. Here, we present an alternative multivariate genomic approach that can be utilized when both phenotypic and environmental data are available for the population. This framework utilizes Bayesian networks (BN) to elucidate interdependencies between local environmental conditions and empirical phenotypes, and jointly estimates the direct and indirect genetic covariances between empirical phenotypes and environmental conditions using a mixed-effects structural equation model (SEM). Direct genomic covariance between empirical phenotypes and environmental conditions may provide insight into whether QTL that affect adaptation to an environmental gradient also affects the observed phenotype. To demonstrate the utility of this approach, we leveraged two existing datasets consisting of 55 climate variables for 1,130 *Arabidopsis* accessions and empirical phenotypes for fitness and phenology collected on 515 accessions in two common garden locations in Europe. BN showed that plasticity for fitness and phenology was highly dependant on local environmental conditions. Moreover, genomic SEM revealed relatively high positive genomic correlation between plasticity in fitness and environmental variables that describe the favorability of the local environment for plant growth, indicating the presence of common QTL or independent QTL that are tightly linked. We believe the frameworks presented in this manuscript can provide new insights into the genetic basis of local adaptation.

Introduction

Identifying traits that confer adaptation to a given environment and elucidating the genetic determinants driving variation for these traits is an important goal for physiologists, evolutionary biologists, and quantitative geneticists. In many cases, particularly those working with agronomic species, these studies involve large-scale phenotypic evaluations in multiple environments, which are later integrated with genomic data using quantitative genetic frameworks. However, when the population is composed of individuals sampled across an environmental gradient, information regarding local environmental conditions at collection sites can be leveraged together with genomic data to identify genetic variants associated with variation for a given environmental factor (Fournier-Level et al., 2011; Blanquart et al., 2013; Yoder et al., 2014; Tiffin and Ross-Ibarra, 2014; Hoban et al., 2016). In recent years, a number of studies have employed similar approaches, termed environmental association analysis (EAA), to study the genetic basis of local adaptation (Fournier-Level et al., 2011; Yoder et al., 2014; Lasky et al., 2015).

EAA seeks to identify genomic variants that are associated with variation in environmental conditions at collection sites (Jones et al., 2013; Dell'Acqua et al., 2014; Yoder et al., 2014; Lasky et al., 2015; Anderson et al., 2016). The rationale for these approaches is that local environmental conditions impose selective pressure on some trait(s), and these traits are regulated in part by variation at a genomic level. Since adaptive traits should be correlated with local environmental conditions, regression of environmental variables on genome-wide single nucleotide polymorphisms (SNPs) may yield markers that are associated with environmental variables and, by proxy, adaptive traits. Several studies have leveraged these, and similar approaches, to elucidate the genetic basis of local adaptation (Fournier-Level et al., 2011; Yoder et al., 2014; Lasky et al., 2015). The only requirements for EAA is genomic data for a georeferenced population and environmental variables recorded at, or close to, collection sites.

Downstream analyses or independent studies are performed to determine if these variants have an effect on the phenotype, or whether they can be used to predict phenotypic variation. For instance, Yoder et al. (2014) utilized a population of 202 wild *Medicago truncatula* accessions to identify genomic associations with annual mean temperature, precipitation in the wettest month, and isothermality. They showed that accessions with a greater number of alleles associated with high precipitation in the wettest month also exhibited higher growth rate in a wet controlled environment. Similarly, Lasky et al. (2015) first identified environment-genotype associations in a panel of *Sorghum* landraces, and used these associations to predict agronomic characteristics in environments with contrasting moisture or edaphic conditions. Thus, these studies provide evidence that EAA can recover genetic determinate that are associated with environmental adaptation, and may influence phenotypic variation for adaptive or agronomically relevant traits.

However, when both phenotypic and environmental data are available for the population, alternative multivariate approaches can be utilized to jointly estimate genomic parameters and elucidate the genetic

interrelationships between local environmental conditions and observable phenotypes. With these approaches we can address whether there is a dependency between the empirical phenotype and the local environmental condition, effectively addressing the question “Is local adaptation to an environmental variable dependant on this trait?” and “What genes have an impact on both local adaptation *and* the empirical phenotype?” Structural equation models (SEM) are powerful frameworks that can be used to model the interdependencies between multiple variables (Wright, 1921; Haavelmo, 1943). When integrated into a quantitative genetics framework, these approaches allow quantitative genetic loci (QTL) or total genomic values to be decomposed into direct and indirect effects based on a predefined graphical model that describes directed relationships between variables (Gianola and Sorensen, 2004; Valente et al., 2013). SEM can be viewed as an extension of a conventional multi-trait (MT) quantitative genetic framework (Valente et al., 2013). Whereas with MT approaches, covariances among observable phenotypes are estimated and used to describe the symmetric linear relationships between variables, SEM extends the multivariate framework to allow recursive (effects from one phenotype affects the outcome of another) and simultaneous (reciprocal) structures among its variables by utilizing phenotypes as predictors for other phenotypes (Goldberger, 1972; Bielby and Hauser, 1977).

In quantitative genetics, SEM has been largely applied to topics in animal breeding and genetics. For instance in one of the first applications of SEM in quantitative genetics in the context of a linear mixed model, de los Campos et al. (2006b) used SEM to elucidate the interrelationship between milk yield and mastitis (inflammation of the udder quantified using somatic cell scores) in dairy cattle. The authors showed that models where milk yield was dependant on mastitis were better supported by the data, indicating that disease was the primary driver of reduced milk production rather than the converse. Since this work, quantitative genetic SEM frameworks have been used to elucidate the genetic interdependencies among meat quality traits, calving traits, fertility metrics, as well as milk yield and mastitis in other species or breeds (de los Campos et al., 2006a,b; Varona et al., 2007; Wu et al., 2007; König et al., 2008; Heringstad et al., 2009; de Maturana et al., 2009, 2010; Jamrozik et al., 2010; Peñagaricano et al., 2015a,b). More recently, the SEM quantitative genetic frameworks have been extended to perform genome-wide associations in chicken and rice (Momen et al., 2018, 2019). Given that many EAA studies assume a causal relationship between an unobserved phenotype and the local environment, SEM provides a framework where in these relationships can be explicitly encoded in the model – when empirical phenotypes are available for the same population. Moreover, these frameworks provide a means to examine the covariance in genetic effects that act directly on the empirical phenotype and the environmental variable (Valente et al., 2013).

Direct applications of quantitative genetic SEM frameworks to EAA is not trivial. For one, SEM requires a putative causal networks that describes the dependencies among and between environmental variables and empirical phenotypes (Gianola and Sorensen, 2004). In most cases, these networks are not

only unknown, but learning the structure may even be an objective of the study itself. Secondly, the environmental data often consist of dozens or hundreds of variables that are highly correlated (Ferrero-Serrano and Assmann, 2019; Lasky et al., 2015). Thus, prior to applying SEM to EAA we must reduce the dimensionality of the environmental data and determine an appropriate network structure. One popular approach for dimensional reduction is factor analysis (FA) (de Los Campos and Gianola, 2007). The underlying rationale for FA is that relationships among variables are due to some underlying unobserved process. The goal of FA is to define a reduced set of unobserved, latent variables that maximize the correlation among groups of related observed variables. In quantitative genetics, FA is routinely applied to multi-environmental trials and high-dimensional multi-trait applications (Kelly et al., 2007; Meyer, 2009; de Los Campos and Gianola, 2007; Runcie and Mukherjee, 2013; Yu et al., 2019). Thus, when applied to high dimensional environmental data, FA may yield a reduced set of underlying variables that capture major patterns of local environments. When the underlying causal structure is unknown, Bayesian network (BN) approaches can be utilized to elucidate the probabilistic dependencies among variables (Scutari, 2009; Scutari and Denis, 2014). These dependencies are expressed using a directed acyclic graph where each variable is depicted as a node and directed edges join dependant nodes. Although BN approaches learn dependencies from the data itself, these approaches can yield insightful information regarding the causal relationships among variables. Such approaches have been leveraged to understand the genetic interdependencies among complex traits and have been utilized to elucidate potential causal structures that can be used in SEM quantitative genetic frameworks (Valente et al., 2013; Yu et al., 2019; Momen et al., 2018, 2019). Thus, both FA and BN can be leveraged to reduce the dimensionality of local environmental variables and elucidate the relations between traits or latent factors.

The objective of this study is to demonstrate the utility of SEM quantitative genetic frameworks for studying the genetic interrelationships between local environmental conditions and empirical phenotypes associated with fitness and phenology. To this end, we utilized two publicly available data sets that describe environmental variables at collection sites for 1,130 diverse *Arabidopsis thaliana* accessions, and empirical phenotypes in two precipitation regimes at two common garden locations in Europe (Ferrero-Serrano and Assmann, 2019; Exposito-Alonso et al., 2019). Several studies have shown that adaptation is polygenic (Pritchard and Di Rienzo, 2010; Pritchard et al., 2010; Flood and Hancock, 2017). With this in mind, we sought to forego single marker inferences and instead predict total genomic values for each individual, which are the cumulative additive genetic value for a given phenotypic variable. We seek to decompose total genetic effects for these variables into direct and indirect effects, effectively allowing us to address the following questions: "Are genetic effects for empirical phenotypes dependant on the genetic drivers for adaptation to local environmental conditions (and vice versa)?" and "How much of the total genomic value for an empirical phenotype is due to genetic effects from upstream

phenotypic variables?".

145

Materials and Methods

146

Environmental variables

147

This study utilized a publicly available data set of local environmental conditions for 1,130 *Arabidopsis* accessions. The original data, compiled by Ferrero-Serrano and Assmann (2019), consisted of 205 environmental variables for 829 unique collection sites. Categorical variables were removed from the data set, as well as variables that had missing values in $\geq 20\%$ of the accessions. After this filtering, 139 climate variables remained. Prior to FA, we removed variables that showed high collinearity, as variables with very high correlation can interfere with factor analysis. In total, these quality control steps provided data for 55 environmental variables for 1,130 accessions.

