# Weighted AMMI to study genotype-by-environment interaction and QTL-by-environment interaction

Paulo Canas Rodrigues[1]

Marcos Malosetti[2]

Hugh G. Gauch[3]

Fred van Eeuwijk[2]

## 1    Introduction

A differential response of genotypes across environments (often, location by year combinations) is frequent in multi-environment trials (METs) and is known as genotype-by-environment interaction (GEI). Data from METs are often summarized in two-way tables with genotypes in the rows and environments in the columns. GEI occurs in various forms, with the most extreme consisting of crossovers, when there is a change of ranking of genotypes across environments, e.g., a genotype that is superior under well watered conditions may yield poorly under dry conditions. The study and understanding of GEI is a major challenge in the improvement of complex traits like yield across environmental gradients.

The additive main effects and multiplicative interaction (AMMI) model (Gauch 1992) is one of the most widely used statistical methods to understand and structure interactions between genotypes and environments. In essence the AMMI model applies the singular value decomposition (SVD) to the residuals of the analysis of variance (ANOVA). However, if there is a strong GEI in the data, we also expect the trials to have heterogeneous error variances, and this is not taken into account by the standard AMMI model. Therefore, we propose a generalization of the AMMI model that is able to take into account heterogeneity of error variance by using a weighted low-rank SVD, the weighted AMMI (WAMMI) model. Although this generalization occurs in a fixed model, WAMMI offers a reasonable approximation to mixed model methodology for GEI, which is considered to be more appropriate in case of heterogeneous error variances.

A natural follow up to the analysis of GEI, is the study of the genetic factors underlying GEI: QTL (quantitative trait locus) and environment interaction, QEI. Gauch et al. (2011)

---

[1] CMA – FCT/Universidade Nova de Lisboa, Portugal. e-mail: paulocanas@gmail.com
[2] Wageningen University, The Netherlands.
[3] Cornell University, USA.

proposed the AQ analysis where the AMMI model is used to obtain predicted values for genotype-by-environment combinations, which are then used in QTL mapping. We believe that the weighting by the (reciprocal) of error variances in the AMMI model will not only improve the analysis of GEI, but equally so that of subsequent QTL analysis. Thus, we present a generalization of the AQ analysis that is able to account for heterogeneity in both genetic variances, captured by the interaction principal components in AMMI, and error variances, by weighting; we replace the AMMI model by the WAMMI model. The weighted version of the AQ analysis, WAQ, can be conducted in three stages: (i) compute the weights for each environment based on the error variances; (ii) fit the WAMMI model to the GEI data table; and (iii) perform the QTL scans using the predictions from the WAMMI model as response variable. In the spirit of AMMI, with our WAMMI approach we expect to separate signal, GEI and QEI patterns, from noise.

WAQ is compared with the QTL analyses on the actual data, with the AQ analysis (Gauch *et al.* 2011), and with a QTL mixed models approach (BOER *et al.* 2007; MALOSETTI *et al.* 2004). Two data sets were used. The first one deals with yield simulated from a backcross pepper (*Capsicum annuum*) population using a crop growth (physiological) genotype-to-phenotype model (Rodrigues *et al.* 2012). The motivation for using a crop growth model to transform genotypic information to phenotypic information was that we wanted a biologically realistic data set, while still wanting to know the underlying genetic architecture. The second data set concerned yield for the well-known Steptoe x Morex barley (*Hordeum vulgare* L.) population, originating from the North American Barley Genome Mapping Project (HAYES *et al.* 1996; HAYES *et al.* 1993).

WAMMI is applicable to a wide range of fields to which also AMMI has been applied, including more than 200 articles referenced by the ISI web of knowledge within the last ten years. In addition applications in plant breeding, crop sciences and genetics, AMMI was applied to microarray experiments (Crossa *et al.* 2005), rDNA studies (ADAMS *et al.* 2002), plant and microbial populations' growth across several environmental conditions (CULMAN *et al.* 2009; CULMAN *et al.* 2008), and animal sciences (BARHDADI and DUBE 2010).

## 2    Materials and methods

### a.    Plant materials

The primary data set of this study is a two-way table with $I = 200$ genotypes and $J = 12$ environments (Table 1) of the complex trait yield. These data were simulated by assuming that the final yield set equals the signal plus the noise. The signal for genotype $i$ in

environment $j$ was simulated from a eco-physiological genotype-to-phenotype crop growth model (CGM) for pepper (Rodrigues *et al.* 2012) and is a function of physiological parameters and environmental characterizations. The model can be written as:

$$Yield_{i,j} = Signal + Noise$$

$$= \frac{FTF_i \times [1 - W_j \times (T_j - T_{FTF})]}{FDMC_i} \times LUE_{i,j} \times \sum_{t=t_0}^{t_f} [1 - exp(-K_i \times LAI_{i,j,t})] \times I_{j,t} + \varepsilon_{i,j} \quad (1)$$

where $t_0$ and $t_f$ represent the beginning and the end of the growing season, in days, $LAI_{i,j,t}$ the leaf area index for genotype $i$, environment $j$ and day $t$, and $I_{j,t}$ is the photosynthetic active radiation incident on the crop for environment $j$ on day $t$, and $\varepsilon_{i,j}$ is the error (or noise) for the $Yield_{i,j}$.

The model (1) is a function of seven physiological parameters: $K$ (light extinction coefficient); $LUE$ (maximum light use efficiency); $B$ (slope for the leaf area increase with temperature sum, used to define $LAI$); $FTF$ (fraction of dry weight partitioned to the fruits); $FDMC$ (fruit dry matter content); $W$ (slope of the linear reduction in harvest index with temperature above 15°C); and $Z$ (slope of the linear reduction in $LUE$ for temperatures below 20°C, used to define $LUE$); and three environmental variables: temperature, radiation and country (Table 1). More details can be found in Rodrigues et al. (2012).

The main motivation for using a nonlinear physiological genotype-to-phenotype model instead of a statistical model is to ensure that the simulated data is close to a biologically credible model, where we have full information on the biological background.

Each of the seven physiological parameters (component traits) was simulated as a sum of a number of QTLs (Table 2) plus a residual effect. This was done for 200 simulated pepper genotypes, characterized by 237 markers covering all the 12 chromosomes (Barchi *et al.* 2007; Barchi *et al.* 2009). In this simulation several QTLs were placed along the 12 chromosomes of the pepper genome. The exact positions of these QTLs and their heritability are described in Table 2. The simulations were made using the package qtl (BROMAN and SEN 2009) of the statistical software R. More details on the model and physiological parameters can be found in Rodrigues et al. (2012).

The noise for yield $\varepsilon_{i,j}$ in equation (1) was simulated from a Gaussian distribution with zero mean and variance $\sigma_{\varepsilon_j}^2, j = 1, \dots, J$, depending on the environment (Table 1) and the chosen heritability for yield ($h^2$, Table 2), i.e.

$$\sigma_{\varepsilon_j}^2 = \frac{1 - h^2}{h^2} \sigma_{g_j}^2,$$

where $\sigma_{\varepsilon_j}^2$ and $\sigma_{g_j}^2$ are the error and genetic variance (from the eco-physiological genotype-to-phenotype model) for the environment $j$, $j = 1, \dots, J$ (Table 1). The final yield data is the result of the sum of the signal and the noise as in equation (1). This simulation was repeated 100 times resulting in 100 two-way tables with 200 genotypes and 12 environments.

**Table 1.** The 12 environments used in the simulated yield data for pepper. Description of the environments considered in the genotype-to-phenotype crop growth model. The first column represents the code for the environments which is used in the text and figures. The countries were chosen to represent different environmental and practical conditions (Rodrigues et al., 2012a). Radiation has two levels (years) based on historical data. Temperature contains three levels of daily average temperature. The heritability for the environments was set to be $h^2 = 0.5$. The mean genetic and error variances, for the 100 simulated data sets, are reported in the last two columns.

| Environment | Country | Radiation | Temperature | Genetic variance | Error variance |
|---|---|---|---|---|---|
| NL1-15 | Netherlands | Lower | 15ºC | 21.27 | 21.32 |
| NL1-20 | Netherlands | Lower | 20ºC | 39.29 | 39.60 |
| NL1-25 | Netherlands | Lower | 25ºC | 39.22 | 39.19 |
| NL5-15 | Netherlands | Higher | 15ºC | 35.08 | 34.82 |
| NL5-20 | Netherlands | Higher | 20ºC | 65.14 | 64.81 |
| NL5-25 | Netherlands | Higher | 25ºC | 65.25 | 64.77 |
| SP1-15 | Spain | Lower | 15ºC | 9.91 | 9.79 |
| SP1-20 | Spain | Lower | 20ºC | 19.06 | 19.27 |
| SP1-25 | Spain | Lower | 25ºC | 19.66 | 19.96 |
| SP5-15 | Spain | Higher | 15ºC | 13.16 | 13.24 |
| SP5-20 | Spain | Higher | 20ºC | 25.43 | 25.75 |
| SP5-25 | Spain | Higher | 25ºC | 26.29 | 26.46 |

**Table 2.** Genetic architecture of the simulated yield data for pepper (signal). The first columns give the name of the parameter related to the QTL, the code for the closest marker, the chromosome, the position and its heritability when included in the physiological genotype-to-phenotype model. The last column gives a summary on which environments the QTLs are expected to be detected, being the $FTF$ and $FDMC$ much weaker and harder to detect than the $LUE$.

| Parameter | Marker | Chromosome | Position (cM) | Heritability | Importance[1] |
|---|---|---|---|---|---|
| $K$ | D1M6 | 1 | 38.0 | 0.95 | None |
| $LUE$ | D2M13 | 2 | 55.0 | 0.16 | All |
| $B$ | D3M10 | 3 | 87.3 | 0.95 | None |
| $FTF$ | D4M9 | 4 | 83.1 | 0.80 | All |
| $FDMC$ | D5M25 | 5 | 103.3 | 0.80 | All |
| $W$ | D6M12 | 6 | 36.7 | 0.20 | 20ºC; 25ºC |
| $LUE$ | D7M5 | 7 | 42.5 | 0.16 | All |
| $LUE$ | D8M7 | 8 | 38.8 | 0.16 | All |
| $W$ | D9M15 | 9 | 100.4 | 0.20 | 20ºC; 25ºC |
| $LUE$ | D10M5 | 10 | 43.1 | 0.16 | All |
| $Z$ | D11M11 | 11 | 62.4 | 0.80 | 15ºC |

[1]Based on a sensitivity analysis with heritability of 1 for all environments (Table 1).

The second data set in our study is a subset of the grain yield data from the Steptoe x Morex (SxM) cross, produced by the North American barley genome mapping project (HAYES *et al.* 1996; HAYES *et al.* 1993). The data contain 150 doubled haploid genotypes evaluated in 16 environments during 1991 and 1992, in USA and Canada. The genotypes were characterized by 116 markers covering all seven chromosomes. For inclusion in our data set, environments needed to have either a complete replication (block) or a complete replication block and an additional partial replication. The 13 chosen environments are presented in Table 3. The trials conducted in 1991 had a full replicate/block and a second one containing only 50 genotypes. For trials in 1992 two complete replications were available.

