

How do the type of QTL effect and the form of the residual term influence QTL detection in multi-parent populations? A case study in the maize EU-NAM population - Supplementary material

Vincent Garin · Valentin Wimmer ·
Sofiane Mezmouk · Marcos Malosetti ·
Fred van Eeuwijk

the date of receipt and acceptance should be inserted later

Vincent Garin
Biometris, Wageningen University and research Center, P.O Box 100, 6700AC Wageningen,
The Netherlands.
Tel.: 0031 317 83 574
E-mail: vincent.garin@wur.nl

Valentin Wimmer
KWS SAAT SE, Einbeck, Germany

S1: Nei and Li genotype similarity coefficient

$GS = 2N_{ij}/(N_i + N_j)$ where N_{ij} is the number of common sites between i and j and $N_{i(j)}$ is the total number of observed sites in individual $i(j)$ (Nei and Li, 1979). In our situation, considering the same number of SNPs in parent i and j makes GS equivalent to the simple matching coefficient (SM).

S2: Simple matching coefficients table between parents of the EU-NAM Dent panel

	B73	D06	D09	EC169	F252	F353	F618	Mo17	UH250	UH304	W117
B73	1	0.566	0.556	0.685	0.561	0.557	0.59	0.554	0.579	0.584	0.542
D06		1	0.866	0.643	0.611	0.586	0.599	0.559	0.836	0.624	0.568
D09			1	0.585	0.628	0.618	0.589	0.559	0.762	0.648	0.564
EC169				1	0.556	0.558	0.589	0.558	0.613	0.591	0.538
F252					1	0.633	0.578	0.576	0.609	0.648	0.579
F353						1	0.589	0.567	0.575	0.761	0.565
F618							1	0.552	0.595	0.606	0.559
Mo17								1	0.578	0.586	0.564
UH250									1	0.614	0.563
UH304										1	0.587
W117											1

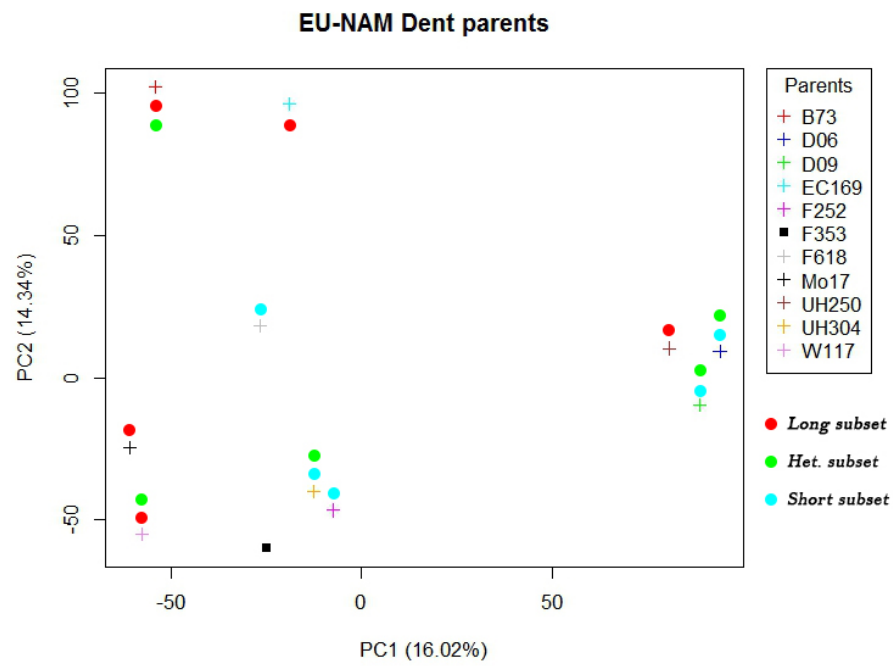
Average genetic similarity score per subset