148

149

150

151

152

153

154

Empirical Phenotypes

155

Since the objective of the study was to examine the genomic interrelationships between local climate conditions and phenotypic plasticity in contrasting environments, we sought a data set that provided phenotypes recorded in the same germplasm in contrasting and ecologically-relevant conditions. To this end, we used data from a recent study by Exposito-Alonso et al. (2019) in which 515 of the 1,130 accessions were phenotyped for fitness, germination time and flowering time in two locations within the natural range of *Arabidopsis thaliana* and two simulated precipitation regimes. The experimental design and collection of phenotypic data is explained in great detail by Exposito-Alonso et al. (2019). Briefly, the 515 accessions were grown in open-ended rain-out shelters in Tuebingen, Germany and Madrid, Spain. The open-ended design allows for the temperature and humidity conditions within the structure to be similar to the natural environment. Within each location the plants were grown in a split-plot design. Two simulated precipitation regimes, which were designed to mimic natural rainfall at Tuebingen and Madrid, were randomly assigned to each subplot. The interquartile range for soil water content (SWC) in the low-precipitation treatment was 11.38-22.51% with a median of 16.1% in Madrid and 10.76-20.09% with a median of 14.7% in Tuebingen. The interquartile range for the high precipitation regime was 20.73-29.02% with a median of 24.6% in Madrid, and 22.62-33.00% with a median of 27.8% in Tuebingen. Median midday photosynthetically active radiation (PAR) values inside the shelters were $45.7 \text{ mol}\cdot\text{m}^{-2}\cdot\text{day}^{-1}$ in Madrid and $30.9 \text{ mol}\cdot\text{m}^{-2}\cdot\text{day}^{-1}$ in Tuebingen. Temperatures outside the structures ranged from $5.34\text{-}12.39^\circ\text{C}$ with a median of 8.5°C in Madrid and $2.44\text{-}9.54^\circ\text{C}$ with a median of 5.6°C in Tuebingen. These ranges are very consistent with temperatures recorded in the structures (Exposito-Alonso et al., 2019).

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

We estimated the macroenvironmental sensitivity for each accession and each empirical phenotype that was recorded by Exposito-Alonso et al. (2019) using the Finlay-Wilkinson (FW) approach (Finlay and Wilkinson, 1963). FW essentially expresses the plasticity of an accession grown across multiple

environments as a function of the overall population performance in each environment. The FW model is given by

$$y_{ij} = \mu + g_i + E_j + h_i E_j + e_{ij}$$

where y_{ij} is the phenotype for accession i in environment j , μ is the overall mean, g_i is the main accession effect, E_j is the main environment effect, h_i is the slope for accession i on the overall environment means, and e_{ij} is the residual for accession i in environment j . Here, y_{ij} are best linear unbiased estimates for the accession effect in each environment from a model that accounts for systematic experimental effects (Exposito-Alonso et al., 2019). The FW model was fit using the FW package in R (Lian, 2014). The slope from this model was used as a metric for phenotypic plasticity in all downstream analysis.

Genotyping data

Imputed SNP markers were obtained for all 1,135 accessions from 1001genomes (https://1001genomes.org/data/GMI-MPI/releases/v3.1/SNP_matrix_imputed_hdf5/) (Weigel and Mott, 2009; Alonso-Blanco et al., 2016). We extracted marker information for the 1,130 accessions with climate data, and removed SNPs with low minor allele frequencies (MAF < 0.05). Moreover, SNPs in high linkage disequilibrium (LD) ($r > 0.85$) were removed using the PLINK indep function with a 50 SNP window, a step size of 5 SNPs, and a variance inflation factor (VIF) of 3.6. The VIF is computed as $\frac{1}{1-r^2}$. Thus, a VIF of 3.6 corresponds to a $r \approx 0.85$. After these filtering steps, 426,567 SNPs remained.

Factor analysis of environmental variables

To reduce the dimensionality of the 55 environmental variables, and define a reduced subset that captures potential undefined/unobserved variables that give rise to the original covariance, we utilized a combination of FA techniques, specifically exploratory and confirmatory factor analysis (EFA and CFA, respectively). Factor analysis seeks to identify a smaller set of latent variables that capture the underlying interrelationships between the original, manifest variables. The relationships between latent and manifest variables is given by

$$\mathbf{Y} = \mathbf{\Gamma}\mathbf{F} + \mathbf{s} \quad (1)$$

where \mathbf{Y} is an $t \times n$ matrix of phenotypes with $n = 515$ indicating the number of accessions and $t = 55$ indicating the number of traits; \mathbf{F} is an $l \times n$ matrix of factor scores that describe the values for each latent factor (l) for each accession; $\mathbf{\Gamma}$ is an $t \times l$ matrix that shows how each trait (t) loads onto each latent factor; and \mathbf{s} is a $t \times n$ matrix that represents the specific effects for each trait and accession. Thus, FA expresses a set of manifest variables as a function of common, latent factors.

While both EFA and CFA are based on a similar framework, EFA allows manifest variables to load onto multiple latent factors and CFA does not. Thus, EFA is most often used to determine the appropriate number of latent factors and examine how manifest variables load on to them, and CFA is used to test hypothesis regarding the relationships between manifest and latent factors and to estimate factor loading scores. We determined the appropriate number of factors using parallel analysis. Parallel analysis is a simulation-based method that was originally proposed by Horn (1965) to determine the optimal number of latent factors. Briefly, parallel analysis randomly simulates data sets with similar properties to the observed data and uses these data to extract eigenvalues. Scree plots are used to plot and compare eigenvalues from the simulated data and eigenvalues from the observed data. The optimal number of factors is determined as the maximum number of factors that have observed eigenvalues that are larger than eigenvalues from simulated data. Parallel analysis was performed using the `fa.parallel` function in the `psych` package R (Revelle, 2018). We used the minimum residual method with 1,000 iterations. Once the optimal number of factors was determined (11 latent factors), EFA was performed using the factor analysis function, `fa()`, with varimax rotation and the minimum residual method with 1,000 iterations.

CFA was used to estimate factor scores for each accession and latent environmental variable. Since CFA only allows manifest variables to load onto a single latent variable, we used EFA results to determine which latent factor had the largest absolute loading for each manifest variable. Although EFA identified 11 latent factors, one latent factor was omitted from CFA because all manifest variables that loaded onto this latent factor had higher loadings for other latent factors. CFA was fit using the `sem` package in R according to the loadings provided in Figure 1 (Fox et al., 2017). Factor scores were computed with the ‘regression’ method using the `fscores()` function in the `sem` package (Fox et al., 2017).

Structure learning using Bayesian network

We next sought to elucidate the genomic interrelationships between plasticity and latent factor scores from CFA for local environmental conditions following an approach described by Yu et al. (2019). To this end, we first predicted genomic values for each accession and trait using a Bayesian multi-trait model (MTM). The MTM is given by

$$\mathbf{Y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

where \mathbf{Y} is an $n \times t'$ matrix of phenotypes composed of factor scores for latent environmental factors and plasticity for empirical phenotypes ($t' = 13$), where $n = 515$ is the number of individuals and t' is the number of phenotypes (ten latent local environmental variables and three empirical phenotypes, $t' = 13$); \mathbf{X} and \mathbf{Z} are incidence matrices that relate phenotypes to vectors of systematic effects (\mathbf{b}) and additive genetic effects \mathbf{u} , respectively; and \mathbf{e} is the error term. Moreover, we assume $\mathbf{u} \sim N(0, \Sigma_{\mathbf{u}} \otimes \mathbf{G})$ and $\mathbf{e} \sim N(0, \Sigma_{\mathbf{e}} \otimes \mathbf{I}_{n \times n})$, where \mathbf{G} is a genomic relationship matrix constructed following VanRaden (2008),

$\Sigma_{\mathbf{u}}$ is a $t' \times t'$ covariance matrix for additive genetic effects. The MTM was fit using the MTM package in R with 10,000 Markov chain Monte-carlo (MCMC) samples of which the first 2,000 are discarded and every fifth sample was retained (de los Campos and Grüneberg, 2016).

Bayesian network (BN) learning approaches assume that the samples are independent. However, when predicting additive genomic values using MTM, dependencies are between breeding values for accessions are introduced from \mathbf{G} . Therefore prior to BN learning, we followed an approach described by Töpner et al. (2017) to remove dependencies. Briefly, \mathbf{G} was decomposed into Cholesky factors by $\mathbf{G} = \mathbf{L}\mathbf{L}'$, where \mathbf{L} is a lower triangle matrix with dimensions $n \times n$. We define a $nt' \times nt'$ matrix, \mathbf{M} via $\mathbf{M} = \mathbf{I}_{t' \times t'} \otimes \mathbf{L}$. By multiplying the nt' vector of genomic values (\mathbf{u}) by the inverse of \mathbf{M} , we are provided with a vector of transformed genomic values ($\mathbf{u}^* = \mathbf{M}^{-1}\mathbf{u}$) that follow a distribution given by $N(0, \Sigma_{\mathbf{g}} \otimes \mathbf{I}_{n \times n})$. Thus, the transformed genomic values are independent between accessions.