### b. AMMI analysis

The additive main effects and multiplicative interaction (AMMI) model (Gauch 1988; Gauch 1992) combines together the features of analysis of variance (ANOVA) and singular value decomposition (SVD), the SVD is applied to the residuals from the additive ANOVA, i.e. to the GEI. In the ANOVA part, the additive main effects are estimated, whereas the SVD models the interaction via $N$ axes, or $N$ interaction principal components, IPCs, $N \leq \min(I - 1, J - 1)$, with $I$ the number of genotypes (rows) and $J$ the number of environments (columns). The model is usually written as (Gauch 1992)

$$y_{i,j} = \mu + \alpha_i + \beta_j + \sum_{n=1}^{N} \lambda_n \gamma_{n,i} \delta_{n,j} + \varepsilon_{i,j}, \tag{2}$$

where $y_{i,j}$ is the yield of genotype $i$ in environment $j$, $\mu$ the grand mean, $\alpha_i$ the genotype deviations from $\mu$, $\beta_j$ the environment deviations from $\mu$, $\lambda_n$ is the singular value for the IPC axis $n$, $\gamma_{n,i}$ and $\delta_{n,j}$ the genotype and environment IPC scores (i.e. the left and right singular vectors) for axis $n$, and $\varepsilon_{i,j}$ the residual containing both multiplicative terms not included in the model (2) as well as an experimental error. A matrix formulation of equation (2) can be given by

$$\boldsymbol{y} = \mathbf{1}_I \mathbf{1}_J^T \mu + \boldsymbol{\alpha} \, \mathbf{1}_I^T + \mathbf{1}_J^T \boldsymbol{\beta} + \boldsymbol{U} \boldsymbol{D} \boldsymbol{V}^T + \boldsymbol{\varepsilon}, \tag{3}$$

where $\boldsymbol{y}$ is the $(I \times J)$ two-way data table, $\mathbf{1}_I \mathbf{1}_J^T \mu$ is a $(I \times J)$ matrix with the grand mean $\mu$ in all positions, $\boldsymbol{\alpha} \, \mathbf{1}_I^T$ is a $(I \times J)$ matrix with the genotype deviations from the grand mean (equal rows), $\mathbf{1}_J^T \boldsymbol{\beta}$ is a $(I \times J)$ matrix with the environment deviations from the grand mean (equal columns), $\boldsymbol{U}$ is a $(I \times N)$ matrix whose columns contain the left singular vectors of the multiplicative part of the data, i.e. $\tilde{\boldsymbol{y}} = \boldsymbol{Y} - \boldsymbol{\mu} - \boldsymbol{\alpha} - \boldsymbol{\beta}$, $\boldsymbol{D}$ a $(N \times N)$ diagonal matrix containing the singular values of $\tilde{\boldsymbol{y}}$ in the diagonal, $\boldsymbol{V}$ is a $(J \times N)$ matrix whose columns

contain the right singular vectors of $\widetilde{\boldsymbol{y}}$, and $\boldsymbol{\varepsilon}$ is the $(I \times J)$ matrix with the residuals. With this procedure we are aiming at a low rank $(N)$ approximation to the matrix $\widetilde{\boldsymbol{y}}$, i.e. the interaction.

**Table 3.** The 13 environments used in the SxM analysis. The first column gives the code for the environment (location year combination) which is used in the text and figures. The information about full replication of partially replication, genetic and error variances, and heritability are presented in the next columns.

| Environment | Full replication | Genetic variance | Error variance | Heritability |
|---|---|---|---|---|
| ID91 | No | 0.94 | 0.74 | 0.56 |
| ID92 | Yes | 0.55 | 0.42 | 0.57 |
| MAN92 | Yes | 0.38 | 0.20 | 0.66 |
| MIN92 | Yes | 0.35 | 0.37 | 0.49 |
| MTd91 | No | 0.39 | 0.11 | 0.78 |
| MTd92 | Yes | 0.43 | 0.31 | 0.58 |
| MTi91 | No | 0.36 | 0.23 | 0.61 |
| MTi92 | Yes | 0.43 | 0.31 | 0.58 |
| NY92 | Yes | 0.20 | 0.67 | 0.23 |
| ONT92 | Yes | 0.21 | 0.31 | 0.40 |
| OR91 | No | 0.26 | 1.66 | 0.14 |
| WA91 | No | 0.72 | 0.71 | 0.50 |
| WA92 | Yes | 0.20 | 0.25 | 0.44 |

The number of interaction terms in the model, $N$, has to be chosen wisely as it will affect all the subsequent results (GAUCH *et al.* 2008; YANG *et al.* 2009). In this paper we use a cross-validation based method proposed by Krzanowski (1987). By considering the model for the GEI data table $x_{i,j} = \sum_{n=1}^{N} \lambda_n \gamma_{n,i} \delta_{n,j} + \varepsilon_{i,j}$, we are able to compute the average squared discrepancy between the actual and predicted values:

$$PRESS(n) = \frac{1}{IJ} \sum_{i=1}^{I} \sum_{j=1}^{J} \left( \hat{x}_{i,j}^{(n)} - x_{i,j} \right)^2, \tag{4}$$

and, consequently,

$$W_n = \frac{PRESS(n-1) - PRESS(n)}{I+J-2n} \div \frac{PRESS(n)}{D_r}, \tag{5}$$

where $D_r$ can be obtained by sequential subtraction from $(I-1)(j-1)$ of $I+J-2N$. The $W_n$ represent the increase in predictive information supplied by the $n$-th component, divided by the average predictive information in each of the remaining components (Krzanowski 1987). Krzanowski (1987) suggested that the optimal number of components is the highest number of $n$ such that $W_n$ is greater than 0.6.

### c. Weighted AMMI analysis

When the two-way data table $y$ has missing cells and/or the error variance is not constant across environments, the cells of the table should have different weights for their squared residuals in the estimation procedure for the model parameters. To account for heterogeneity of error variances across environments, our proposal is to replace the standard low-rank SVD in equation (3) by a weighted low-rank SVD (GABRIEL and ZAMIR 1979). The approach we use here is based on an expectation-maximization (EM) procedure and, while the sum of squares of the difference between two consecutive iterations, $X^{(t+1)}$ and $X^{(t)}$, is greater than some small value, e.g. $10^{-9}$, we run

$$X^{(t+1)} = SVD\big(W\odot y + (1 - W)\odot X^{(t)}\big) \tag{6}$$

where $W$ is a $(I \times J)$ matrix with weights, $W_{i,j}$, $0 \le W_{i,j} \le 1$, $1$ is a $(I \times J)$ matrix with ones in all positions, $\odot$ the Hadamard (or entrywise) product of matrices, and $t$ is the iteration number (SREBRO and JAAKKOLA 2003). $X$ should be initialized to $X^{(0)} = y$ or to $X^{(0)} = 0$. The outputs of this procedure are the matrices $U_N$, $D_N$ and $V_N$ such that $\tilde{y} \approx U_N D_N V_N{}'$, being $N$ the rank of approximation. The R code for this algorithm, with detailed explanation, can be found in the File S1 of the supplementary material.

Applying the weighted low-rank SVD (6) to the matrix $\tilde{y}$ and replacing in equation (3) will result in the weighted AMMI (WAMMI) model. This generalization of the AMMI model is now able to account for differences in error variances across environments and/or missing cells, and can be applied to all data sets where the AMMI model has been used. Of course, we need to be able to estimate the error variance for an environment, so we need at least partial replication per environment. It should be remarked that in this paper we first estimate the main effects without using weights, then produce the residuals from additivity and finally approach these residuals by a weighted SVD. We could also have used the weights already in the estimation of the main effects, but this approach may be less robust. We feel that this topic merits further study.

With partial replication, cell means based on more replication will have smaller variances than those with less replication. The scheme of weights should reflect the number of replications per cell. The a $(I \times J)$ matrix with weights, $W_{i,j}$, $0 \le W_{i,j} \le 1$, can be calculated from the Hadamard (or entrywise) product of two matrices: (i) a matrix in which the entries are column wise constant, being the inverse of the error variance; and (ii) a matrix with the proportion of replications per cell, i.e.

$$W = \begin{bmatrix} \dfrac{1/\sigma_1^2}{m} & \dfrac{1/\sigma_2^2}{m} & \cdots & \dfrac{1/\sigma_J^2}{m} \\ \dfrac{1/\sigma_1^2}{m} & \dfrac{1/\sigma_2^2}{m} & \cdots & \dfrac{1/\sigma_J^2}{m} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{1/\sigma_1^2}{m} & \dfrac{1/\sigma_2^2}{m} & \cdots & \dfrac{1/\sigma_J^2}{m} \end{bmatrix} \odot \begin{bmatrix} \dfrac{Nrep_{1,1}}{Nrep} & \dfrac{Nrep_{1,2}}{Nrep} & \cdots & \dfrac{Nrep_{1,J}}{Nrep} \\ \dfrac{Nrep_{2,1}}{Nrep} & \dfrac{Nrep_{2,2}}{Nrep} & \cdots & \dfrac{Nrep_{2,J}}{Nrep} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{Nrep_{I,1}}{Nrep} & \dfrac{Nrep_{I,2}}{Nrep} & \cdots & \dfrac{Nrep_{150,J}}{Nrep} \end{bmatrix}, \tag{7}$$

where $I$ is the number of genotypes, $J$ the number of environments $m = max_j(1/\sigma_{\varepsilon_j}^2)$, $\sigma_{\varepsilon_j}^2, j = 1, ..., J$, is the error variance for environment $j$, $Nrep_{i,j}, i = 1, ... I, j = 1, ..., J$, is the number of replications for genotype $i$ in environment $j$, and $Nrep$ is the maximum number of replications in the data set.

### d. Weighted AQ analysis

Gauch et al. (2011) suggested a new approach for detecting and understand QEI, the AQ analysis, where the QTL scans are made based on AMMI predictions (instead of direct QTL scans on the actual data). In this paper we make use of the above proposed weighted version of the AMMI model, the WAMMI, to generalize the AQ analysis to account for both heterogeneous genetic (SVD) and error variances (weighs) across environments. We use a fixed effects model as an alternative to the QTL mixed models approach (BOER *et al.* 2007; MALOSETTI *et al.* 2004) that can be fitted with standard statistical software for linear models. The weighted AQ analysis can be conducted in three stages: (i) compute the weights for each environment based on the error variances, i.e. the weights are given by the inverse of the error variance in each environment and are (usually) constant for all genotypes in an environment; (ii) fit the WAMMI model to the GEI data table and obtain the predicted values for each combination of genotype and environment; and (iii) perform the QTL scans using the WAMMI predicted values as response variable for each environment separately. This approach can potentially improve the power for QTL detection as it uses improved genotypic predictions as response variable that showed to be better than the means from the ANOVA model. The environments can then be ordered by AMMI and WAMMI parameters that summarize GEI and QEI information to reveal consistent patterns and systematic trends that often can be explained in terms of environmental conditions (Gauch 1992; Gauch *et al.* 2011).