$$\bar{S}M_{short} = 0.639$$

$$\bar{S}M_{het.} = 0.613$$

$$\bar{S}M_{long} = 0.573$$

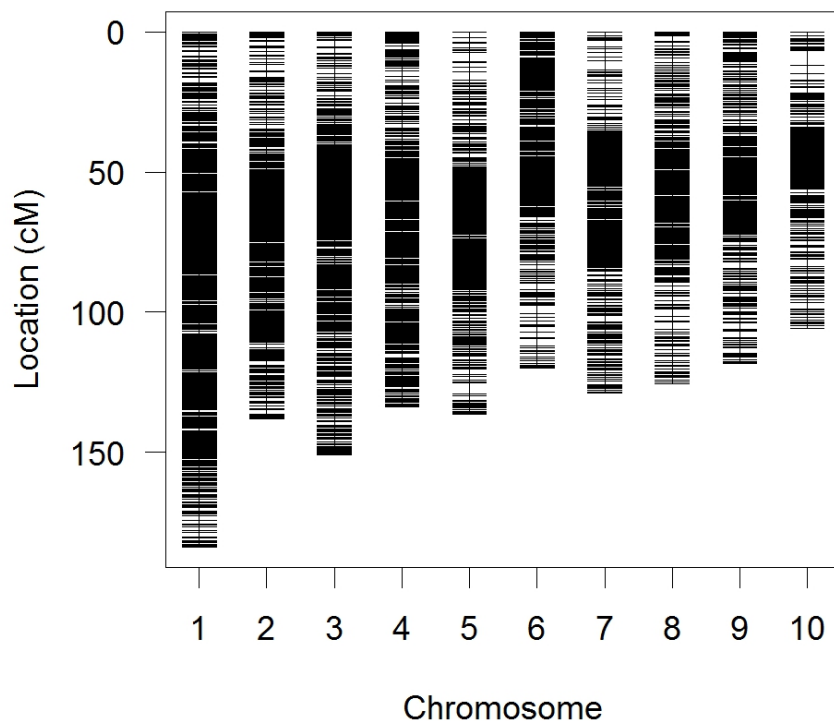
S3: Principal component bi-plot of the EU-NAM Dent parents



S4: Genetic marker map of the different subsets

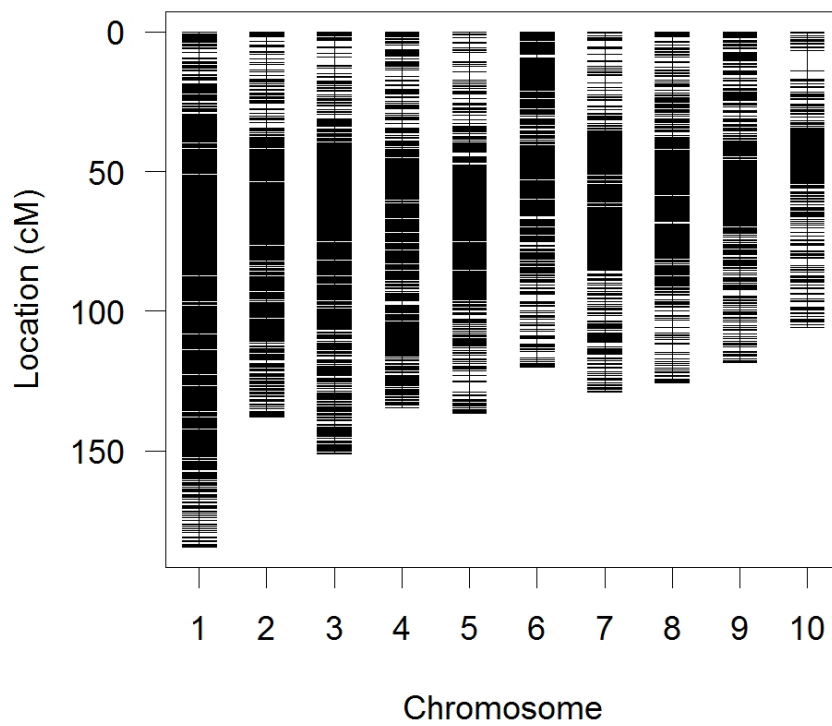
Short subset

Chromosome	N	Length(cM)	Average spacing	maximum spacing
1	985	183.8	0.2	1.8
2	633	137.9	0.2	2.0
3	702	150.9	0.2	2.4
4	617	133.9	0.2	2.2
5	563	136.5	0.2	3.8
6	482	119.9	0.2	3.1
7	470	128.9	0.3	3.3
8	496	125.6	0.3	2.6
9	449	118.4	0.3	2.4
10	340	105.8	0.3	5.3
Overall	5737	1341.6	0.2	5.3