BN are a class of graphical models that represent the probabilistic dependencies between a set of random variables as a directed acyclic graph (\mathcal{G}) (Scutari and Denis, 2014). \mathcal{G} is composed of nodes (V) that represent random variables and edges (E) that depict probabilistic dependencies between nodes. BN follow the Markov property, which states that given its parents, a node is conditionally independent of all nodes that are non-descendants (Scutari and Denis, 2014). The joint probability distribution for k random variables ($X_V = (X_1, \dots, X_k)$) is given by

$$P(X_V) = P(X_1, \dots, X_k) = \prod_{V=1}^k P(X_V | \Pi_{X_V})$$

where parent nodes to X_v is indicated by Π_{X_v} (Scutari and Denis, 2014).

The vector of transformed genomic values (\mathbf{u}^*) was used as input for BN learning using the `bnlearn` package (Scutari, 2009). Structure learning was performed using four algorithms: hill-climbing (HC), tabu-search, max-min hill-climbing (MMHC), and general 2-phase restricted maximization (RSmax2). HC and tabu are score-based, greedy algorithms which seek to maximize the goodness-of-fit (i.e., network score). These algorithms begin with an empty network structure and add, remove, or reverse edge each edge until a maximum score is reached. The latter two algorithms, MMHC and RSmax2, are hybrid learning algorithms, which essentially restrict the score-based approach described above on a subset of nodes within the network (Tsamardinos et al., 2006). For each algorithm, we used a combination of bootstrapping and model averaging to identify robust networks and quantify uncertainty in linkages and the direction of each edge. Five hundred bootstrapping replicates were used and edges that were present in less than 85% of the networks were removed, and the models were averaged. We compared networks from each algorithm using the Bayesian information criteria (BIC) and selected the 'best' network according to the network that produced the highest BIC since `bnlearn` rescales BIC values by -2.

Genomic structural equation model

256

Work by Gianola and Sorensen (2004) provided a basis to introduce SEM into classical quantitative genetics frameworks. SEM utilize a system of linear equations to model the interrelationships between multiple dependant variables. Once introduced into the quantitative genetics frameworks pioneered by Henderson (1984), these approaches provide a means to partition multiple phenotypes into direct and indirect genetic components according to a predefined network structure (Gianola and Sorensen, 2004; Valente et al., 2013; Bello et al., 2018). In matrix form, the structural equation model is given by

$$\mathbf{Y} = \mathbf{\Lambda}\mathbf{Y} + \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

where all matrices are defined according to the MTM described above. However, note that the response variable \mathbf{Y} appears on both the right and left-hand side of the equation, meaning that some phenotypes will serve as covariates for other phenotypes. The effect of an upstream phenotype on a downstream phenotype is determined by the direction and magnitude of elements in the coefficient matrix ($\mathbf{\Lambda}$). $\mathbf{\Lambda}$ is typically a lower triangle matrix with zeros in the diagonal and upper triangle. We assume $\mathbf{u} \sim N(0, \mathbf{\Sigma}_{\mathbf{u}_0} \otimes \mathbf{G})$ and $\mathbf{e} \sim N(0, \mathbf{\Sigma}_{\mathbf{e}_0} \otimes \mathbf{I}_{n \times n})$, where $\mathbf{\Sigma}_{\mathbf{u}_0}$ and $\mathbf{\Sigma}_{\mathbf{e}_0}$ represent the genomic and residual covariances for total effects.

Given a simple, hypothetical causal structure for three phenotypes ($y_1 \rightarrow y_2$ and $y_1 \rightarrow y_3$), we can decompose each phenotype into genetic and non-genetic components using the following system of equations

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u}_1 + \mathbf{e}_1 \\ \mathbf{y}_2 &= \lambda_{y_1 \rightarrow y_2} \mathbf{y}_1 + \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u}_2 + \mathbf{e}_2 \\ \mathbf{y}_3 &= \lambda_{y_1 \rightarrow y_3} \mathbf{y}_1 + \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u}_3 + \mathbf{e}_3 \end{aligned}$$

Since y_1 has no variables leading to it, the total genomic effects for y_1 are given by $\mathbf{u}_{1_{\text{total}}} = \mathbf{u}_1$. For y_2 we have an indirect effect coming from y_1 , therefore the total genomic value is given by $\mathbf{u}_{2_{\text{total}}} = \lambda_{y_1 \rightarrow y_2} \mathbf{u}_1 + \mathbf{u}_2$. For y_3 , total genomic values are given by $\mathbf{u}_{3_{\text{total}}} = \lambda_{y_1 \rightarrow y_3} \mathbf{u}_1 + \mathbf{u}_3$. Solving the mixed model equation provides solutions for direct genomic values and estimates the genetic and residual (co)variances for direct effects among traits ($\mathbf{\Sigma}_{\mathbf{u}_0}$ and $\mathbf{\Sigma}_{\mathbf{e}_0}$, respectively). Covariances for total genomic and residual effects can be computed through a simple transformation on the appropriate covariances matrix for direct effects. The total genomic covariance is given by $\mathbf{\Sigma}_g = (\mathbf{I}_{t \times t} - \mathbf{\Lambda})^{-1} \mathbf{\Sigma}_{\mathbf{u}_0} (\mathbf{I}_{t \times t} - \mathbf{\Lambda})^{-1'}$. We fit SEM using the ten latent environmental variables and the plasticity measures for three empirical phenotypes according to the learned structure described above. The model was fit using the MTM package with 10,000 MCMC samples with the first 2,000 samples discarded and every fifth sample retained (de los Campos and Grüneberg, 2016).

Data availability

278

Local environmental variables were obtained from the Arabidopsis ClimTools repository
(<https://github.com/CLIMtools>) (Ferrero-Serrano and Assmann, 2019), and empirical phenotypes for
common garden locations were obtained from Exposito-Alonso et al. (2019). Scripts used for analyses of
these data are available Arabidopsis EFA repository
(<https://github.com/malachycampbell/ArabidopsisEFA>) and are documented to ensure
reproducibility. Supplemental figures and files are available at FigShare ().

279

280

281

282

283

284

Results

285

To examine the genomic relationship between local environments across the native range of *Arabidopsis thaliana* we utilized a publicly available panel of 1,135 diverse *Arabidopsis* accessions. These materials were collected from 829 non-redundant sites across Europe, Asia, Africa and N. America, and are discussed in detail by Ferrero-Serrano and Assmann (2019). The collection site for each accession is provided as Supplemental Figure S1. We utilized an existing dataset of 205 climatic, edaphic, and remote sensing variables to characterize the local environmental conditions at each of the collection sites. These variables describe precipitation, temperature, and vegetative productivity patterns, as well as soil physical and chemical characteristics (Ferrero-Serrano and Assmann, 2019).

286

287

288

289

290

291

292

293

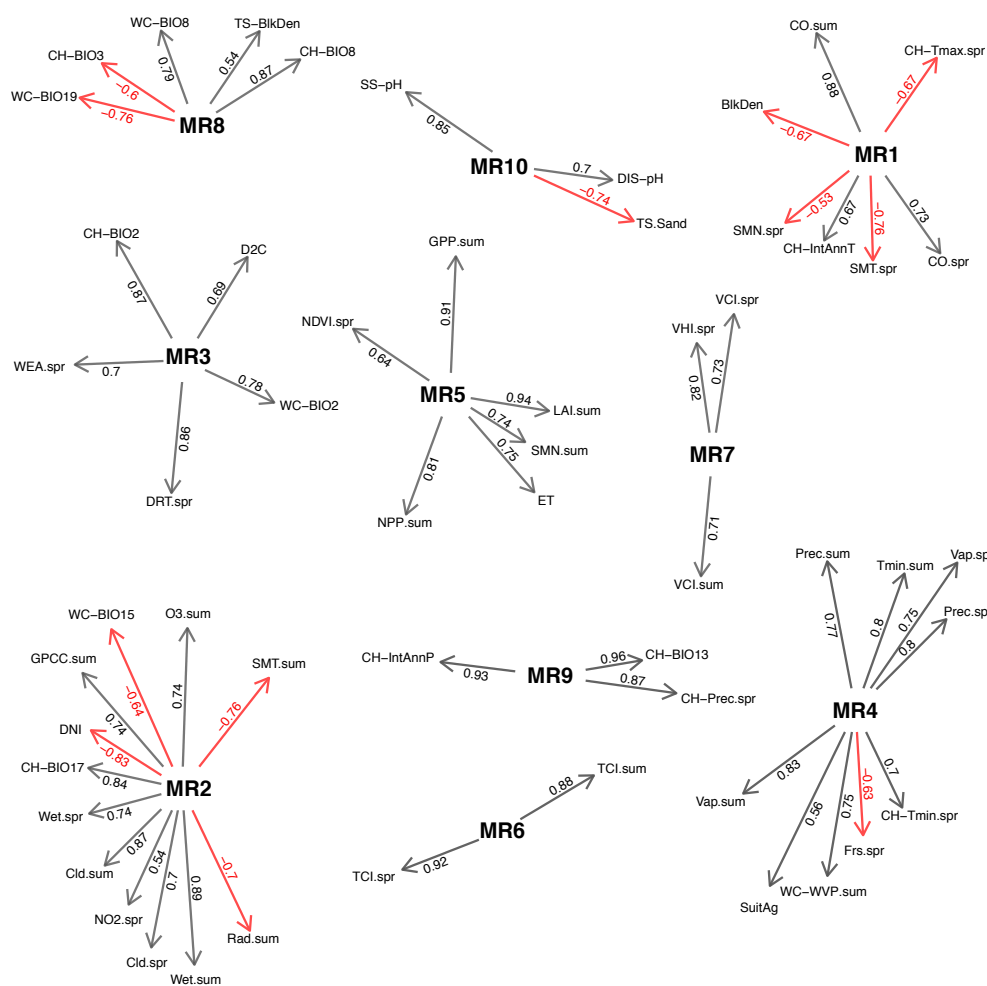


Figure 1. Factor loadings for manifest local environmental variables. Variables in bold type face are latent factors identified using factor analysis, while nodes emanating from these are manifest environmental variables. Edges colored in grey indicate the manifest variable has a positive loading on the latent factor, while those in red indicate negative loadings.