### e. Linear mixed model

As a kind of bench mark for QTL analysis, we analysed the simulated pepper and Steptoe x Morex barley data also by a QTL analysis based on mixed models, as described by Boer et al. (2007), and implemented in Genstat 14 (Payne *et al.* 2011). The input for the QTL analysis
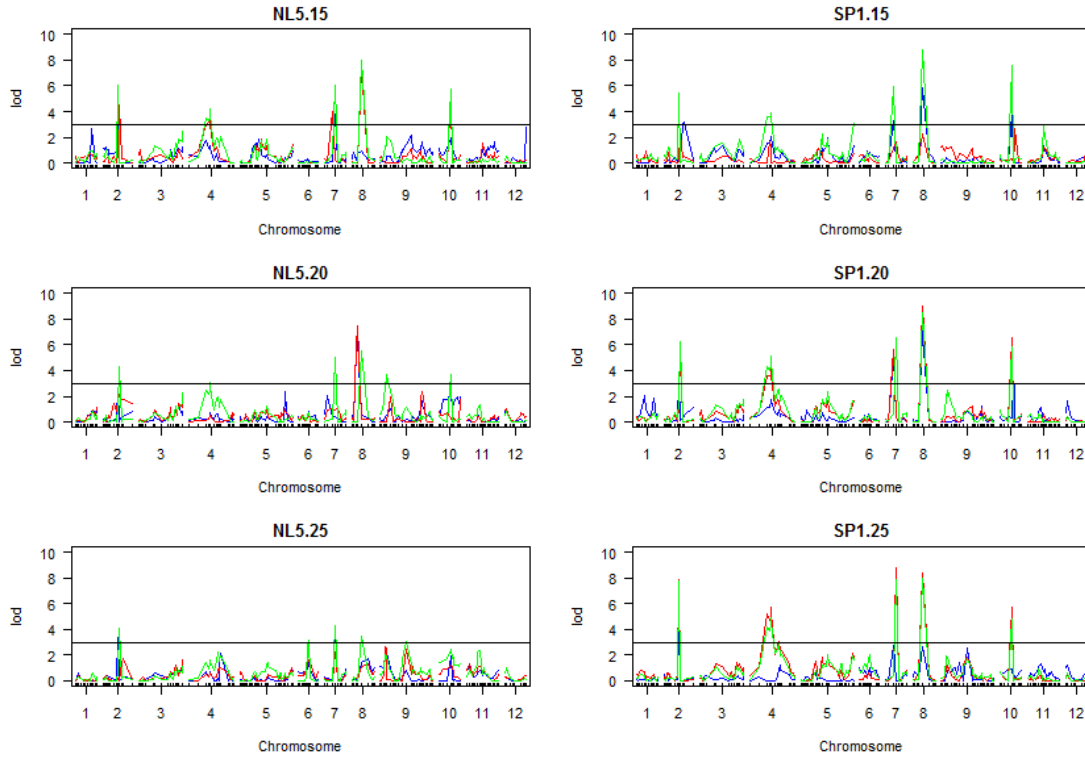
consists of the genotype-by-environment means and corresponding weights, defined in the way also used in WAMMI and WAQ. The QTL analysis fits fixed environment specific QTLs, i.e., actually the sum of QTL main effect and QEI, to the genotypic main effects and GEI, say GGEI, jointly. The model contains a multivariate normal distribution for the GGEI effects allowing heterogeneity of genetic variances and correlations.

## 3   Results for the simulated pepper data

### a.   Preliminary analysis

Table 2 gives the simulation conditions and, therefore, the "true" genetic architecture for the pepper population under study. Figure 1 (blue line) depicts the single trait single environment QTL scans for 6 environments, of the complex trait yield simulated from the physiological genotype-to-phenotype model with seven physiological parameters (Rodrigues *et al.* 2012). The six environments were chosen to represent the three levels of temperature with the lowest and highest error variances. Comparing the "true" genetic architecture in Table 2 and the QTLs detected in Figure 1 and Figure S1 (QTL scans for all 12 environments) for the actual data, only those associated with the parameters *LUE* (22 out of the expected 48 = 12 environments times 4 chromosomes, Table 2) and *W* (three out of the expected 16 = 8 environments with temperatures of 20°C or 25°C times 2 chromosomes, Table 2) were found, which represents a poor outcome of this single trait single environment analysis.

**Figure 1.** QTL scans for 6 environments of the yield data for pepper simulated from the genotype-to-phenotype model with seven physiological parameters (Rodrigues et al. 2012a). Each row represents a different level of temperature. The plots on the left correspond to the highest error variance in this simulated data table and the ones on the right to the lowest. The blue line represents the scans for the actual data, the red for the AMMI2 predicted values, and the green for the WAMMI2 predicted values. All the scans are based on composite interval mapping. The horizontal lines correspond to the thresholds for a LOD score of 3. These scans are based on one randomly chosen realization out of the 100 simulations. The codes for the captions of the individual scans are described in Table 1.
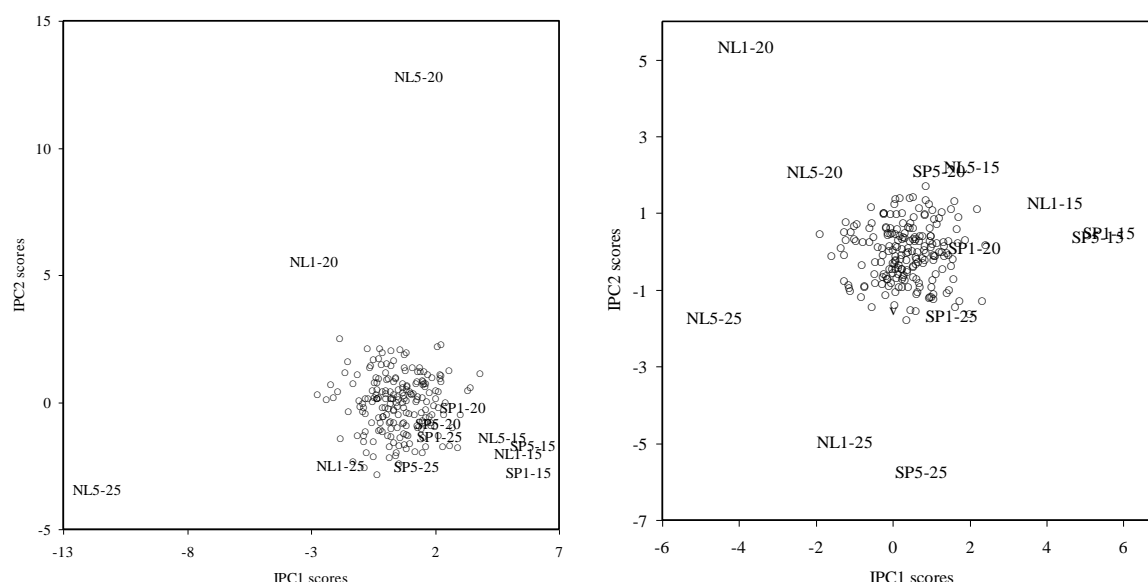
### b. AMMI analysis

Table 4 gives the ANOVA for the model AMMI5 based on one randomly chosen realization of a genotype-by-environment two-way data table. Similar results are obtained for other two-way data tables simulated from the model in use (Rodrigues *et al.* 2012). The genotypes, environments and GEI account for 31.5, 34.4 and 34.1% of the treatment sum of squares (SS). Two interaction principal components were chosen for the AMMI model as in Rodrigues et al. (2012). This choice was confirmed by the cross-validation proposed by Krzanowski (1987): the $W_n$ values from equation (5) for the first five components are 10.371; 1.385; 0.579; 0.475; 0.306, and show the "best" model to have two principal components because only two $W_n$ values are above the cut-off of 0.6.

The AMMI2 biplot is depicted on the left hand side of Figure 2. In this figure, the environments with higher genetic variance (NL5-20 and NL5-25, Table 1) are farthest away from the origin, showing an extreme behaviour when compared with the remaining environments. However these environments with higher genetic variance also have higher error variance (Table 2), which should be down-weighted to produce a more trustable result. This can be achieved by giving smaller weights to the environments with higher error variance.

**Table 4.** ANOVA of the AMMI5 model for the simulated yield data for pepper. Results based on one randomly chosen realization of the genotype-to-phenotype crop growth model. The columns of the table show the source of variation, the degrees of freedom (df), the sums of squares (SS) and the mean squares (MS).

| Source | df | SS | MS |
|---|---|---|---|
| Total | 2399 | 256089 | 106.7 |
| Genotypes | 199 | 80774 | 405.9 |
| Environments | 11 | 88054 | 8004.9 |
| GEI | 2189 | 87261 | 39.9 |
| IPC1 | 209 | 18122 | 86.7 |
| IPC2 | 207 | 14740 | 71.2 |
| IPC3 | 205 | 11470 | 56.0 |
| IPC4 | 203 | 10074 | 49.6 |
| IPC5 | 201 | 7916 | 39.4 |
| IPC6—IPC11 | 1164 | 24938 | 21.4 |



**Figure 2.** AMMI2 (left) and WAMMI2 (right) biplots for one randomly chosen realization. The abscissa represents the first multiplicative term and the ordinate the second. The open dots represent the 200 genotypes and the codes for the 12 environments are defined in Table 1.

### c. Weighted AMMI analysis

To avoid considering environments with high error variance as outliers (Gauch *et al.* 2011) or letting them influence (too much) the results, the weighted AMMI analysis described above was used where the contribution (i.e. the weight) of a given environment to the model fit is the inverse of its error variance (Table 1). The WAMMI biplot is given in Figure 2 (right). In this plot the environments SP5-20 and SP5-25 ceased to show extreme behaviour. There is also a visible pattern in the environments: (i) the right hand side presents more Spanish environments whereas the left hand side has more Dutch environments; and (ii) the right top corner shows environments with temperatures of 15ºC, the left bottom corner shows environments with temperatures of 25ºC, and in between are placed the environments with temperatures of 20ºC.

### d. AQ analysis and weighted AQ analysis

The AQ analysis is the AMMI analysis followed by QTL scans on the AMMI predicted values (Gauch *et al.* 2011). The weighted AQ (WAQ) analysis is a generalization of the AQ analysis, where the AMMI analysis is replaced by the weighted AMMI analysis proposed before. This WAMMI and WAQ analyses are particularly useful to analyse data sets whose environments show high heterogeneity in their error variances.

Figure 1 shows the AQ (red line) and WAQ (green line) analyses for models with two IPCs. There is a clear improvement from the QTL scans of the actual data to the AQ and WAQ analysis in both the number of detected QTLs and higher LOD scores. As in the biplots of Figure 2, the improvement from the unweighted to the weighted method is visible in Figure 1 for AQ and WAQ analysis. As an example, all peaks of environment SP1-15 (lowest error variance, Table 1) are below the LOD 3 threshold when using AQ analysis but five (true) QTLs are detected when using the WAQ analysis. This happens because, being SP1-15 the most accurate environment (lowest error variance), its weight is expected to be underestimated by the AMMI2 model but corrected with the WAMMI2.