EU-NAM Dent short subset map

Heterogeneous subset

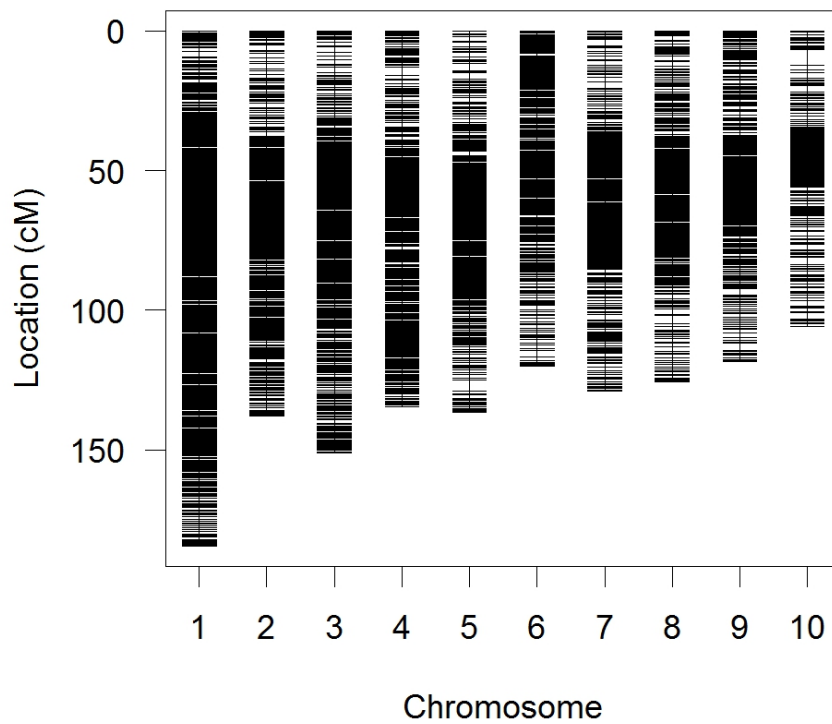
Chromosome	N	Length(cM)	Average spacing	maximum spacing
1	1055	184.5	0.2	1.9
2	630	137.9	0.2	2.1
3	711	151.0	0.2	2.9
4	620	134.6	0.2	2.2
5	573	136.6	0.2	3.8
6	526	119.9	0.2	2.6
7	493	128.9	0.3	2.6
8	527	125.6	0.2	2.9
9	467	118.4	0.3	2.4
10	332	105.8	0.3	7.3
Overall	5934	1343.2	0.2	7.3

EU-NAM Dent hetero subset map

Long subset

Chromosome	N	Length(cM)	Average spacing	maximum spacing
1	1121	184.5	0.2	1.9
2	641	137.9	0.2	2.0
3	743	151.0	0.2	1.9
4	647	134.6	0.2	3.1
5	598	136.6	0.2	4.1
6	550	119.9	0.2	2.1
7	524	128.9	0.2	2.2
8	560	125.6	0.2	2.9
9	484	118.4	0.2	2.7
10	344	105.8	0.3	5.7
Overall	6212	1343.2	0.2	5.7

EU-NAM Dent long subset map



S5: Test statistic of the QTL effect

Wald test derivation

The significance of the estimated QTL effects $\hat{\beta}_Q$ can be estimated using the Wald test (Wald, 1943). From model (2), the phenotype values \mathbf{y} have an expectation equivalent to $\mathbf{X}\beta$ and their variance is \mathbf{R} (McCulloch and Searle, 2001, 6.5, 6.6). In such a situation, we can derive a generalized estimate for β and its variance as follow (Rao et al, 2008, 4.65, 4.66):

$$\hat{\beta} = (\mathbf{X}'\mathbf{R}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{R}^{-1}\mathbf{y} \quad (1)$$

$$V(\hat{\beta}) = (\mathbf{X}'\mathbf{R}^{-1}\mathbf{X})^{-1} \quad (2)$$

In its general form, the Wald statistic is equal to (McCulloch and Searle, 2001, 5.39):

$$W = (\hat{\beta} - \hat{\beta}_0)'[V(\hat{\beta})]^{-1}(\hat{\beta} - \hat{\beta}_0) \quad (3)$$