Factor analysis reveals the underlying structure of local environments

294

An initial inspection of the environmental variables showed a high degree of correlation between variables (Supplemental Fig. S2). Given the size of the data set, as well as the high degree of correlation between variables, we sought to reduce the 55 variables to a smaller set of factors that capture the underlying theoretical structure of the environments. To this end, we performed EFA on the set of 55 variables to explore the underlying structure of local environmental conditions and define a reduced set of variables that capture unobserved processes (latent factors) that drive these relationships. Confirmatory factor analysis was used to determine the contribution of each environmental variable to the latent factor and quantify how each accession contributed to each latent factor. EFA revealed that the 55 variables could be reduced to a set of 11 latent factors (Supplemental Fig. S3). Although 11 latent factors were defined, variables loading onto factor 11 had stronger loading on other latent factors. Thus this latent factor was omitted from downstream analysis. The loadings from EFA are provided as Supplemental File S1.

295

296

297

298

299

300

301

302

303

304

305

In theory, these latent factors should represent unobserved processes that give rise to the observed variables, and in the context of the current study, may describe processes that shape local environments. Factor loadings from CFA are shown in Figure 1. A complete listing of latent factors, the manifest variables that load onto them, and the interpretation of latent factors is provided in Supplemental File S2. Twelve environmental variables loaded onto the second latent factor (MR2). The manifest variables describe the frequency of wet days, cloud coverage, solar radiation, precipitation seasonality and precipitation of the driest quarter. Variables associated with precipitation and cloud cover largely showed positive contributions to MR2, while those associated with solar radiation showed negative contributions. Thus, MR2 is likely a description of how bright and dry an environment is. Three latent factors were defined which captured the favorability of local environments to plant growth. For instance, two metrics for vegetation condition index (VCI) which quantifies vegetation cover in a period of time to relative extremes and vegetative health index (VHI) that represents the favorability of the environment for vegetation activity showed positive loadings onto MR7. Moreover, the two manifest variables that represent temperature condition index (TCI) which loaded positively onto MR6. While MR6 and MR7 are largely associated with indices that describe the potential impact of environmental conditions on plant health, MR5 captures the productivity of the environment as manifest variables associated with gross primary productivity, evapotranspiration, normalized difference vegetation index, and net primary productivity were loaded onto this latent factor. Several other latent factors were identified that captured precipitation patterns at each local environment. For instance, the ninth latent factor (MR9) largely captures precipitation and precipitation variability between years. Environmental variables representing the amount of precipitation in the wettest month, precipitation in the spring, and interannual precipitation showed strong positive contributions to MR9.

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

Examining plasticity in fitness and phenology in contrasting environments

The ability of plants to exhibit plasticity in phenotypic traits is important strategy for adaptation to environmental constraints. With this in mind, we sought to elucidate the genetic interrelationships between plasticity in phenological traits and fitness, and local environmental characteristics. We utilized an existing dataset consisting of phenological (time to germination and flowering) traits and fitness recorded on 515 diverse *Arabidopsis* accessions grown in common garden experiments in Tuebingen and Madrid (Exposito-Alonso et al., 2019). At each common garden location, accessions were grown under simulated high and low rainfall conditions, with high rainfall conditions mimicking the natural precipitation in Tuebingen and low rainfall conditions mimicking the precipitation at Madrid (Exposito-Alonso et al., 2019).

The distribution of phenological and fitness traits at each precipitation-location combinations are shown in Figure 2. Significant differences between precipitation-location combinations were observed for fitness and flowering time ($p < 0.0001$). In general, accessions flowered later at Tuebingen compared to Madrid, while low precipitation seemed to delay flowering in both locations indicating that temperature and daylength differences between locations may be the largest driver of differences in flowering time between locations. In general, the accessions exhibited higher fitness in the two high-rainfall treatments compared to low rainfall treatments. Fitness was highest for the high rainfall treatment in Madrid, while the low precipitation treatment at Madrid showed the lowest average fitness. The environment in the high rainfall treatment at Madrid is characterized by simulated rainfall that is similar to the natural precipitation at the common garden location in Tuebingen. Thus, the ample water availability (27.8% SWC) combined with the warm temperatures (median temperature 8.5°C) in Madrid are highly favorable for growth and reproduction in *Arabidopsis*. However, when warm temperatures are combined with inadequate rainfall (16.1% SWC), the overall performance is reduced greatly, as observed for the low average fitness observed in low precipitation in Madrid (M_1).

To estimate environmental plasticity for fitness and phenological traits, we estimated reaction norms for each accession using the FW approach (Finlay and Wilkinson, 1963). Briefly, the FW approach expresses the plasticity for each individual grown across a range of environments as a function of the average population performance at each environment. For each individual, the slope of the linear model expresses the plasticity (or macroenvironmental sensitivity) with respect to average plasticity of the population. The plasticity for each accession with respect to mean performance at each environment is shown in Figure 2D-F.

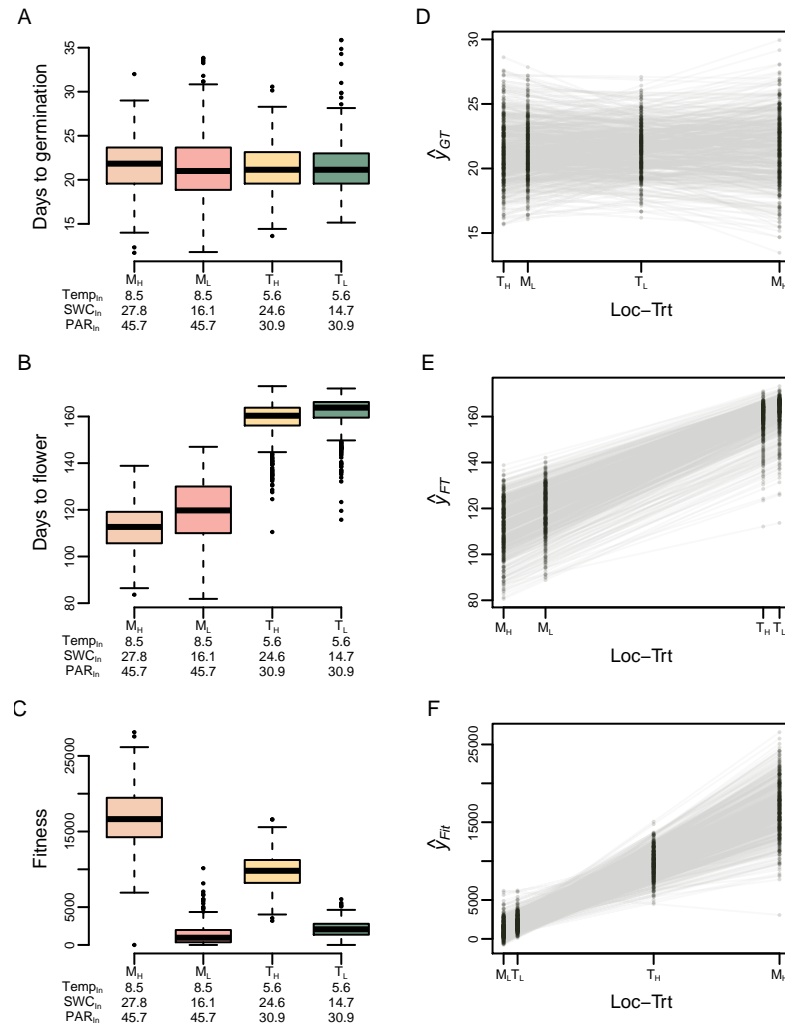


Figure 2. Distribution and plasticity of fitness and phenological traits across contrasting environments. (A-C) Distribution of adjusted means for fitness and phenological traits. Median values for environmental conditions within shelters are shown beneath each boxplot. Temp refers to median temperature in °C, SWC indicates soil water content, and PAR refers to photosynthetically active radiation ($\text{mol m}^{-2} \text{day}^{-1}$). The predicted phenotypic values (\hat{y}) of each accession in each location-treatment (Loc-Trt) combination is shown in panels D-F and were obtained using the FW approach. ‘M’ refers to common garden in Madrid and ‘T’ indicates common garden in Tuebingen, while the subscripts L and H refer to the low and high precipitation treatment, respectively.

Elucidating genetic dependencies between local environmental factors and fitness related traits

To elucidate the genetic interdependencies between local environmental conditions, and fitness and phenological plasticity, we inferred the potential causal genetic relationships between environmental

359

360

361

362

factors and observed phenotypes using four BN structure learning algorithms. Structure learning was performed using the ten latent environmental factors described above and reaction norm slopes for phenological traits and fitness, and the “best” structure was selected based on BIC scores. Of the four algorithms evaluated, the “best” network was given by tabu algorithm (Table 1). Since the primary objective of this study is to elucidate the relationships between local environmental conditions and empirical phenotypes, we focused interpretations of the network on relationships within the Markov blanket for plasticity traits (Figure 3).