### e. The 100 simulated data sets and comparison between methods

A more detailed comparison for all the 100 simulated data sets is presented in Figures 3 and 4. As expected, the worst performance (in terms of detected QTLs) is obtained by the QTL scans of the actual data. Better are the QTL scans on the AMMI2 predicted values (AQ analysis), which, however, do not detect QTLs for some environments in some runs. The WAQ analysis and QTL mixed model framework are the best options in the presence of
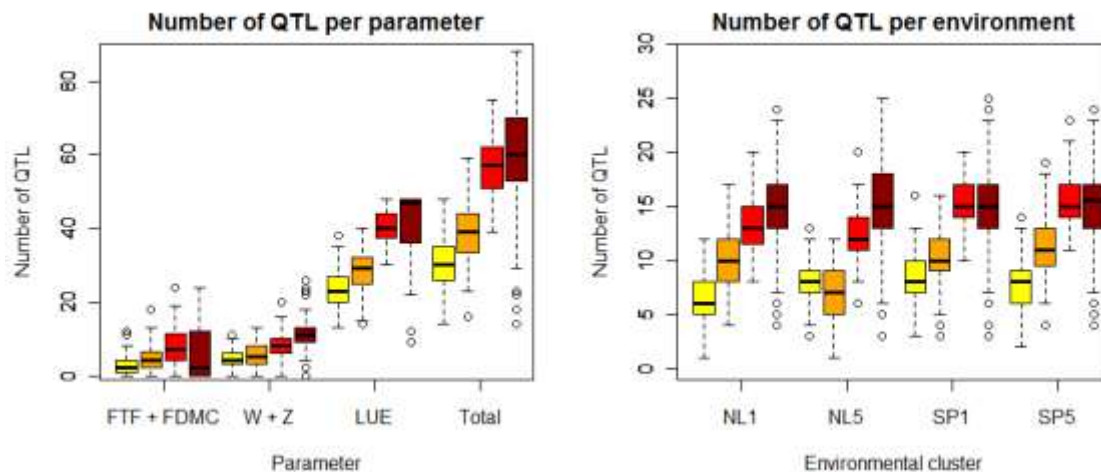
heterogeneity of error variance across environments. Although the mixed model detects slightly more QTLs, the fixed effects WAQ analysis shows less variance for the number of QTLs (Figures 3 and 4).

The analysis and interpretation was clearly improved by using the error variances in each environment, which leads us to conclude that the WAMMI biplot is also an improved version (closer to the reality) of the AMMI biplot (Figure 2).
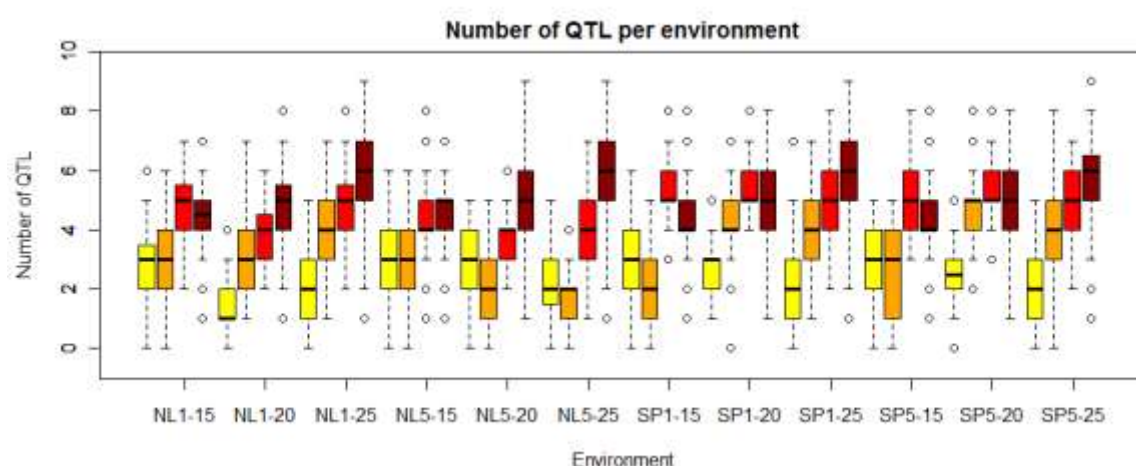
## 4    Results for the barley experiment
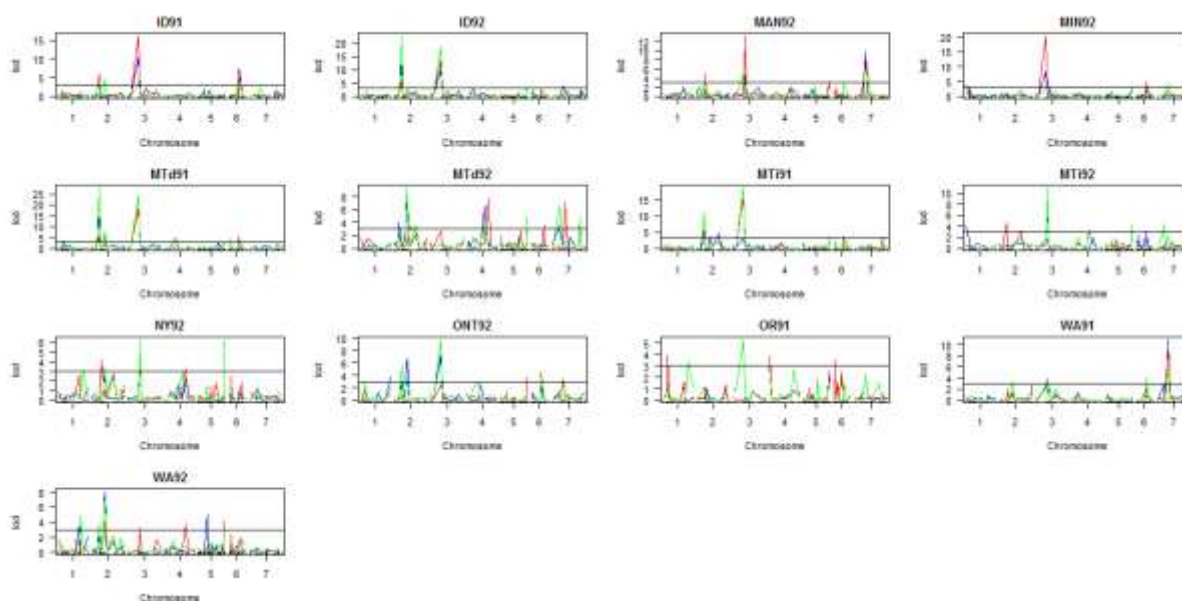
### a.    Preliminary analysis

Two previous studies have applied the AMMI model to the SxM yield data to improve and better understand QTL detections (GAUCH *et al.* 2011; ROMAGOSA *et al.* 1996). Here we used the genotype-by-environment means for 13 environments (Table 3), where the experiment was partially replicated, instead of the means for the original 16 environments. Table S1 gives a short summary of findings in the literature about detected QTLs on the SxM yield data. Figure 5 (blue line) depicts the QTL scans for the actual data of the 13 environments.



**Figure 3.** Summary of the number of detected QTLs for the actual data (yellow), AMMI2 predicted values (orange), WAMMI2 predicted values (red) and linear mixed model (dark red). The graph on the left hand side shows the box plots for the number of QTLs per model parameter (Table 2), and on the right hand side the number of QTLs per environmental cluster (Table 1). These values are for QTLs detected when considering an interval of 20 cM centred on the right QTL position. These plots are for a heritability of 0.5 in all environments.

**Figure 4.** Number of QTLs detected per environment for an expected maximum of 7 (environments with temperature of 15ºC or 20ºC, Table 2) or 8 (environments with temperature of 25ºC, Table 8). The box plots are presented for the actual data (yellow), AMMI2 predicted values (orange), WAMMI2 predicted values (red) and linear mixed model (dark red). These plots are for a heritability of 0.5 in all environments.



**Figure 5.** QTL scans for the 13 environments for the means of the SxM yield data (blue line), AMMI3 predicted values (red line), and WAMMI3 predicted values (green line). The results are for composite interval mapping and the threshold was set to a LOD score of 3.

### b. AMMI analysis

Table 5 gives the ANOVA for the AMMI5 model. The genotypes, environments and GEI account for 9.2, 67.4 and 23.4% of the treatments sum of squares (SS). The amount of noise
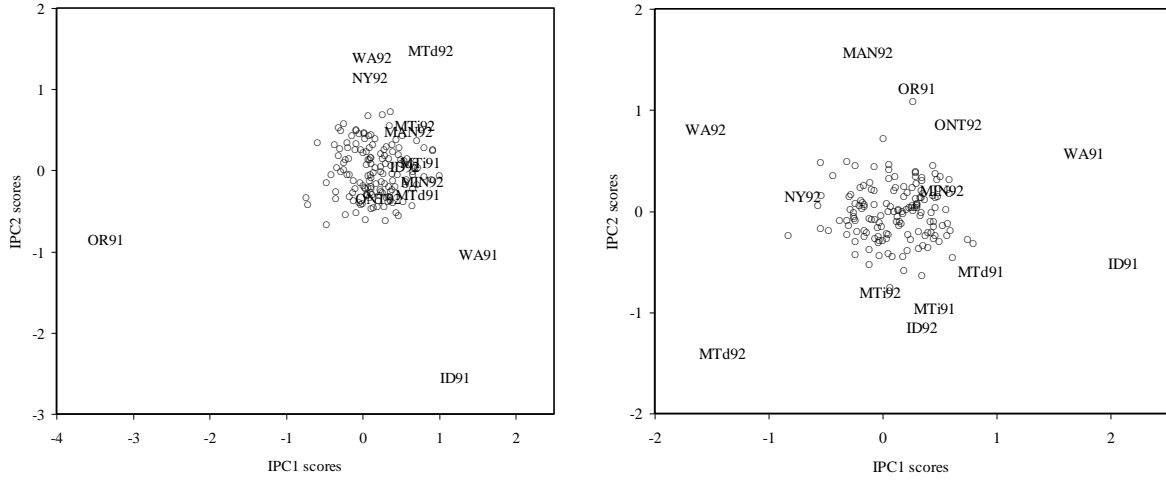
in the GEI can be estimated by the product of the interaction degrees of freedom (df) with the error mean square (MS), namely 768.8, which by difference from the total of 2157 implies a GEI signal of 1388.2, or 64.4% (Gauch 1992; Voltas *et al.* 2002). IPC1 captures a SS of 566, IPC2 412 and IPC3 287, which includes the most of the signal and little noise because the first principal components tend to capture more signal and less noise (Gauch 1992).