Under the null hypothesis we assume $\hat{\beta} = \hat{\beta}_0 = [0 \ 0 \ \dots \ 0]'$. After substituting (1) and (2) in (3), we can rewrite the Wald statistic like that:

$$\begin{aligned} W &= \mathbf{y}'\hat{\mathbf{R}}^{-1}\mathbf{X}(\mathbf{X}'\hat{\mathbf{R}}^{-1}\mathbf{X})^{-1}\mathbf{X}'\hat{\mathbf{R}}^{-1}\mathbf{y} \\ &= \mathbf{y}'\hat{\mathbf{R}}^{-1}\mathbf{X}\hat{\beta} \\ &= \mathbf{y}'\hat{\mathbf{R}}^{-1}\mathbf{H}\mathbf{y} \\ &= \mathbf{y}'\hat{\mathbf{R}}^{-1}\hat{\mathbf{y}} \end{aligned} \quad (4)$$

where,

$$\mathbf{H} = \mathbf{X}(\mathbf{X}'\hat{\mathbf{R}}^{-1}\mathbf{X})^{-1}\mathbf{X}'\hat{\mathbf{R}}^{-1}$$

is the generalized hat matrix.

The previous expression W represent a global Wald test coefficient $W(\hat{\beta})$ including both effect of the cross intercepts $\hat{\beta}_c$ and the QTL effects $\hat{\beta}_Q$. The significance of the QTL effects $W(\hat{\beta}_Q)$ can be obtained doing the difference between $W(\hat{\beta})$ and $W(\hat{\beta}_c)$, the Wald statistic of a model including only the cross intercept terms.

$$W(\hat{\beta}_Q) = W(\hat{\beta}) - W(\hat{\beta}_c) \quad (5)$$

$W(\hat{\beta}_Q)$ is therefore proportional to $W(\hat{\beta})$ described in (4). $W(\hat{\beta}_Q)$ tests the hypothesis of all QTL effects equal zero versus at least on component of $\hat{\beta}_Q$ being non null. $W(\hat{\beta}_Q)$ follows a χ^2 distribution with degree of freedom equal to the rank of \mathbf{X}_Q (the number of estimated QTL effects).

Interpretation

Expression (4) allows to summarise the main features of the proposed QTL models: parsimony versus goodness of fit and accurate form of the residual term.

Parsimony vs goodness of fit

Since the degree of freedom of the Wald test depends on the number of estimated parameters, more parsimonious models like the ancestral or the bi-allelic model will automatically increase the significance level of $W(\hat{\beta}_Q)$. The use of more parsimonious models is however not always a guaranty of better results (Bardol et al, 2013; Steinhoff et al, 2011). An important criteria to be balanced with parsimony is the necessity to infer allele effects that capture correctly the trait variability. To illustrate this, we can further reduce expression $W(\hat{\beta})$ (4) and substitute it in (5) to draw the following relationship:

$$W(\hat{\beta}_Q) \propto \mathbf{y}'\hat{\mathbf{y}} = \sum_{n=1}^N y_n \hat{y}_n \quad (6)$$

From this expression we can see that the more the vector \mathbf{y} and the vector $\hat{\mathbf{y}}$ vary in the same direction the higher will be $W(\hat{\beta}_Q)$. This simply means that the form of the QTL incidence matrix \mathbf{X}_Q should be chosen to give genetic estimates $\hat{\beta}_Q$ allowing a projection $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ capturing the highest proportion of the trait variability. If these variations are due to parental or cross-specific effects, corresponding genetic effect estimates $\hat{\beta}_Q$ capturing these local variation should perform better at the price of a higher number of parameter to estimate. On the other hand if the effects are similar through the MPP, a reduced number of parameters will capture this variability and allows gains in power by a lower number of degree of freedom.