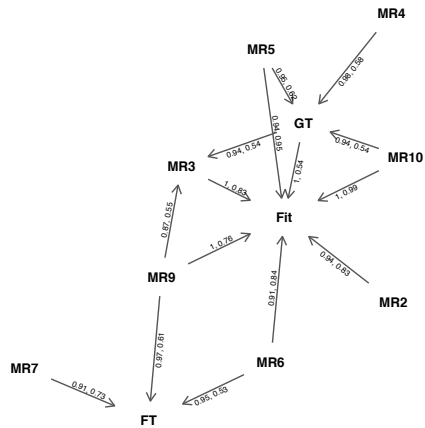
Table 1. Evaluation of four Bayesian structure learning algorithms. Bayesian network structures were learned using the ten latent environmental variables and plasticity for phenological traits (germination and flowering time) and fitness. The “best” network was selected based on the highest Bayesian information criteria (BIC) and Gaussian BIC values (gBIC). Algo.: algorithm; HC: hill-climbing; MMHC: min-max hill-climbing; RSmax2: general 2-phase restricted maximization

Algo.	gBIC	BIC
HC	-1963.02	-1963.02
tabu	-1962.58	-1962.58
MMHC	-2480.54	-2480.54
RSmax2	-2681.40	-2681.40

Although the learned structure is complex, several interesting features are apparent. First, of the 29 edges in the network, 41.4% (12 edges) describe relationships from environmental variables to empirical phenotypes, while only 3.45% (1 edge) describe relationships from plastic responses to environmental variables. These results suggest that genomic values for empirical phenotypes are highly dependant on genetic factors associated with adaptation to local environmental conditions. In addition, 51.7% edges (15 edges) were from environmental variables to other environmental variables, and only a single edge was from plastic responses to other plastic responses. Thus, genetic relationships between environmental variables or plastic responses are far more common than relationships from plastic responses to environmental variables.

In addition to overall topological features of the BN, several nodes were identified that were heavily influenced by other variables. For instance, plasticity in flowering time (FT) showed the largest number of indirect effects, suggesting that plasticity in flowering time is highly dependant on genetic effects from adaptation to local environments. A total of seven variables were leading to FT, while three were leading to both plasticity in germination time (GT) and fitness (Fit). Several variables were identified that had indirect effects on many variables. For instance, MR5 and MR9, which describe overall plant productivity, and precipitation and interannual precipitation variability, respectively, each showed indirect effects on four nodes.

A



B

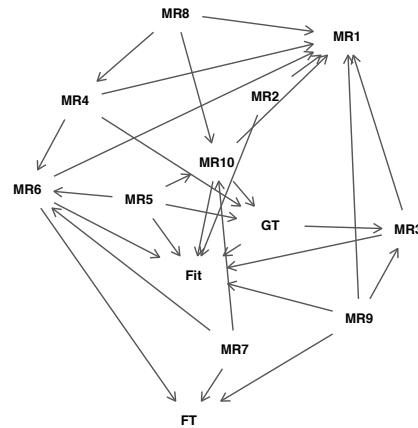


Figure 3. Visual depiction of probabilistic dependencies between environmental variables and empirical phenotypes. The network shown in panel A depicts the Markov blanket for empirical phenotypes and the full network is shown in panel B. A model averaging approach with 500 bootstrap samples was used to learn Bayesian network. The two numbers above each directed edge in panel A shows the proportion of bootstrap samples with the given edge and the proportion of samples with the given direction. The environmental variables are indicated with the “MR” prefix, while the empirical phenotypes are defined as follows: Fit: fitness plasticity; FT: flowering time plasticity; GT: germination time plasticity.

Structural equation modeling

The BN described above represents the probabilistic dependencies between plastic responses and local environmental conditions (Scutari and Denis, 2014). While this approach may provide insights into how variables act on one another, it does not tell *how much* of an effect one variable has on another. To estimate the magnitude of direct (QTL acting directly on focal trait) and indirect (QTL effects transmitted on focal trait by upstream trait) relationships among variables, we performed SEM using the learned structure described above. We leveraged this approach to decompose total genomic values for each environmental variable and empirical phenotype into direct and indirect effects, and examine the covariance between total genomic values and direct genomic values. The matrix of structural equation coefficients is shown in Table 2, and the genomic correlation matrix of direct and total effects is shown in Figure 4.

The utilization of plastic responses for phenological traits was motivated by several studies that suggest changes in an individual’s life cycle may be an important mechanism for adaptation to specific

387

388

389

390

391

392

393

394

395

396

397

398

399

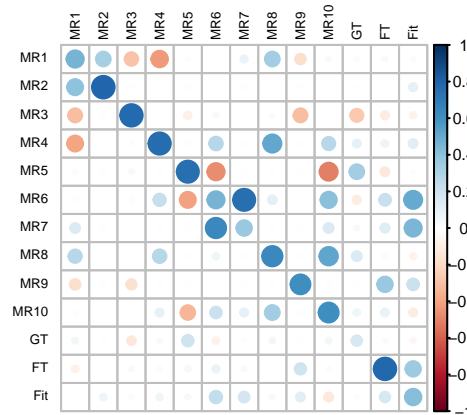


Figure 4. Genomic heritability and correlation for direct and indirect genetic effects The genomic heritability for total additive genetic effects (h^2) are shown in the diagonal. The upper triangle of the matrix shows the genomic correlation for total effects, while the lower triangle shows the genomic correlation for direct genetic values. Fit: fitness plasticity; FT: plasticity in flowering time; GT: plasticity in germination time

environmental constraints (Anderson et al., 2012; Vitasse et al., 2013; Augspurger, 2008; Chuine, 2010). 400
 While total genomic covariances provide insight into the relationships between total genetic values for 401
 two phenotypes, examination of the direct genomic covariances between traits may be more important in 402
 the context of the current study, as the covariance of direct genomic effects is driven by QTL that have 403
 an effect on both environmental adaptation and plasticity or QTL that affect each trait independently 404
 but are in tight LD (Valente et al., 2013). For direct genomic effects, the strongest positive genomic 405
 correlation between plastic responses and environmental variables was observed for Fit and MR6 406
 ($r_{g_{direct}} = 0.24$), which is a composite of temperature conditioning indices with lower values indicate a 407
 potential for high temperature stress on vegetative biomass. Fit also showed positive direct genomic 408
 correlation with MR7 ($r_{g_{direct}} = 0.18$), a variable composed of indices quantifying plant health, and MR9 409
 ($r_{g_{direct}} = 0.13$), which quantifies precipitation and interannual variability in precipitation. Collectively, 410
 these results indicate that the accessions that harbor alleles for reduced sensitivity of fitness to 411
 temperature gradients likely also harbor alleles associated with adaptation to warm, low rainfall 412
 environments. 413

In addition to Fit, relatively strong positive direct genomic correlation was observed between FT and 414
 MR9 ($r_{g_{direct}} = 0.20$), as well as GT and MR5 ($r_{g_{direct}} = 0.21$). However, the slope for FT largely 415
 represents the sensitivity of flowering time to differences in photoperiod and/or temperature for an 416
 accession, with lower values indicating more similar flowering times between common garden locations. 417

Table 2. Structural coefficients estimated using structural equation modeling. Path coefficients for network structure pictured in Figure 3 was estimated using a structural equation modeling approach. The columns indicate upstream nodes, while the rows indicate downstream nodes. Elements with ‘-’ indicate pairs of nodes that are not linked by an edge. Coefficient matrices for structures learned using a Bayesian network approach are typically have zero elements in the diagonal and upper triangle, however the coefficient matrix below has been reordered so that environmental variables are grouped and ordered by name. Fit: fitness plasticity; FT: plasticity in flowering time; GT: plasticity in germination time. Variables with the ‘MR’ prefix indicate latent environmental variables.

	MR1	MR2	MR3	MR4	MR5	MR6	MR7	MR8	MR9	MR10	GT	FT	Fit
MR1	-	0.35	-0.31	-0.47	-	0.13	-	0.6	-0.23	-0.18	-	-	-
MR2	-	-	-	-	-	-	-	-	-	-	-	-	-
MR3	-	-	-	-	-	-	-	-	-0.17	-	-0.2	-	-
MR4	-	-	-	-	-	-	-	0.38	-	-	-	-	-
MR5	-	-	-	-	-	-	-	-	-	-	-	-	-
MR6	-	-	-	0.23	-0.43	-	0.81	-	-	-	-	-	-
MR7	-	-	-	-	-	-	-	-	-	-	-	-	-
MR8	-	-	-	-	-	-	-	-	-	-	-	-	-
MR9	-	-	-	-	-	-	-	-	-	-	-	-	-
MR10	-	-	-	-	-0.33	-	0.12	0.41	-	-	-	-	-
GT	-	-	-	0.01	0.18	-	-	-	-	0.1	-	-	-
FT	-	-	-	-	-	0.13	-0.06	-	0.31	-	-	-	-
Fit	-	0.1	-0.01	-	0.13	0.34	-	-	0.13	-0.19	-	0.11	-

Therefore, it is unclear whether the non-zero direct genomic covariance between these variables indicates a common mechanism, or potential confounding of photoperiod insensitive accessions originating from more southern locations.

418
419
420

Discussion

Environment association analyses have become popular approaches to elucidate the genetic basis of local adaptation in the absence of fitness measurements in multilocation common garden trials (Fournier-Level et al., 2011; Yoder et al., 2014; Lasky et al., 2015). The aim of EAA is to identify genes or loci that may impact traits that confer fitness along an environmental gradient. However, when fitness is measured in multiple common garden locations along an environmental gradient, the change in fitness as a function of mean population performance provides a single metric that describes the impact of the environment on fitness. Moreover, when this metric is introduced as the response variable in genome-wide association studies, strong associations indicate the presence of gene(s) that may influence fitness along the environmental gradient.