**Table 5.** ANOVA of the AMMI5 model for the SxM yield data.

| Source | df | SS | MS |
|---|---|---|---|
| Total | 3399 | 9829 | 2.89 |
| Treatments | 1949 | 9202 | 4.72 |
| Genotypes | 149 | 844 | 5.66 |
| Environments | 12 | 6201 | 516.77 |
| GEI | 1788 | 2157 | 1.21 |
| IPC1 | 160 | 566 | 3.54 |
| IPC2 | 158 | 412 | 2.61 |
| IPC3 | 156 | 287 | 1.84 |
| IPC4 | 154 | 227 | 1.47 |
| IPC5 | 152 | 137 | 0.90 |
| IPC6—IPC11 | 1008 | 528 | 0.52 |
| Intra Block Error | 1450 | 626 | 0.43 |

Within the two studies where the AMMI model was applied to the SxM yield data, Romagosa et al. (1996) found QTLs in the first four IPCs. Subsequently Gauch et al. (2011) considered the AMMI3 based on the Ockham's valley for the root mean squared prediction error following from a jackknife procedure. For this particular data set the cross-validation procedure described before and introduced by Krzanowski (1987) was considered. When computing the $W_n$ values for the first five components of the SxM yield we obtain: 11.019; 0.141; 0.825; 0.675; 0.395. This results in the same limitations as in the example presented by Krzanowski (1987), i.e. the $W_n$ values are not monotonic. Therefore adopting the suggestion provided by Krzanowski (1987) where the $W_n$ values are ordered (11.019; 0.825; 0.675; 0.395; 0.395; 0.141;…) and the number of components should correspond to the number of $W_n$ values greater than 0.6 (Krzanowski 1987), we consider three principal components. I.e. we used an AMMI model with three interaction principal components.

The first two axes of the AMMI3 model are depicted in Figure 6 (left). As before, the environments with higher error variance tend to be placed away from the origin. A similar pattern was found by Gauch et al. (2011) where the environment OR91 was considered as an outlier.

**Figure 6.** Biplots for the first two axes of AMMI3 (left) and WAMMI3 (right) models, for the SxM yield data. The abscissa represents the first multiplicative term and the ordinate the second. The open dots represent the 150 genotypes and the codes for the 13 environments are defined in Table 3.

### c. Weighted AMMI analysis

Since the data in use is partially replicated, the cell means bring more information when there are two observations for the genotype environment combination. Therefore we have used the matrix of weights as defined by equation (7), with $m = max_j(1/\sigma^2_{\varepsilon_j})$, $\sigma^2_{\varepsilon_j}, j = 1, ..., 13$, is the error variance for environment $j$, $Nrep_{i,j}, i = 1, ... 150, j = 1, ..., 13$, is the number of replications for genotype $i$ in environment $j$, and $N_{rep} = 2$ is the maximum number of replications in the data set.

Figure 6 (right) shows the first two axes of the WAMMI3 model, weighted by $\boldsymbol{W}$ as in (7), and represents 75.7% of the total variance explained by the WAMMI3 model. As in the first example (simulated data), the environments with higher influence in the AMMI analysis have a more homogeneous distribution when using the WAMMI3 model (right hand side of Figure 6).

### d. AQ analysis and weighted AQ analysis

Figure 5 shows the QTL scans for the AMMI3 (red line) and WAMMI3 (green line) predicted values. The LOD scores show an increase when the QTL scans are made for the AMMI3 predicted values instead of the actual data. The same pattern is observed for most of
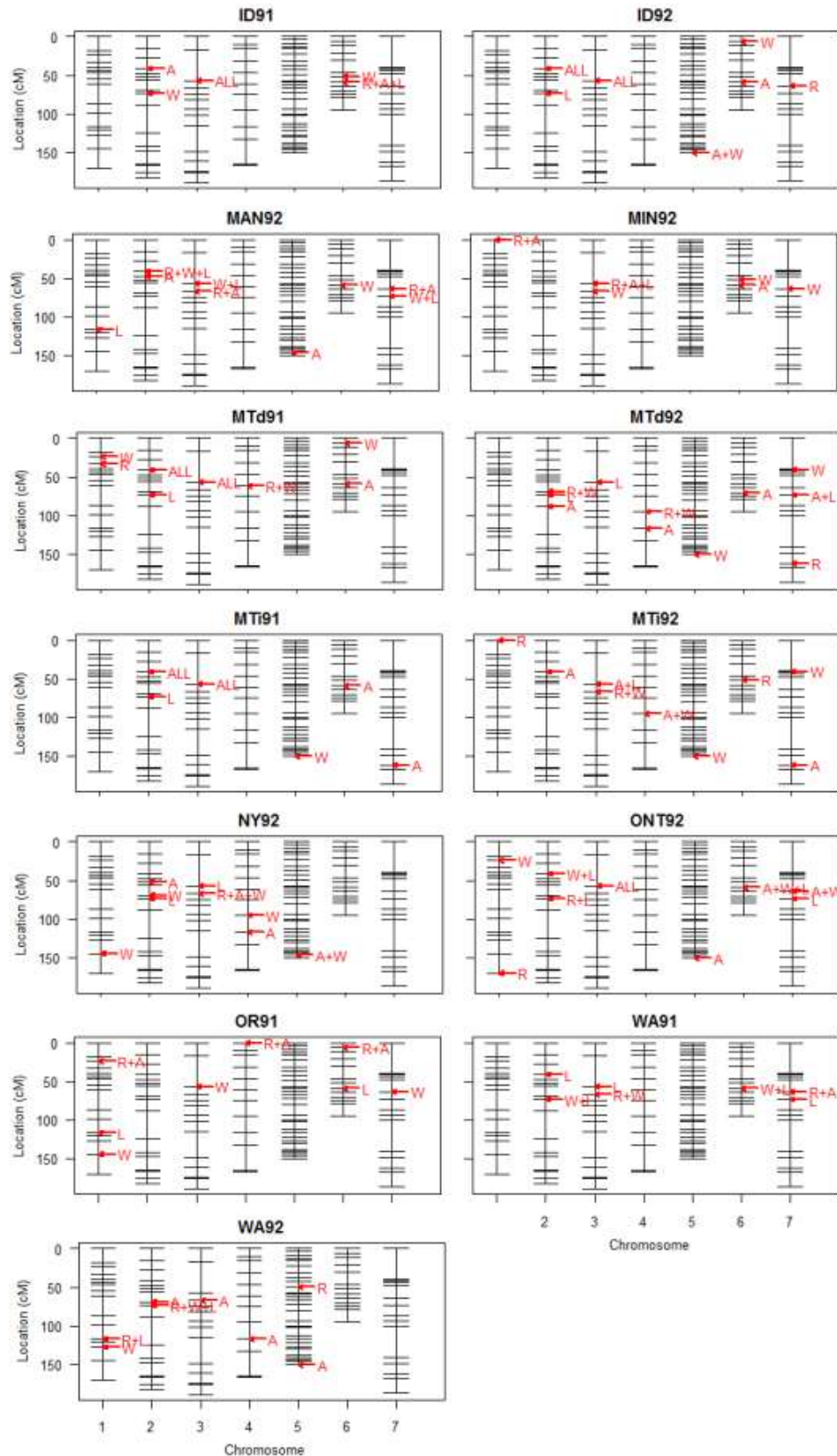
the environments. When the AQ analysis is replaced by the WAQ analysis the LOD scores become higher for the three environments with lowest LOD scores in the actual data and AMMI3 predicted values: OR91, NY92 and WA92. The two QTLs on chromosome 2 (Malosetti et al., 2004) are now visible in Figure 5 (red and green lines) for the most of the environments (more clear for WAQ analysis).

### e.  Weighted AQ analysis and comparison with QTL mixed linear models

Figure 7 presents a general comparison between the four approaches used here: direct QTL scans of the actual data; AQ analysis; WAQ analysis; and QTL mixed model framework. The exact positions can be found in Table S1. Most of the QTLs, detected with the WAQ analysis and the QTL mixed model framework , were either found in previous analyses or are very close to those QTLs (GAUCH *et al.* 2011; HAYES *et al.* 1993; LACAZE *et al.* 2009; LARSON *et al.* 1996; MALOSETTI *et al.* 2004; ROMAGOSA *et al.* 1999; ROMAGOSA *et al.* 1996; ZHU *et al.* 1999). Figure S6 shows the genome scan for the SxM yield data using the QTL mixed model framework.

On chromosome 1, between 116.7 and 170.1 cM, the linear mixed model identifies a QTL in two environments and WAQ analysis in four environments. For chromosome 2 two QTLs are identified by linear mixed models (in six and eight environments, respectively) and WAQ analysis (in five and five environments, respectively). The QTL on chromosome 3 is identified in 11 of the 13 environments by the linear mixed model and WAQ analysis. No QTL is detected by the linear mixed model on chromosomes 4 and 5. However, the WAQ analysis identifies a QTL on chromosome 4 at the same place as Lacaze et al. (2009) did. A QTL on chromosome 5 is detected for five environments by the WAQ analysis. Between 53.1 cM and 72.5 cM of chromosome 6 there is a QTL detection for linear mixed model and WAQ analysis, in five and four environments, respectively. A QTL is detected between 45.6 cM and 78.2 cM of chromosome 7 for linear mixed model and WAQ analysis, in seven and four environments, respectively.

When comparing the methods under study we can conclude again that the QTL scans of the actual data have the worst performance in terms of QTL detection (Figures 5 and 6). The AQ analysis finds more QTLs but the WAQ analysis is more similar with the results from the QTL mixed model analysis (Figure 6 and Table S1), which we believe are more credible.

**Figure 7.** Genetic map with the information of the place where a QTL was detected for each of the four approaches: QTL scans of the actual data (R); QTL scans of the AMMI3 predicted values (A); QTL scans of the WAMMI3 predicted values (W); and linear mixed model framework (L), for the SxM yield data in all 13 environments. "ALL" means that the QTL was detected with all the four approaches.

## 5  Discussion

### a.  Weighted AMMI analysis

The WAMMI model proposed here is a generalization of the standard AMMI analysis (Gauch 1992) that is able to account for heterogeneity of error variances across environments in a multiple-environment trial. This extension also allows the generalisation of the AQ analysis where the QTL scans are based on the AMMI predicted values (Gauch *et al.* 2011), which makes the results become similar to the often used QTL mixed model framework (Boer *et al.* 2007; Malosetti *et al.* 2008; Malosetti *et al.* 2004).

In this paper we used an algorithm based on the EM procedure proposed by Srebro and Jaakkola (2003) to conduct the weighted low-rank approximation. However, many alternatives can be found in the literature: maximum likelihood principal component analysis (Wentzell *et al.* 1997); a steepest descent algorithm and a Newton-like algorithm (Manton *et al.* 2003); and the use of a weighted rank correlation coefficient instead of the usual Pearson's (da Costa *et al.* 2011), among others. We chose the EM approach because of its easy implementation and good behaviour for our type of data.