Accurate form of the residual term

The test statistic is also influenced by the chosen VCOV. As we can see in expression (4) each element composing $W(\hat{\beta}_Q)$ is weighted by the estimated $\hat{\mathbf{R}}$. The more $\hat{\mathbf{R}}$ will reflect the correct form of the residual term the more accurate will be the QTL detection process. In the HRT case, substituting (4) in (5), we can write the following relationship:

$$W(\hat{\beta}_Q) \propto \sum_{n=1}^N \frac{y_n \hat{y}_n}{\sigma_r^2} \quad (7)$$

In this situation each elements is weighted by an average level of uncertainty σ_r^2 which may not be representative of crosses particularities. In the CSRT situation, we have:

$$W(\hat{\beta}_Q) \propto \sum_{c=1}^{n_c} \sum_{n=1}^{N_c} \frac{y_n \hat{y}_n}{\sigma_{r_c}^2} \quad (8)$$

Here, the different elements will be weighted by cross-specific variance residual terms $\sigma_{r_c}^2$ that take into account the potential differences of uncertainty between crosses. The more $\sigma_{r_c}^2$ are heterogeneous the more a CSRT model is needed to

handle this variability. In such case, a HRT model will tend to be more liberal when uncertainty is in fact larger (cross 1 in Table 1) and more conservative when effects are in reality more certain (cross 3 in Table 1).

Table 1 Illustration of the difference between the HRT and CSRT assumption in an heterogeneous MPP and its effect on the QTL test statistic ($\frac{\beta}{\sigma^2}$).

	$\sigma_{r_c}^2$ ("True")	σ_r^2 ("Average")	Test (HRT)		Test (CSRT)
cross 1	190		$\frac{\beta_1}{100}$	>	$\frac{\beta_1}{190}$
cross 2	100	100	$\frac{\beta_2}{100}$	=	$\frac{\beta_2}{100}$
cross 3	10		$\frac{\beta_3}{100}$	<	$\frac{\beta_3}{10}$

S6: CSRT model approximation

The idea is to first estimate the variance covariance structure $\hat{\mathbf{R}}$ and then use it in a generalized estimate of the Wald test (4):

$$W(\hat{\beta}) = \mathbf{y}' \hat{\mathbf{R}}^{-1} \mathbf{X} (\mathbf{X}' \hat{\mathbf{R}}^{-1} \mathbf{X})^{-1} \mathbf{X}' \hat{\mathbf{R}}^{-1} \mathbf{y} \quad (9)$$

Since in the CSRT model $\hat{\mathbf{R}}$ contains only diagonal elements we can simply invert $\hat{\mathbf{R}}$ by doing $\hat{\mathbf{R}}^{-1} = 1/\hat{\mathbf{R}}$.

SIM

$$\text{Cr(Int) model:} \quad \mathbf{y} = \mathbf{X}_c \beta_c + \mathbf{r} \quad (10)$$

$$\text{Cr(Int) + Q model:} \quad \mathbf{y} = \mathbf{X}_c \beta_c + \mathbf{X}_Q \beta_Q + \mathbf{r} \quad (11)$$

To estimate the significance of the QTL effect we can use the following incremental Wald statistics

$$W(\hat{\beta}_Q) = W([\hat{\beta}_c | \hat{\beta}_Q]) - W(\hat{\beta}_c) \quad (12)$$

The procedure to estimate the significance of the QTL effect genome-wide is the following:

1. Estimate $\hat{\mathbf{R}}$ from model (10).
2. Compute $W(\hat{\beta}_c)$ substituting $\hat{\mathbf{R}}$ and $\mathbf{X} = \mathbf{X}_c$ in (9).
3. Compute at each position $W([\hat{\beta}_c | \hat{\beta}_Q])$ substituting $\hat{\mathbf{R}}$ and $\mathbf{X} = [\mathbf{X}_c | \mathbf{X}_Q]$ in (9).
4. Compute the p-value of the QTL effect using (12) and $W(\hat{\beta}_Q) \sim \chi_{df}^2$ with $df = \text{Rank}(\mathbf{X}_Q)$.