In the current study, we seek to integrate both data types in the SEM framework to examine the genetic interdependencies and covariances between changes in fitness and phenology in multiple environments and local environmental conditions. However, whereas most EAA estimate the effects of individual loci, we predict the total genetic values (i.e. the summation of QTL effects for a given genotype) for each variable. Thus, in cases where collection sites and common garden locations follow the same gradients, we expect covariance in genetic signals that impact both variables directly. Consistent with this expectation, we observed non-zero genetic covariance between local environmental conditions and changes in fitness between common garden locations. For instance, Fit showed positive correlation of direct genetic effects for MR6, as well as MR7. The latent variables MR6 and MR7, capture the favorably of the local environment for plant growth. Thus, higher values indicate environments that have favorable conditions for plant growth and, on a whole, are highly productive. Moreover, Fit describes the changes in fitness driven largely by water availability, with higher values indicating greater fitness in high-rainfall treatment in Madrid and low values indicating low fitness in low-rainfall treatment in Madrid (Figure 2). Thus, the positive genomic correlation of direct effects indicates that accessions harboring alleles for high fitness in simulated, high-productivity environments will also tend to harbor alleles associated with adaptation to highly productive local environments. Although weaker than the direct genomic correlation for MR6 and MR7, Fit also showed positive genomic correlation with a latent environmental variable that largely captured precipitation and precipitation variability of the local environment, with higher values indicating higher precipitation ($MR9; r_{direct} = 0.13$). Collectively, these results indicate that fitness in response to some local environmental conditions may be regulated common genetic mechanisms that affect fitness in simulated environments. However in either case (e.g., local environment associations or common garden fitness), the traits that impact fitness are largely unknown.

Phenotypic plasticity is an important process that allow plants to quickly modify physiology, morphology, or phenology in response to changes in the environment (Bradshaw, 1965). Individuals that exhibit greater plasticity may be better positioned to respond to new environmental constraints, as novel

phenotypes brought on by environmental change may provide persistence in the short-term 456
(West-Eberhard, 2005; Matesanz et al., 2010; Nicotra et al., 2010; Valladares et al., 2014). However, 457
phenotypic plasticity is not always advantageous (DeWitt et al., 1998; Ghalambor et al., 2007). For 458
instance, Scheepens and Stöcklin (2013) showed that increased temperature leads to early flowering, but 459
reduced seed set in *Campanula thyrsoides*. Thus, it is important to couple observations of plasticity 460
across an environmental gradient with measurements of fitness in the same environments to determine 461
whether phenotypic plasticity can be a mechanism underlying fitness. Here, we utilized measures of 462
fitness and empirical phenotypes in four environments. Correlations for reaction norms for fitness and 463
phenological traits showed a significant, albeit weak, correlations between Fit and GT ($r = -0.09$, $p =$ 464
 0.417) and Fit and FT ($r = 0.18$, $p < 0.0001$), indicating that these changes in fitness are associated 465
with changes in phenology. In the case of FT, the positive correlation indicates that accessions that show 466
greater plasticity in flowering time (positive slope for FT meaning delayed flowering in Germany relative 467
to Madrid) tend to exhibit greater fitness in high-precipitation regimes relative to low-precipitation 468
regimes (i.e., positive slope for Fit). Correlation provides a simple means to measure the relationships 469
between two traits within a population. However, a non-zero correlation does not necessarily indicate 470
that the outcome/expression for one characteristic is dependant on another. BN approaches on the other 471
hand, have been developed to elucidate probabilistic dependencies among a group of interrelated 472
variables (Pearl, 2014; Scutari and Denis, 2014). The BN shown in Figure 3 shows an directed edge from 473
GT to Fit, indicating that changes in fitness across the common garden environments is dependant on 474
changes in germination time. However, no edges were found between FT and Fit, indicating that 475
although these two characteristics covary, changes in Fit may not be dependant on changes in FT. 476

Although it is seemingly a natural tendency to view these dependencies as causal relationships, it is 477
important not to over interpret results from BN. While BN are a powerful approach to assess the 478
interdependencies between variables, structure learning with BN imposes several constraints that may 479
limit its applications in biology. One major limitation is that BN do not allow feedback loops or cyclical 480
relationships in the structure, which are pervasive throughout biology especially at a molecular level 481
(Scutari and Denis, 2014). Thus, if the underlying causal relationships between traits involves feedback 482
loops, the structure learned with BN will likely be inaccurate (Valente et al., 2013). Thus, the network 483
might reflect highly probable relationships between variables, but may not represent the true causal 484
relationships that give rise to the data. Secondly, in the current study, BN were constructed using a 485
mixture of observational and experimental data. In the absence of randomization, dependencies observed 486
in Bayesian networks constructed using observational data may be driven by unobserved confounders, 487
thereby making causal claims based on the data problematic (Bello et al., 2018, see for review). 488
Nevertheless, causal relationships can be learned from the data and should be used to generate 489
hypothesis for further studies. In our study, BN revealed dependencies between plasticity in fitness and 490

several environmental variables. Fitness in a given environment is largely the consequence of a trait or 491
traits that confer adaptation to a set of environmental conditions. In other words, fitness is not a 492
mechanism for local adaptation, but rather is a measure of adaptation. Thus, we expect that fitness in a 493
given location/precipitation regime should be highly dependant on mechanisms that were selected by 494
environmental pressures in the accessions' local environments, and this expectation is largely confirmed 495
by the network learned from the data (Figure 3). However, covariance in direct effects for other variables, 496
such as between FT and MR9, may not be so easy to explain. The latent environmental variable MR9 497
largely captures precipitation and precipitation variability, as the manifest variables spring precipitation, 498
precipitation of the wettest month, and interannual precipitation variability load onto MR9. The positive 499
direct covariance between MR9 and FT suggest that accessions that harbor alleles associated with 500
adaptation to environments with high precipitation will also tend to harbor alleles associated with higher 501
plasticity in flowering time. However, plasticity in flowering time is largely driven by differences in day 502
length and temperature between common garden locations rather than by precipitation regimes (Figure 503
2). Thus, it is questionable whether the direct genomic covariance is due to QTL that affect adaptation 504
to precipitation gradients and the sensitivity of flowering time to photoperiod and/or temperature, or if 505
this is driven by unaccounted, confounding effects within the data. Projection of phenotypic values for 506
FT on collection sites show clusters of accessions originating from the Northern Iberian peninsula and 507
Southern Sweden with low plasticity for flowering time (Supplemental Figure S4). Moreover, these 508
regions also exhibit low values for MR9. Further studies or alternative experimental designs are necessary 509
to determine whether this covariance is due to common effects on adaptation to precipitation gradients 510
and plasticity in FT, or are due to sampling bias. Thus, while BN can provide important insight into the 511
interrelationships between traits, when these networks are constructed using observational data we 512
should view these results with caution rather than to discount inferred relationships as spurious. 513

While BN describe the probabilistic dependencies among variables, they only provide insight into the 514
structure of relationships in the data. In many cases, we are interested in understanding how genetic 515
effects for an upstream trait affect the outcome of a downstream trait. SEM provides a means to 516
estimate path coefficients according to a predefined network structure, as well as partition phenotypic 517
values into genetic values that affect a trait directly (i.e., direct genetic values) and genetic values that 518
are due to genetic effects acting directly on upstream variables (Gianola and Sorensen, 2004; Valente 519
et al., 2013). In some sense, estimates of the structural coefficients may seem like the most attractive 520
component of SEM, as these describe how intervention on an upstream variable (e.g., a latent 521
environmental variable) will impact the outcome of the downstream variable (e.g., empirical phenotype) 522
given the direct effects for the downstream variable remain unchanged (Gianola and Sorensen, 2004). 523
However in the current study, we have data that is a combination of latent environmental variables and 524
empirical phenotypes. Thus, a more biologically meaningful question is whether QTL that have a direct 525

effect on adaptation to an environmental gradient also have a direct impact on some observable 526
phenotype. Non-zero covariance in direct effects between local environmental conditions indicates the 527
presence of common QTL, or independent QTL that are tightly linked (Valente et al., 2013). Thus, 528
identification of such QTL can provide important insights into the common mechanisms that impact 529
adaptation to local environments and plasticity. 530

Acknowledgements

531

Funding for this research was provided by the National Science Foundation (United States) through
Award No. 1736192 to GM.

532

533

Supplemental Materials

534

- **Supplemental Figure S1.** Geographic locations for all 1,035 *Arabidopsis* accessions. 535
- **Supplemental Figure S2.** Heatmap for 55 manifest environmental variables. 536
- **Supplemental Figure S3.** Projection of phenotypic values for flowering time plasticity on
collection sites for 515 accessions. 537
538
- **Supplemental File S1.** Factor loading from exploratory factor analysis. 539
- **Supplemental File S2.** Description of latent and manifest variables, and their loadings from
confirmatory factor analysis. 540
541

References

- Alonso-Blanco, C., Andrade, J., Becker, C., Bemm, F., Bergelson, J., Borgwardt, K. M., Cao, J., Chae, E., Dezaan, T. M., Ding, W., et al. (2016). 1,135 genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana*. *Cell*, 166(2):481–491.
- Anderson, J. E., Kono, T. J., Stupar, R. M., Kantar, M. B., and Morrell, P. L. (2016). Environmental association analyses identify candidates for abiotic stress tolerance in glycine soja, the wild progenitor of cultivated soybeans. *G3: Genes, Genomes, Genetics*, 6(4):835–843.
- Anderson, J. T., Inouye, D. W., McKinney, A. M., Colautti, R. I., and Mitchell-Olds, T. (2012). Phenotypic plasticity and adaptive evolution contribute to advancing flowering phenology in response to climate change. *Proceedings of the Royal Society B: Biological Sciences*, 279(1743):3843–3852.
- Augspurger, C. K. (2008). Early spring leaf out enhances growth and survival of saplings in a temperate deciduous forest. *Oecologia*, 156(2):281–286.
- Bello, N. M., Ferreira, V. C., Gianola, D., and Rosa, G. J. (2018). Conceptual framework for investigating causal effects from observational data in livestock. *Journal of Animal Science*, 96(10):4045–4062.
- Bielby, W. T. and Hauser, R. M. (1977). Structural equation models. *Annual Review of Sociology*, 3(1):137–161.
- Blanquart, F., Kaltz, O., Nuismer, S. L., and Gandon, S. (2013). A practical guide to measuring local adaptation. *Ecology Letters*, 16(9):1195–1205.
- Bradshaw, A. D. (1965). Evolutionary significance of phenotypic plasticity in plants. In *Advances in genetics*, volume 13, pages 115–155. Elsevier.
- Chaine, I. (2010). Why does phenology drive species distribution? *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365(1555):3149–3160.
- de Los Campos, G. and Gianola, D. (2007). Factor analysis models for structuring covariance matrices of additive genetic effects: a bayesian implementation. *Genetics Selection Evolution*, 39(5):481.
- de los Campos, G., Gianola, D., Boettcher, P., and Moroni, P. (2006a). A structural equation model for describing relationships between somatic cell score and milk yield in dairy goats. *Journal of Animal Science*, 84(11):2934–2941.
- de los Campos, G., Gianola, D., and Heringstad, B. (2006b). A structural equation model for describing relationships between somatic cell score and milk yield in first-lactation dairy cows. *Journal of Dairy Science*, 89(11):4445–4455.