### b.  AMMI model selection

Much research has been done about the choice of the "optimal" number of principal components in a PCA in general, and the number of multiplicative terms in the AMMI model in particular. Besides the cross-validation proposed by Krzanowski (1987) and used in this paper to decide on the number of multiplicative terms in the AMMI model, there are many options widely in AMMI literature. These include the signal to noise ratio (Gauch 1992), the Ockham's valley (Gauch 2006; MacKay 1992), a cross-validation (Gauch 1988; Gauch 1992; Piepho 1994) which can be performed with , e.g. the software MATMODEL version 3.0 (Gauch 2007), and significance tests. Gollob (1968) was the first proposing $F$ tests to help choosing the number of multiplicative terms. Other $F$ tests were also suggested: $F_{GH2}$ (Cornelius 1993; Cornelius *et al.* 1992), and $F_R$ (Piepho 1995). The cross-validation is usually seen as a conservative method because it refers to the modelling of a subset of the original data, which is expected to be less accurate than to model all data after the model choice (Annicchiarico 1997; Cornelius 1993). Although it is commonly used in AMMI literature, the cross-validation procedure proposed by Gauch (1992) is based on within experiment variation, which may not be the most suitable form of variation for a cross-validation procedure for a multiple environment data set.  Gollob's $F$ test is too liberal (Piepho 1997), $F_{GH2}$ is less conservative than $F_R$ (Annicchiarico 1997) and both tend to retain

a higher number of multiplicative terms than the cross-validation. However, these $F$ tests are used for determining the number of non-null multiplicative terms, which is different from finding the optimal number of terms for a prediction purposes (Cornelius 1993; Piepho 1997). The number of multiplicative terms for a predictive model should then be lower than the found significant by a significance $F$ test (Piepho 1997). Other alternatives widely used in PCA but not in AMMI modelling are parallel analysis (Horn 1965); minimum average partial (MAP) de Velicer (Velicer 1976), and very simple structure (VSS) (REVELLE and ROCKLIN 1979). These methods can only be applied to the two-way table of means, and should be applied to the multiplicative part of the data, i.e. after removing the genetic and environmental main effects.

The variety of possible methods is wide and so is the outcome. From an exhaustive analysis of related literature usually two or three axes are used to model the data because one component is (usually) not enough to capture the entire signal present in the data, and more than three components are already capturing a big amount of noise and are more difficult to visualize graphically. Moreover, in multi-environment trials, usually, there is no further information beyond two or three principal components.

### c. The influence of the heritability in the results

In field crops, the range of heritabilities is wide, and may vary from about 0.3 for yield in cereals in open environments (CLARKE and TOWNLEYSMITH 1986; SAEED $et$ $al.$ 2007), to more than 0.7 for tomato (Reif $et$ $al.$ 2009) or pepper (DO REGO $et$ $al.$ 2011; SOOD $et$ $al.$ 2009) in greenhouse experiments. When the heritability of the environments under study decreases the WAQ analysis tend to out-perform the QTL mixed model framework ($h^2$=0.3, Figures S2 and S3); whereas for higher heritability of the environments the QTL mixed model framework out-performs the WAQ analysis ($h^2$=0.8, Figures S4 and S5) but detects some QTLs which are likely to be false positives (e.g. a few detections on chromosome 12 and the very unlikely scenario of 10 QTLs found in several environments, Figure S5, Table 2). Both WAQ analysis and QTL mixed model framework out-perform AQ analysis, and all of them out-perform the QTL scans of the actual data. For all these comparison we should bear in mind that the thresholds for the WAQ analysis are fully comparable with the thresholds for AQ analysis and with the thresholds for the QTL scans on the actual data. However, because of the different methodologies and different software, the thresholds for WAQ analysis are not fully comparable with the QTL mixed model framework, but an approximation for illustration purposes. It should be remarked that the mixed model QTL mapping was used
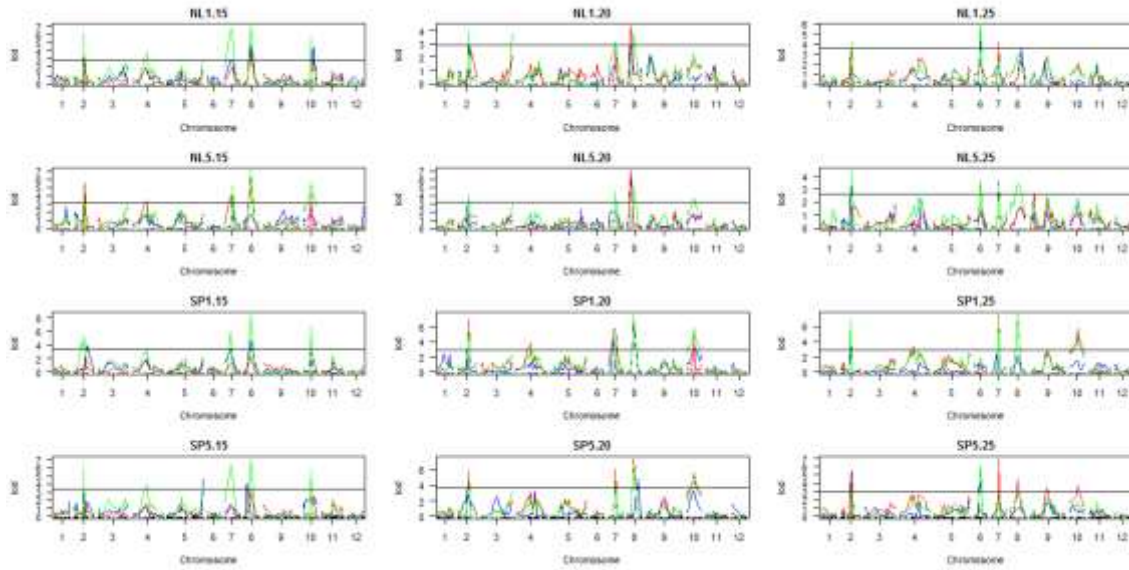
with default multiple testing corrections as set in Genstat 14 (Payne *et al.* 2011). Some playing around with those settings might have produced results closer to WAQ.
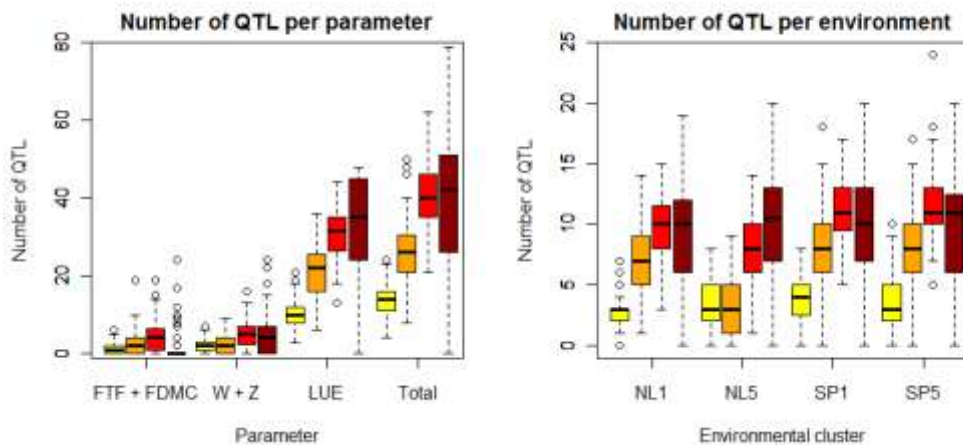
### d. Alternatives to the QTL mixed model methodology

Boer et al. (2007) suggested the possibility of having different methodologies performing as well as the QTL mixed model approach. They named Bayesian based methods and penalized regression as possibilities for similar analyses. The AQ analysis, i.e. use the AMMI expected values in the QTL scans across environments, first proposed by Gauch et al. (2011), is also an alternative to the QTL mixed model framework. However, the AQ analysis is not general enough to account for different error variances across environments which this paper generalizes by introducing the weighted SVD in the AMMI analysis.

The results presented in this paper are very encouraging because of several factors: (i) the WAQ analysis can be performed with the package qtl (BROMAN and SEN 2009) of the open source R software (Team 2009); (ii) the computation time to obtain the QTL scans and its summary is much shorter than the QTL mixed model framework in GenStat (Payne *et al.* 2011); and (iii) the results are very similar with the QTL mixed model output (Figures 3, 4, 7 and Table S1). It is also remarkable how the inclusion of the information about the error variances improves that much the results when the heritability of the trait/environment decreases (Figures 3, S2 and S4), comparing with the AQ analysis and the QTL scans of the actual data, and makes them very similar to the QTL mixed model methodology (Figure 3, 4, 7 and Table S1). So, the WAQ analysis is easy to apply with open source software and faster to run when compared with the QTL mixed linear model framework. Moreover, the WAMMI model and WAQ analysis are fully applicable to a wide range of fields such as plant breeding, crop sciences, genetics, microarray experiments (Crossa *et al.* 2005), rDNA studies (Adams *et al.* 2002); plant and microbial populations' growth across several environmental conditions (CULMAN *et al.* 2009; CULMAN *et al.* 2008) and animal sciences (BARHDADI and DUBE 2010).
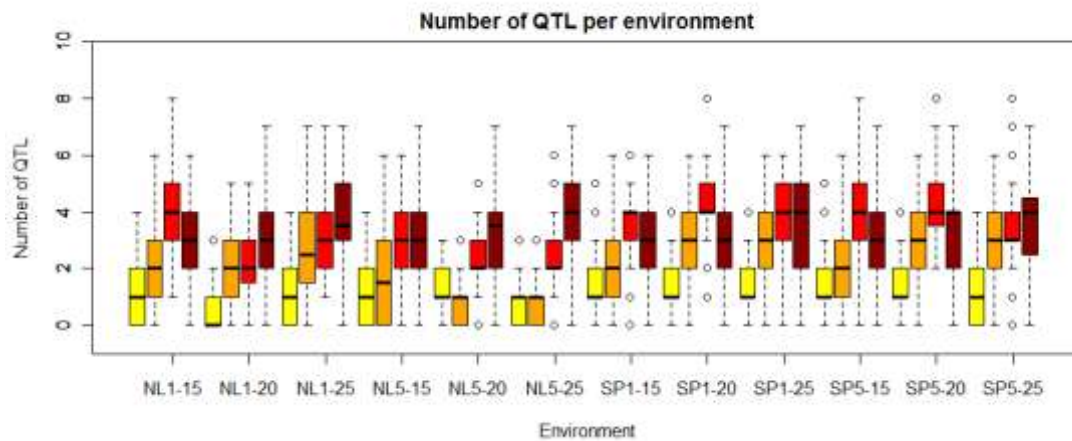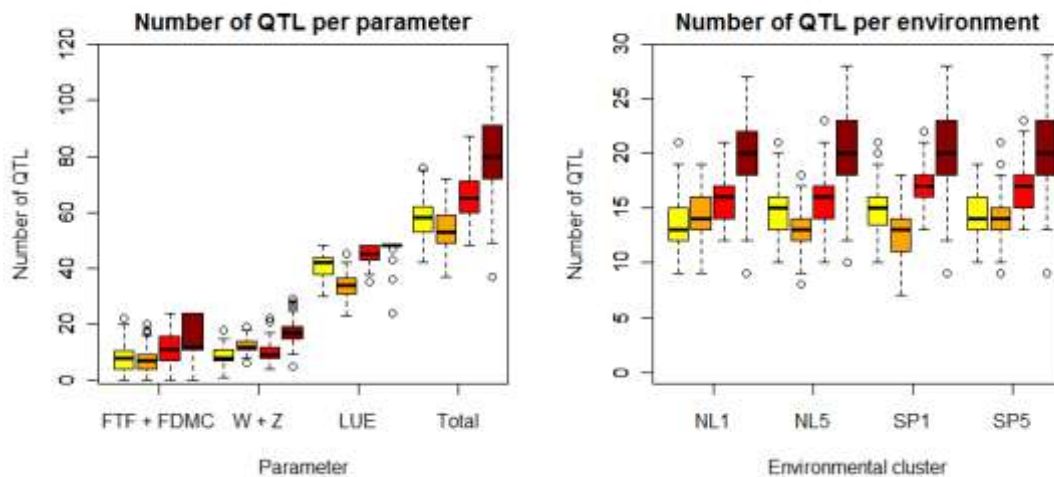
## 6    Supplementary material



**Figure S1.** QTL scans for the 12 environments of the yield data for pepper simulated from the physiological genotype-to-phenotype model with seven physiological parameters (Rodrigues *et al.* 2012). Each column represents a different level of temperature. The first row corresponds to the highest error variance in this realization and the second to the lowest. The black line represents the scans for the actual data, the blue for the AMMI2 predicted values, and the red for the WAMMI2 predicted values. All the scans are based on the composite interval mapping. The horizontal lines correspond to the thresholds for a LOD score of 3. These scans are based on one randomly chosen realization out of the 100 simulations.