CIM

$$\text{Cr(Int) + cof model:} \quad \mathbf{y} = \mathbf{X}_c\boldsymbol{\beta}_c + \mathbf{X}_{cof}\boldsymbol{\beta}_{cof} + \mathbf{r} \quad (13)$$

$$\text{Cr(Int) + cof + Q model:} \quad \mathbf{y} = \mathbf{X}_c\boldsymbol{\beta}_c + \mathbf{X}_{cof}\boldsymbol{\beta}_{cof} + \mathbf{X}_Q\boldsymbol{\beta}_Q + \mathbf{r} \quad (14)$$

To estimate the significance of the QTL effect we can use the following incremental Wald statistics

$$W(\hat{\boldsymbol{\beta}}_Q) = W([\hat{\boldsymbol{\beta}}_c|\hat{\boldsymbol{\beta}}_{cof}|\hat{\boldsymbol{\beta}}_Q]) - W([\hat{\boldsymbol{\beta}}_c|\hat{\boldsymbol{\beta}}_{cof}]) \quad (15)$$

The procedure to estimate the significance of the QTL effect genome-wide is the following:

1. Estimate $\hat{\mathbf{R}}_j$ for the different combinations of cofactor $\mathbf{X}_{cof,j}$ using model (13).
2. Compute $W([\hat{\boldsymbol{\beta}}_c|\hat{\boldsymbol{\beta}}_{cof,j}])$ for the different combinations of cofactor substituting $\hat{\mathbf{R}}_j$ and $\mathbf{X} = [\mathbf{X}_c|\mathbf{X}_{cof,j}]$ in (9).
3. Compute at each position $W([\hat{\boldsymbol{\beta}}_c|\hat{\boldsymbol{\beta}}_{cof,j}|\hat{\boldsymbol{\beta}}_Q])$ substituting $\hat{\mathbf{R}}_j$ and $\mathbf{X} = [\mathbf{X}_c|\mathbf{X}_{cof,j}|\mathbf{X}_Q]$ in (9).
4. Compute the p-value of the QTL effect using (15) and $W(\hat{\boldsymbol{\beta}}_Q) \sim \chi_{df}^2$ with $df = \text{Rank}(\mathbf{X}_Q)$.

S7: Multi QTL effects model

The multi-QTL effect model uses a forward regression to build up a model in which different loci are allowed to have different types of QTL effects. At each step one QTL is added, and each of the different types of effects are compared. To select this position, we compute genome wide profiles using the types of QTL effects given by the user. Each of these profiles uses a single type of effect for all the tested (QTL) position.

$$\mathbf{y} = \mathbf{X}_c\beta_c + \mathbf{X}_{Q1}\beta_{Q1} + \mathbf{r} \quad (16)$$

Where, the QTL position $\mathbf{X}_{Q1}\beta_{Q1}$ is parental, ancestral, or bi-allelic. We calculate therefore as many QTL profile as types of QTL effects chosen by the user. From each of these profiles, the most significant position based on the $-\log_{10}(p)$ value statistic is selected (e.g., $\mathbf{X}_{Q1.par}$, $\mathbf{X}_{Q1.anc}$, $\mathbf{X}_{Q1.biall}$). Note that the selected QTL positions might or not be at the same position. The one that increases the most the R_{adj}^2 is selected as QTL. The selected position with its type of QTL effect is added to the cofactors list and the selection process starts again. If at step 1 we selected a bi-allelic QTL, then at step 2 the QTL profiles will be based on the following models:

$$\mathbf{y} = \mathbf{X}_c\beta_c + \mathbf{X}_{q1.biall}\beta_{q1} + \mathbf{X}_{Q2}\beta_{Q2} + \mathbf{r} \quad (17)$$

With again the tested QTL position $\mathbf{X}_{Q2}\beta_{Q2}$ taking, in each QTL profile, the form of one of the QTL effect specified by the user.