-
- de los Campos, G. and Grüneberg, A. (2016). Mtm (multiple-trait model) package [www document]. 573
URL <http://quantgen.github.io/MTM/vignette.html> (accessed 10.25. 17). 574
- de Maturana, E. L., de los Campos, G., Wu, X.-L., Gianola, D., Weigel, K. A., and Rosa, G. J. (2010). 575
Modeling relationships between calving traits: a comparison between standard and recursive mixed 576
models. *Genetics Selection Evolution*, 42(1):1. 577
- de Maturana, E. L., Wu, X.-L., Gianola, D., Weigel, K. A., and Rosa, G. J. (2009). Exploring biological 578
relationships between calving traits in primiparous cattle with a bayesian recursive model. *Genetics*, 579
181(1):277–287. 580
- Dell’Acqua, M., Zuccolo, A., Tuna, M., Gianfranceschi, L., and Pè, M. E. (2014). Targeting 581
environmental adaptation in the monocot model brachypodium distachyon: a multi-faceted approach. 582
BMC Genomics, 15(1):801. 583
- DeWitt, T. J., Sih, A., and Wilson, D. S. (1998). Costs and limits of phenotypic plasticity. *Trends in* 584
Ecology & Evolution, 13(2):77–81. 585
- Exposito-Alonso, M., Burbano, H. A., Bossdorf, O., Nielsen, R., and Weigel, D. (2019). Natural selection 586
on the arabidopsis thaliana genome in present and future climates. *Nature*, pages 1–5. 587
- Ferrero-Serrano, Á. and Assmann, S. M. (2019). Phenotypic and genome-wide association with the local 588
environment of arabidopsis. *Nature Ecology & Evolution*, 3:274–285. 589
- Finlay, K. and Wilkinson, G. (1963). The analysis of adaptation in a plant-breeding programme. 590
Australian Journal of Agricultural Research, 14(6):742–754. 591
- Flood, P. J. and Hancock, A. M. (2017). The genomic basis of adaptation in plants. *Current Opinion in* 592
Plant Biology, 36:88–94. 593
- Fournier-Level, A., Korte, A., Cooper, M. D., Nordborg, M., Schmitt, J., and Wilczek, A. M. (2011). A 594
map of local adaptation in arabidopsis thaliana. *Science*, 334(6052):86–89. 595
- Fox, J., Nie, Z., and Byrnes, J. (2017). *sem: Structural Equation Models*. R package version 3.1-9. 596
- Ghalambor, C. K., McKay, J. K., Carroll, S. P., and Reznick, D. N. (2007). Adaptive versus 597
non-adaptive phenotypic plasticity and the potential for contemporary adaptation in new 598
environments. *Functional Ecology*, 21(3):394–407. 599
- Gianola, D. and Sorensen, D. (2004). Quantitative genetic models for describing simultaneous and 600
recursive relationships between phenotypes. *Genetics*, 167(3):1407–1424. 601
-

-
- Goldberger, A. S. (1972). Structural equation methods in the social sciences. *Econometrica: Journal of the Econometric Society*, pages 979–1001. 602 603
- Haavelmo, T. (1943). The statistical implications of a system of simultaneous equations. *Econometrica*, 11(1):1–12. 604 605
- Henderson, C. (1984). *Applications of linear models in animal breeding*. 606
- Heringstad, B., Wu, X.-L., and Gianola, D. (2009). Inferring relationships between health and fertility in norwegian red cows using recursive models. *Journal of Dairy Science*, 92(4):1778–1784. 607 608
- Hoban, S., Kelley, J. L., Lotterhos, K. E., Antolin, M. F., Bradburd, G., Lowry, D. B., Poss, M. L., Reed, L. K., Storfer, A., and Whitlock, M. C. (2016). Finding the genomic basis of local adaptation: pitfalls, practical solutions, and future directions. *The American Naturalist*, 188(4):379–397. 609 610 611
- Horn, J. L. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, 30(2):179–185. 612 613
- Jamrozik, J., Bohmanova, J., and Schaeffer, L. (2010). Relationships between milk yield and somatic cell score in canadian holsteins from simultaneous and recursive random regression models. *Journal of Dairy Science*, 93(3):1216–1233. 614 615 616
- Jones, M. R., Forester, B. R., Teufel, A. I., Adams, R. V., Anstett, D. N., Goodrich, B. A., Landguth, E. L., Joost, S., and Manel, S. (2013). Integrating landscape genomics and spatially explicit approaches to detect loci under selection in clinal populations. *Evolution*, 67(12):3455–3468. 617 618 619
- Kelly, A. M., Smith, A. B., Eccleston, J. A., and Cullis, B. R. (2007). The accuracy of varietal selection using factor analytic models for multi-environment plant breeding trials. *Crop Science*, 47(3):1063–1070. 620 621 622
- König, S., Wu, X., Gianola, D., Heringstad, B., and Simianer, H. (2008). Exploration of relationships between claw disorders and milk yield in holstein cows via recursive linear and threshold models. *Journal of Dairy Science*, 91(1):395–406. 623 624 625
- Lasky, J. R., Upadhyaya, H. D., Ramu, P., Deshpande, S., Hash, C. T., Bonnette, J., Juenger, T. E., Hyma, K., Acharya, C., Mitchell, S. E., et al. (2015). Genome-environment associations in sorghum landraces predict adaptive traits. *Science Advances*, 1(6):e1400218. 626 627 628
- Lian, L. (2014). *FW: Performs Gibbs Sampler and Least Square models for Finlay-Wilkinson regressions*. R package version 0.0. 629 630
- Matesanz, S., Gianoli, E., and Valladares, F. (2010). Global change and the evolution of phenotypic plasticity in plants. *Annals of the New York Academy of Sciences*, 1206(1):35–55. 631 632
-

-
- Meyer, K. (2009). Factor-analytic models for genotype \times environment type problems and structured covariance matrices. *Genetics Selection Evolution*, 41(1):21. 633
634
- Momen, M., Ayatollahi Mehrgardi, A., Amiri Roudbar, M., Kranis, A., Mercuri Pinto, R., Valente, B., Morota, G., Rosa, G. J., and Gianola, D. (2018). Including phenotypic causal networks in genome-wide association studies using mixed effects structural equation models. *Frontiers in Genetics*, 9:455. 635
636
- Momen, M., Campbell, M. T., Walia, H., and Morota, G. (2019). Utilizing trait networks and structural equation models as tools to interpret multi-trait genome-wide association studies. *Plant Methods*, 15(1):107. 638
639
640
- Nicotra, A. B., Atkin, O. K., Bonser, S. P., Davidson, A. M., Finnegan, E., Mathesius, U., Poot, P., Purugganan, M. D., Richards, C. L., Valladares, F., et al. (2010). Plant phenotypic plasticity in a changing climate. *Trends in Plant Science*, 15(12):684–692. 641
642
643
- Pearl, J. (2014). *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier. 644
- Peñagaricano, F., Valente, B., Steibel, J., Bates, R., Ernst, C., Khatib, H., and Rosa, G. (2015a). Searching for causal networks involving latent variables in complex traits: Application to growth, carcass, and meat quality traits in pigs. *Journal of Animal Science*, 93(10):4617–4623. 645
646
647
- Peñagaricano, F., Valente, B. D., Steibel, J. P., Bates, R. O., Ernst, C. W., Khatib, H., and Rosa, G. J. (2015b). Exploring causal networks underlying fat deposition and muscularity in pigs through the integration of phenotypic, genotypic and transcriptomic data. *BMC Systems Biology*, 9(1):58. 648
649
650
- Pritchard, J. K. and Di Rienzo, A. (2010). Adaptation—not by sweeps alone. *Nature Reviews Genetics*, 11(10):665. 651
652
- Pritchard, J. K., Pickrell, J. K., and Coop, G. (2010). The genetics of human adaptation: hard sweeps, soft sweeps, and polygenic adaptation. *Current Biology*, 20(4):R208–R215. 653
654
- Revelle, W. (2018). *psych: Procedures for Psychological, Psychometric, and Personality Research*. Northwestern University, Evanston, Illinois. R package version 1.8.12. 655
656
- Runcie, D. E. and Mukherjee, S. (2013). Dissecting high-dimensional phenotypes with bayesian sparse factor analysis of genetic covariance matrices. *Genetics*, 194(3):753–767. 657
658
- Scheepens, J. and Stöcklin, J. (2013). Flowering phenology and reproductive fitness along a mountain slope: maladaptive responses to transplantation to a warmer climate in *campanula thyrsoides*. *Oecologia*, 171(3):679–691. 659
660
661
- Scutari, M. (2009). Learning bayesian networks with the bnlearn r package. *arXiv preprint arXiv:0908.3817*. 662
663
-