**Figure S2.** Summary of the number of detected QTLs for the actual data (yellow), AMMI2 predicted values (orange), WAMMI2 predicted values (red) and linear mixed model (dark red). The graph on the left hand side shows the box plots for the number of QTLs per model parameter (Table 2), and on the right hand side the number of QTLs per environmental cluster (Table 1). These values are for QTLs detected when considering an interval of 20 cM centred in the right position. These plots are for a heritability of 0.3 in all environments.
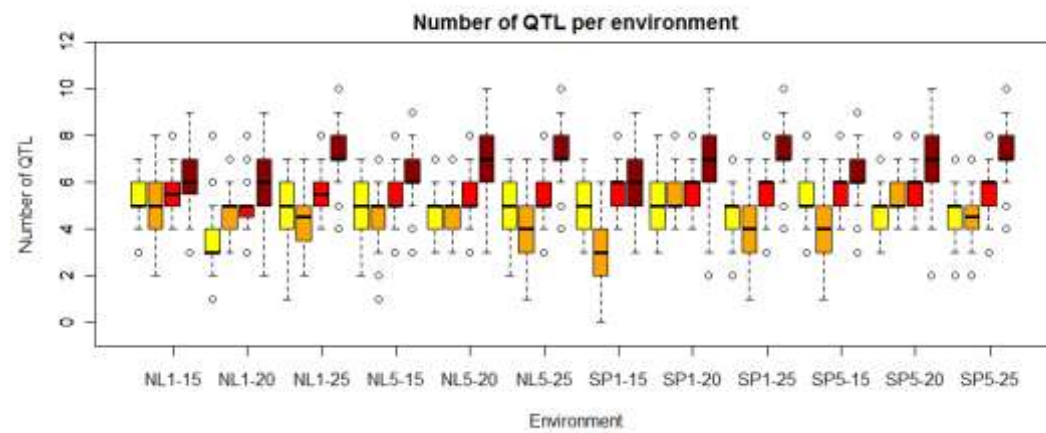
**Figure S3.** The top panel shows the number of QTLs detected per environment for a expected maximum of 7 (environments with temperature of 15ºC or 20ºC, Table 2) or 8 (environments with temperature of 25ºC, Table 2). The bottom panel shows the LOD scores per environment. The box plots are presented for the actual data (yellow), AMMI2 predicted values (orange), WAMMI2 predicted values (red) and linear mixed model (dark red). These plots are for an heritability of 0.3 in all environments.
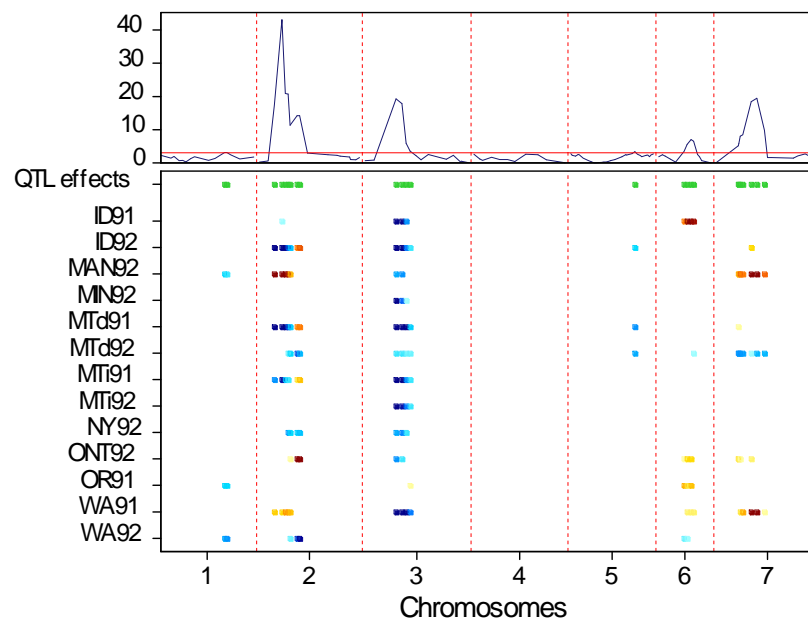


**Figure S4.** Summary of the number of detected QTLs for the actual data (yellow), AMMI2 predicted values (orange), WAMMI2 predicted values (red) and linear mixed model (dark red). The graph on the left hand side shows the box plots for the number of QTLs per model parameter (Table 2), and on the right hand side the number of QTLs per environmental cluster (Table 1). These values are for QTLs detected when considering an interval of 20 cM centred in the right position. These plots are for a heritability of 0.8 in all environments.

**Figure S5.** The top panel shows the number of QTLs detected per environment for a expected maximum of 7 (environments with temperature of 15ºC or 20ºC, Table 2) or 8 (environments with temperature of 25ºC, Table 2). The bottom panel shows the LOD scores per environment. The box plots are presented for the actual data (yellow), AMMI2 predicted values (orange), WAMMI2 predicted values (red) and linear mixed model (dark red). These plots are for an heritability of 0.8 in all environments.



**Figure S6.** Genome scan for the means of the SxM yield data. The $-\log_{10}(p)$-values for the QTL main effects plus QEI are shown. The red horizontal line is the 5% genomewide significance threshold. The green horizontal line in the bottom section summarizes the top panel. The environment specific QTL effects are shown. Blue (red) indicates that parent Steptoe (Morex) has significantly higher yield contribution. The considered variance-covariance (VCOV) structure was the factor analytic with two multiplicative terms.

**Table S1.** Chromosome (Chr) and respective positions (Pos) where a QTL was detected for each of the four approaches: QTL scans of the actual data (R); QTL scans of the AMMI3 predicted values (A); QTL scans of the WAMMI3 predicted values (W); and linear mixed model framework (L), for the SxM yield data in all 13 environments. "ALL" means that the QTL was detected with all the four approaches. The last column gives the reference where the same detection was observed. "None" indicates that no reference was found where a similar QTL was detected.

| Chr | Pos | ID91 | ID92 | MAN92 | MIN92 | MTd91 | MTd92 | MTi91 | MTi92 | NY92 | ONT92 | OR91 | WA91 | WA92 | Reference |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | [0; 33.5] | | | | R+A | R+W | | R | | | W | R+A | | | (Hayes *et al.* 1993) |
| 1 | [116.7; 170.1] | | | L | | | | | | W | W | W+L | | R+W+L | None |
| 2 | [41.2; 52.6] | A | ALL | ALL | | ALL | | ALL | A | A | W+L | | L | | (Gauch *et al.* 2011) (Hayes *et al.* 1993) (Lacaze *et al.* 2009) (Malosetti *et al.* 2004) (Romagosa *et al.* 1996) (Romagosa *et al.* 1999) (Zhu *et al.* 1999) |
| 2 | [68.8; 88.2] | W | L | | | L | R+W+L | L | | W+L | R+L | | W+L | ALL | (Gauch *et al.* 2011) (Malosetti *et al.* 2004) |
| 3 | [73; 83.6] | ALL | ALL | ALL | ALL | ALL | L | ALL | ALL | ALL | ALL | W | R+W+L | A | (Gauch *et al.* 2011) (Hayes *et al.* 1993) (Lacaze *et al.* 2009) (Larson *et al.* 1996) (Romagosa *et al.* 1999) |
| 4 | 1.4 | | | | | | | | | | | R+A | | | None |
| 4 | 63.2 | | | | | R+W | | | | | | | | | (Lacaze *et al.* 2009) |
| 4 | 96.5 | | | | | | R+W | | A+W | W | | | | | None |
| 4 | 118.3 | | | | | | A | | A | | | | | A | None |
| 5 | 49.6 | | | | | | | | | | | | | R | None |
| 5 | [146.3; 150.8] | | A+W | A | | | W | W | W | A+W | A | | | A | None |
| 6 | 8.1 | | W | | | W | | | | | | R+A | | | None |
| 6 | [53.1; 72.5] | ALL | A | W | A+W | A | A | A | R | | A+W+L | L | W+L | | (Gauch *et al.* 2011) (Romagosa *et al.* 1996) (Romagosa *et al.* 1999) |
| 7 | [45.6; 78.2] | | R | ALL | W | | A+W+L | | W | | A+W+L | W | ALL | | (Gauch *et al.* 2011) (Romagosa *et al.* 1996) |

**File S1.** The R code for the weighted low-rank SVD (SREBRO and JAAKKOLA 2003) (adapted from the Marlin's MatLab code in http://www.cs.toronto.edu/~marlin/code/wsvd.m).

```
# Imputs
# Y – (IxJ) data matrix
# W – (IxJ) weight matrix with 0 <= Wij <= 1; I = 1,…,I; j = 1,…,J
# N – rank of approximation
#
# Outputs
# U,D,V such that Y~UDV'

Y<- read.csv("data.csv", header=T)      # Read the original data set (IxJ)
X<- matrix(0,ngen,nenv)                 # Matrix (IxJ) with zero in all positions to initialize
the algorithm
aux<- matrix(1,ngen,nenv)
Xold=Inf*aux
Err=Inf                                 # Initial distance between consecutive iterations –
X(i) and X(i+1)
eps<- 1e-10                             # Maximum admissible distance between consecutive
iterations

while(Err>eps){                        # Repeats the code until the distance between X(i)
and X(i+1) is below eps=1e-10
   Xold=X                              # Update Xold to X(i)
   wsvd<- svd(W*Y + (1-W)*X)           # Weighted SVD
   U<- wsvd$u                          # Left singular vectors
   D<- diag(wsvd$d)                    # Singular values
   V<- wsvd$v                          # Right singular vectors
   D[(N+1):length(wsvd$d),(N+1):length(wsvd$d)]<- 0     # Discard singular values
above N
   X<- U %*% D %*% t(V)                                 # Update X (i.e. compute
X(i+1))
   Err=sum(sum((X-Xold)^2))                             # Update the distance between
consecutive iterations (Err)
}
```

## 7 References

ADAMS, G. C., N. T. WU and B. E. EISENBERG, 2002 Virulence and double-stranded RNA in Sphaeropsis sapinea. **Forest Pathology** 32: 309-329.

ANNICCHIARICO, P., 1997 Additive main effects and multiplicative interaction (AMMI) analysis of genotype-location interaction in variety trials repeated over years. **Theoretical and Applied Genetics** 94: 1072-1077.