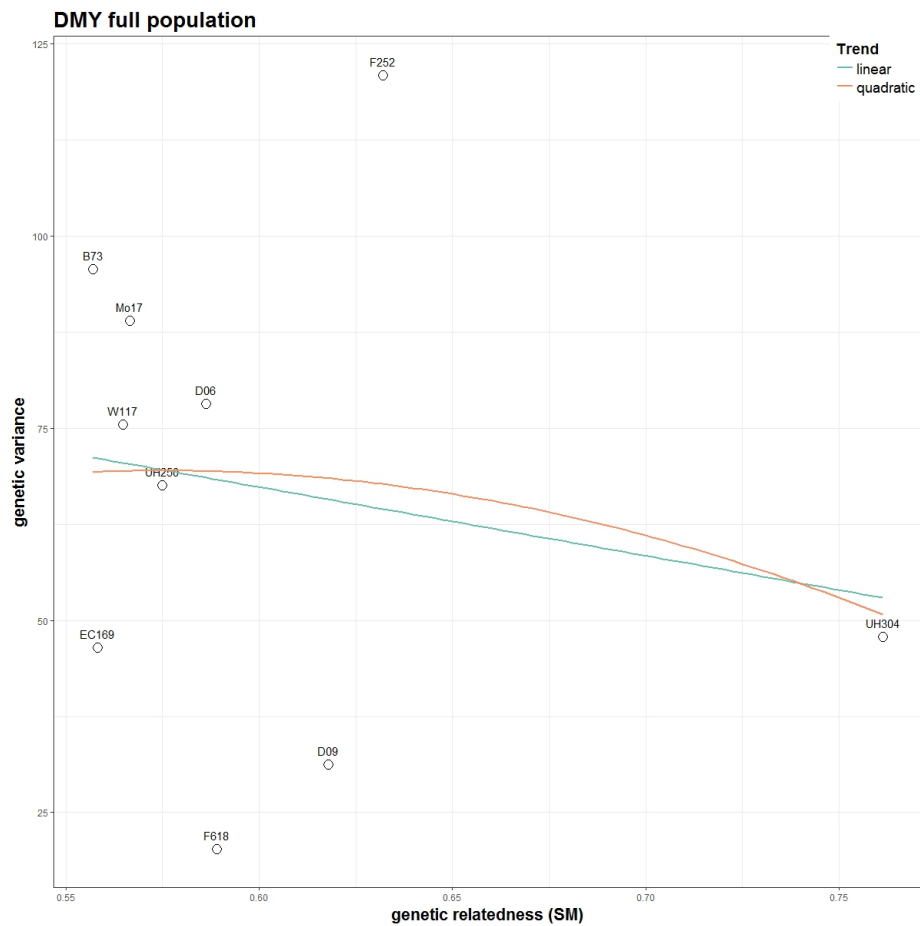
The procedure stops when there is no more significant position. The final list of QTL is tested simultaneously using a backward elimination. The final model could look like that:

$$\mathbf{y} = \mathbf{X}_c\beta_c + \mathbf{X}_{q1.biall}\beta_{q1} + \dots + \mathbf{X}_{q(t-1).par}\beta_{q(t-1)} + \mathbf{X}_{Qt.anc}\beta_{Qt} + \mathbf{r} \quad (18)$$