-
- Scutari, M. and Denis, J.-B. (2014). *Bayesian networks: with examples in R*. Chapman and Hall/CRC, 664
Florida. 665
- Tiffin, P. and Ross-Ibarra, J. (2014). Advances and limits of using population genetics to understand 666
local adaptation. *Trends in Ecology & Evolution*, 29(12):673–680. 667
- Töpner, K., Rosa, G. J., Gianola, D., and Schön, C.-C. (2017). Bayesian networks illustrate genomic and 668
residual trait connections in maize (*zea mays* l.). *G3: Genes, Genomes, Genetics*, 7(8):2779–2789. 669
- Tsamardinos, I., Brown, L. E., and Aliferis, C. F. (2006). The max-min hill-climbing bayesian network 670
structure learning algorithm. *Machine Learning*, 65(1):31–78. 671
- Valente, B. D., Rosa, G. J., Gianola, D., Wu, X.-L., and Weigel, K. (2013). Is structural equation 672
modeling advantageous for the genetic improvement of multiple traits? *Genetics*, 194(3):561–572. 673
- Valladares, F., Matesanz, S., Guilhaumon, F., Araújo, M. B., Balaguer, L., Benito-Garzón, M., Cornwell, 674
W., Gianoli, E., van Kleunen, M., Naya, D. E., et al. (2014). The effects of phenotypic plasticity and 675
local adaptation on forecasts of species range shifts under climate change. *Ecology Letters*, 676
17(11):1351–1364. 677
- VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *Journal of Dairy Science*, 678
91(11):4414–4423. 679
- Varona, L., Sorensen, D., and Thompson, R. (2007). Analysis of litter size and average litter weight in 680
pigs using a recursive model. *Genetics*, 177(3):1791–1799. 681
- Vitasse, Y., Hoch, G., Randin, C. F., Lenz, A., Kollas, C., Scheepens, J., and Körner, C. (2013). 682
Elevational adaptation and plasticity in seedling phenology of temperate deciduous tree species. 683
Oecologia, 171(3):663–678. 684
- Weigel, D. and Mott, R. (2009). The 1001 genomes project for arabidopsis thaliana. *Genome Biology*, 685
10(5):107. 686
- West-Eberhard, M. J. (2005). Developmental plasticity and the origin of species differences. *Proceedings 687
of the National Academy of Sciences*, 102(suppl 1):6543–6549. 688
- Wright, S. (1921). Correlation and causation. *J. Agric. Res.*, 20:557–580. 689
- Wu, X.-L., Heringstad, B., Chang, Y.-M., De los Campos, G., and Gianola, D. (2007). Inferring 690
relationships between somatic cell score and milk yield using simultaneous and recursive models. 691
Journal of Dairy Science, 90(7):3508–3521. 692
-

Yoder, J. B., Stanton-Geddes, J., Zhou, P., Briskine, R., Young, N. D., and Tiffin, P. (2014). Genomic signature of adaptation to climate in medicago truncatula. *Genetics*, 196(4):1263–1275. 693
694

Yu, H., Campbell, M. T., Zhang, Q., Walia, H., and Morota, G. (2019). Genomic bayesian confirmatory factor analysis and bayesian network to characterize a wide spectrum of rice phenotypes. *G3: Genes, Genomes, Genetics*, 9(6):1975–1986. 695
696
697

Supplemental Data

698

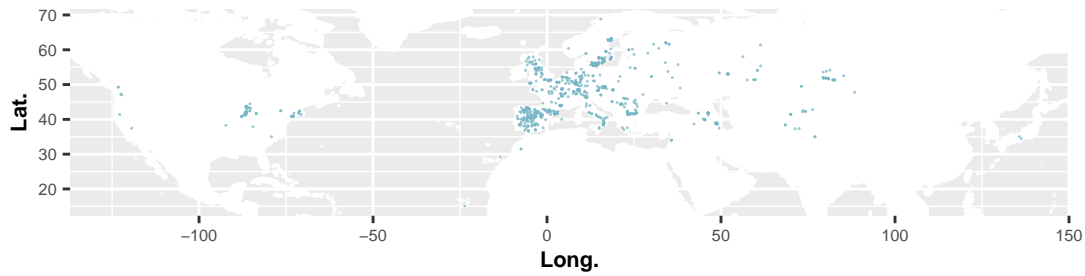


Figure S1. Geographic locations for all 1,035 *Arabidopsis* accessions. The locations for 1,035 accessions used to define latent environmental variables.

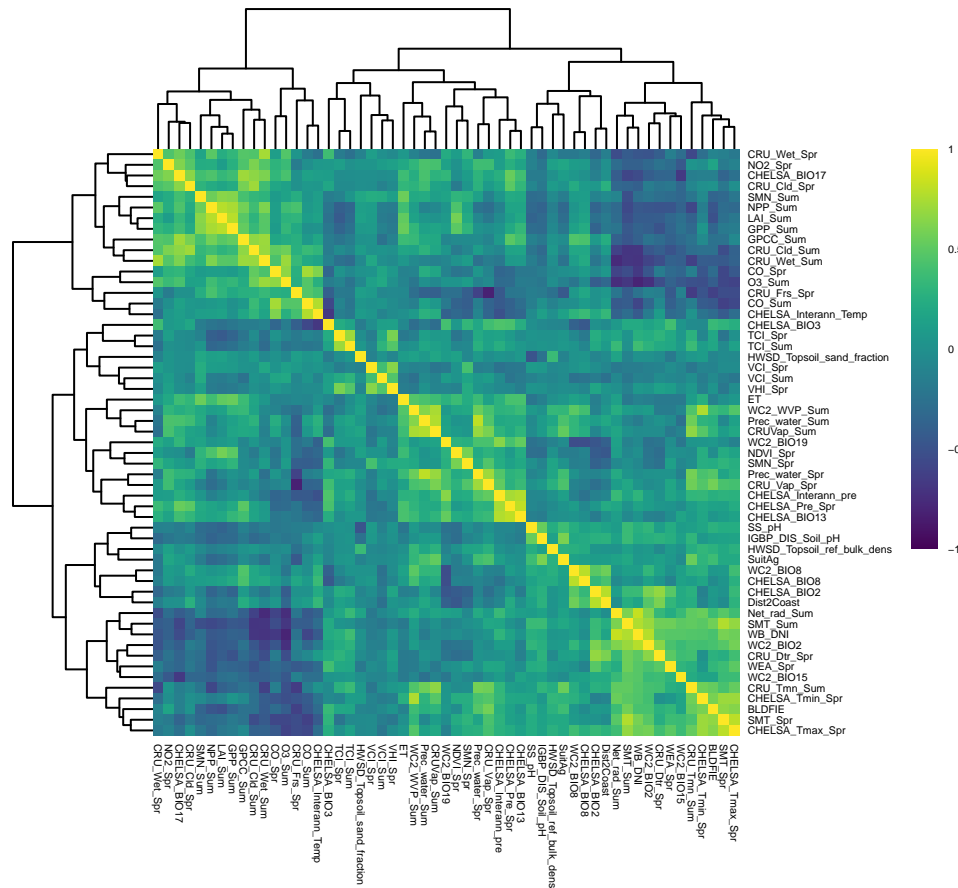


Figure S2. Heatmap for 55 manifest environmental variables. Spearman's method was used to generate the correlation matrix.

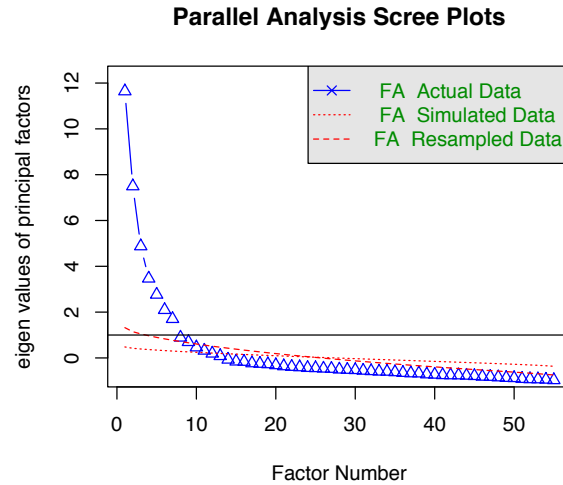


Figure S3. Scree plot indicating the optimal number of latent factors for 55 environmental variables. Parallel factor analysis was performed using the psych package in R. This approach generates scree plots for the observed data and compares the results with scree plots generated from a random data matrix of the same size as the observed data set.

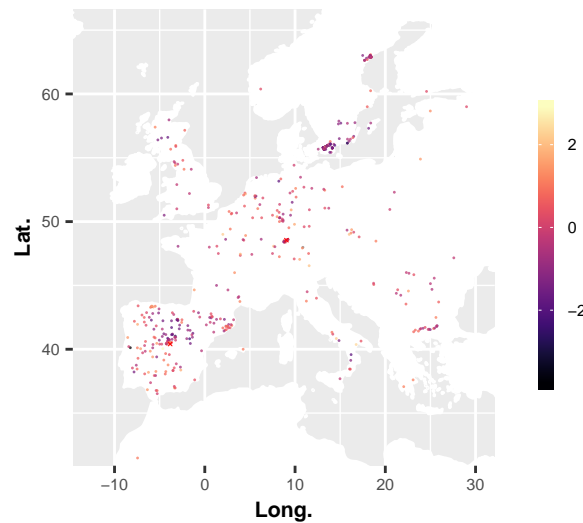


Figure S4. Projection of phenotypic values for flowering time plasticity on collection sites for 515 accessions. Higher plasticity values indicate a greater delay in flowering time Tuebingen relative to Madrid, and is indicated by the continuous color scale on the right. The red 'X' indicates the two common garden locations.