BARCHI, L., J. BONNET, C. BOUDET, P. SIGNORET, I. NAGY *et al.*, 2007 A high-resolution, intraspecific linkage map of pepper (Capsicum annuum L.) and selection of reduced recombinant inbred line subsets for fast mapping. **Genome** 50: 51-60.

BARCHI, L., V. LEFEBVRE, A. M. SAGE-PALLOIX, S. LANTERI and A. PALLOIX, 2009 QTL analysis of plant development and fruit traits in pepper and performance of selective phenotyping. **Theoretical and Applied Genetics** 118: 1157-1171.

BARHDADI, A., and M. P. DUBE, 2010 Testing for Gene-Gene Interaction with AMMI Models. **Statistical Applications in Genetics and Molecular Biology** 9.

BOER, M. P., D. WRIGHT, L. Z. FENG, D. W. PODLICH, L. LUO *et al.*, 2007 A mixed-model quantitative trait loci (QTL) analysis for multiple-environment trial data using environmental covariables for QTL-by-environment interactions, with an example in maize. **Genetics** 177: 1801-1813.

BROMAN, K. W., and S. SEN, 2009 **A Guide to QTL Mapping with R/qtl.** Springer-Verlag, New York.

CLARKE, J. M., and T. F. TOWNLEYSMITH, 1986 Heritability and Relationship to Yield of Excised-Leaf Water-Retention in Durum-Wheat. **Crop Science** 26: 289-292.

CORNELIUS, P. L., 1993 Statistical Tests and Retention of Terms in the Additive Main Effects and Multiplicative Interaction-Model for Cultivar Trials. **Crop Science** 33: 1186-1193.

CORNELIUS, P. L., M. SEYEDSADR and J. CROSSA, 1992 Using the Shifted Multiplicative Model to Search for Separability in Crop Cultivar Trials. **Theoretical and Applied Genetics** 84: 161-172.

CROSSA, J., J. BURGUENO, D. AUTRAN, J. P. VIELLE-CALZADA, P. L. CORNELIUS *et al.*, 2005 Using linear-bilinear models for studying gene expression x teatment interaction in microarray experiments. **Journal of Agricultural Biological and Environmental Statistics** 10: 337-353.

CULMAN, S. W., R. BUKOWSKI, H. G. GAUCH, H. CADILLO-QUIROZ and D. H. BUCKLEY, 2009 T-REX: software for the processing and analysis of T-RFLP data. **BMC Bioinformatics** 10.

CULMAN, S. W., H. G. GAUCH, C. B. BLACKWOOD and J. E. THIES, 2008 Analysis of T-RFLP data using analysis of variance and ordination methods: A comparative study. **Journal of Microbiological Methods** 75: 55-63.

DA COSTA, J. F. P., H. ALONSO and L. ROQUE, 2011 A Weighted Principal Component Analysis and Its Application to Gene Expression Data. **Ieee-Acm Transactions on Computational Biology and Bioinformatics** 8: 246-252.

DO REGO, E. R., M. M. DO REGO, C. D. CRUZ, F. L. FINGER and V. W. D. CASALI, 2011 Phenotypic diversity, correlation and importance of variables for fruit quality and yield traits in Brazilian peppers (Capsicum baccatum). **Genetic Resources and Crop Evolution** 58: 909-918.

GABRIEL, K. R., and S. ZAMIR, 1979 Lower Rank Approximation of Matrices by Least-Squares with Any Choice of Weights. **Technometrics** 21: 489-498.

GAUCH, H. G., 1988 Model Selection and Validation for Yield Trials with Interaction. **Biometrics** 44: 705-715.

GAUCH, H. G., 1992 **Statistical analysis of regional yield trials: AMMI analysis of factorial designs.** Elsevier, Amsterdam.

GAUCH, H. G., 2006 Winning the accuracy game - Three statistical strategies - replicating, blocking and modeling - can help scientists improve accuracy and accelerate progress. **American Scientist** 94: 133-141.

GAUCH, H. G., 2007 **MATMODEL version 3.0: Open source software for AMMI and related analyses**, software available at http://www.css.cornell.edu/staff/gauch.

GAUCH, H. G., H. P. PIEPHO and P. ANNICCHIARICO, 2008 Statistical analysis of yield trials by AMMI and GGE: Further considerations. **Crop Science** 48: 866-889.

GAUCH, H. G., P. C. RODRIGUES, J. D. MUNKVOLD, E. L. HEFFNER and M. SORRELLS, 2011 Two New Strategies for Detecting and Understanding QTL x Environment Interactions. **Crop Science** 51: 96-113.

GOLLOB, H. F., 1968 A Statistical Model Which Combines Features of Factor Analysis and Analysis of Variance Techniques. **Psychometrika** 33: 73-115.

HAYES, P. M., F. Q. CHEN, A. KLEINHOFS, A. KILIAN and D. E. MATHER, 1996 Barley genome mapping  and its applications, pp. 229-249 in **Method of Genome Analysis in Plants**, edited by P. P. JAUHAR. CRC press, Boca Raton, Florida.

HAYES, P. M., B. H. LIU, S. J. KNAPP, F. CHEN, B. JONES *et al.*, 1993 Quantitative Trait Locus Effects and Environmental Interaction in a Sample of North-American Barley Germ Plasm. **Theoretical and Applied Genetics** 87: 392-401.

HORN, J., 1965 A rationale and test for the number of factors in factor analysis. **Psychometrika** 30: 179-185.

KRZANOWSKI, W. J., 1987 Cross-Validation in Principal Component Analysis. **Biometrics** 43: 575-584.

LACAZE, X., P. M. HAYES and A. KOROL, 2009 Genetics of phenotypic plasticity: QTL analysis in barley, Hordeum vulgare. **Heredity** 102: 163-173.

LARSON, S. R., D. KADYRZHANOVA, C. MCDONALD, M. SORRELLS and T. K. BLAKE, 1996 Evaluation of barley chromosome-3 yield QTLs in a backcross of F2 population using STS-PCR. **Theoretical and Applied Genetics** 93: 618-625.

MACKAY, D. J. C., 1992 Bayesian Interpolation. **Neural Computation** 4: 415-447.

MALOSETTI, M., J. M. RIBAUT, M. VARGAS, J. CROSSA and F. A. VAN EEUWIJK, 2008 A multi-trait multi-environment QTL mixed model with an application to drought and nitrogen stress trials in maize (Zea mays L.). **Euphytica** 161: 241-257.

MALOSETTI, M., J. VOLTAS, I. ROMAGOSA, S. E. ULLRICH and F. A. VAN EEUWIJK, 2004 Mixed models including environmental covariables for studying QTL by environment interaction. **Euphytica** 137: 139-145.

MANTON, J. H., R. MAHONY and Y. B. HUA, 2003 The geometry of weighted low-rank approximations. **Ieee Transactions on Signal Processing** 51: 500-514.

PAYNE, R. W., D. A. MURRAY, S. A. HARDING, D. B. BAIRD and D. M. SOUTAR, 2011 *An Introduction to GenStat for Windows (14th Edition).* VSN International, Hemel Hempstead, UK.

PIEPHO, H. P., 1994 Best Linear Unbiased Prediction (Blup) for Regional Yield Trials - a Comparison to Additive Main Effects and Multiplicative Interaction (Ammi) Analysis. **Theoretical and Applied Genetics** 89: 647-654.

PIEPHO, H. P., 1995 Robustness of statistical tests for multiplicative terms in the additive main effects and multiplicative interaction model for cultivar trials. **Theoretical and Applied Genetics** 90: 438-443.

PIEPHO, H. P., 1997 Analyzing genotype-environment data by mixed models with multiplicative terms. **Biometrics** 53: 761-766.

REIF, J. C., B. KUSTERER, H. P. PIEPHO, R. C. MEYER, T. ALTMANN *et al.*, 2009 Unraveling Epistasis With Triple Testcross Progenies of Near-Isogenic Lines. **Genetics** 181: 247-257.

REVELLE, W., and T. ROCKLIN, 1979 Very Simple Structure: an Alternative Procedure for Estimating the Optimal Number of Interpretable Factors. **Multivariate Behavioral Research** 14: 403-414.

RODRIGUES, P. C., E. HEUVELINK, M. C. A. M. BINK, L. F. M. MARCELIS and F. A. VAN EEUWIJK, 2012 A complex trait with unstable QTLs can follow from component traits with stable QTLs: an illustration by a simulation study in pepper. (to be submited).

ROMAGOSA, I., F. HAN, S. E. ULLRICH, P. M. HAYES and D. M. WESENBERG, 1999 Verification of yield QTL through realized molecular marker-assisted selection responses in a barley cross. **Molecular Breeding** 5: 143-152.

ROMAGOSA, I., S. E. ULLRICH, F. HAN and P. M. HAYES, 1996 Use of the additive main effects and multiplicative interaction model in QTL mapping for adaptation in barley. **Theoretical and Applied Genetics** 93: 30-37.

SAEED, A., K. HAYAT, A. A. KHAN, S. IQBAL and G. ABAS, 2007 Assessment of Genetic Variability and Heritability in Lycopersicon esculentum Mill. **International Journal of Agriculture and Biology** 9: 375-377.

SOOD, S., R. SOOD, V. SAGAR and K. C. SHARMA, 2009 Genetic Variation and Association Analysis for Fruit Yield, Agronomic and Quality Characters in Bell Pepper. **International Journal of Vegetable Science** 15: 272-284.

SREBRO, N., and T. JAAKKOLA, 2003 Weighted Low-Rank Approximations, pp. 600-607 in **Twentieth International Conference on Machine Learning** *(ICML-2003)*, edited by T. FAWCETT and N. MISHRA. The AAAI Press, Menlo Park, California, Washington DC.

TEAM, R. D. C., 2009 **R: A Language and Environment for Statistical Computing**, Vienna, Austria.

VELICER, W., 1976 Determining the number of components from the matrix of partial correlations. **Psychometrika** 41: 321-327.

VOLTAS, J., F. VAN EEUWIJK, E. IGARTUA, L. F. GARCÍA DEL MORAL, J. L. MOLINA-CANO *et al.*, 2002 Genotype by environment interaction and adaptation in barley breeding: Basic concepts and methods of analysis, pp. 205-241 in **Barley science: Recent advances from molecular biology to agronomy of yield and quality**, edited by G. A. SLAFER, J. L. MOLINA-CANO, R. SAVIN, J. L. ARAUS and I. ROMAGOSA. Haworth Press, New York.

WENTZELL, P. D., D. T. ANDREWS, D. C. HAMILTON, K. FABER and B. R. KOWALSKI, 1997 Maximum likelihood principal component analysis. **Journal of Chemometrics** 11: 339-366.

YANG, R. C., J. CROSSA, P. L. CORNELIUS and J. BURGUENO, 2009 Biplot Analysis of Genotype x Environment Interaction: Proceed with Caution. **Crop Science** 49: 1564-1576.

ZHU, H., G. BRICENO, R. DOVEL, P. M. HAYES, B. H. LIU *et al.*, 1999 Molecular breeding for grain yield in barley: an evaluation of QTL effects in a spring barley cross. **Theoretical and Applied Genetics** 98: 772-779.