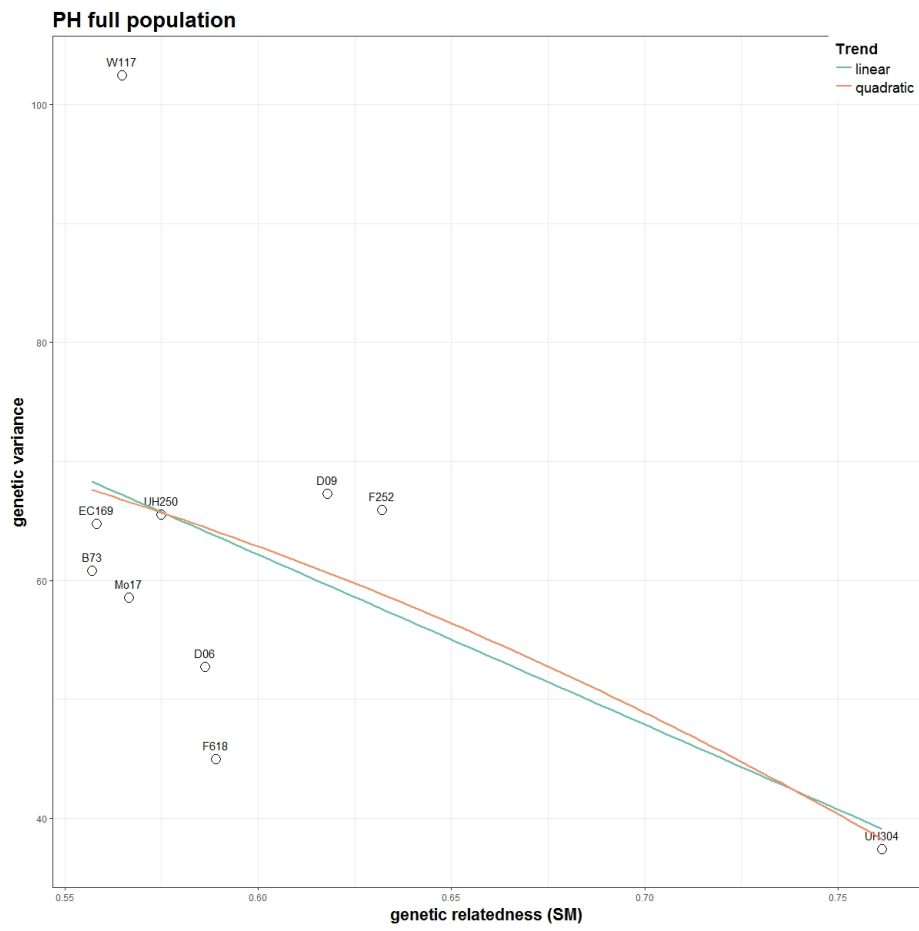
S8: Genetic variance versus parental relatedness tables

The following graphs represent the scatter plots of estimated genetic variance per family on SM coefficient between the central parent and the peripheral one (see S1 and S2). The two trends represent the linear and quadratic trend respectively.

Dry matter yield



Plant height



S9: Permutation threshold results

Significance thresholds determined by 1000 genome-wide permutations taking the $-\log_{10}(\text{p-value})$ of the upper 95% Wald statistic under the empirical null distribution as the critical value for rejection.

Dry mater yield

	parental	ancestral	bi-allelic	Average - MQE
short	3.89	4.09	4.66	4.21
het.	4.03	4.27	5.01	4.44
long	3.82	3.97	4.56	4.12
Average	3.91	4.11	4.74	

Plant height

	parental	ancestral	bi-allelic	Average - MQE
short	3.86	4.03	4.65	4.18
het.	3.97	4.22	5.18	4.46
long	3.77	4.14	4.85	4.25
Average	3.87	4.13	4.89	

References

- Bardol N, Ventelon M, Mangin B, Jasson S, Loywick V, Couton F, Derue C, Blanchard P, Charcosset A, Moreau L (2013) Combined linkage and linkage disequilibrium qtl mapping in multiple families of maize (*zea mays* l.) line crosses highlights complementarities between models based on parental haplotype and single locus polymorphism. *Theoretical and applied genetics* 126(11):2717–2736
- McCulloch CE, Searle SR (2001) *Generalized, linear, and mixed models*. Wiley Online Library
- Nei M, Li WH (1979) Mathematical model for studying genetic variation in terms of restriction endonucleases. *Proceedings of the National Academy of Sciences* 76(10):5269–5273
- Rao CR, Toutenburg H, Shalabh, Heumann C (2008) *Linear Models and Generalizations: Least Squares and Alternatives*. Springer
- Steinhoff J, Liu W, Maurer HP, Würschum T, Friedrich C, Longin H, Ranc N, Reif JC (2011) Multiple-line cross quantitative trait locus mapping in european elite maize. *Crop science* 51(6):2505–2516
- Wald A (1943) Tests of statistical hypotheses concerning several parameters when the number of observations is large. *Transactions of the American Mathematical society* 54(3):426–482