




Omics-based hybrid prediction in maize

Matthias Westhues¹ · Tobias A. Schrag¹ · Claas Heuer^{2,3} · Georg Thaller² · H. Friedrich Utz¹ · Wolfgang Schipprack¹ · Alexander Thiemann⁴ · Felix Seifert⁴ · Anita Ehret² · Armin Schlereth⁵ · Mark Stitt⁵ · Zoran Nikoloski⁵ · Lothar Willmitzer⁵  · Chris C. Schön⁶ · Stefan Scholten⁴  · Albrecht E. Melchinger¹ 

Received: 9 May 2017 / Accepted: 9 June 2017
© Springer-Verlag GmbH Germany 2017

Abstract

Key message Complementing genomic data with other “omics” predictors can increase the probability of success for predicting the best hybrid combinations using complex agronomic traits.

Abstract Accurate prediction of traits with complex genetic architecture is crucial for selecting superior candidates in animal and plant breeding and for guiding decisions in personalized medicine. Whole-genome prediction has revolutionized these areas but has inherent limitations in incorporating intricate epistatic interactions. Downstream “omics” data are expected to integrate interactions

within and between different biological strata and provide the opportunity to improve trait prediction. Yet, predicting traits from parents to progeny has not been addressed by a combination of “omics” data. Here, we evaluate several “omics” predictors—genomic, transcriptomic and metabolic data—measured on parent lines at early developmental stages and demonstrate that the integration of transcriptomic with genomic data leads to higher success rates in the correct prediction of untested hybrid combinations in maize. Despite the high predictive ability of genomic data, transcriptomic data alone outperformed them and other predictors for the most complex heterotic trait, dry matter yield. An eQTL analysis revealed that transcriptomic data integrate genomic information from both, adjacent and distant sites relative to the expressed genes. Together, these findings suggest that downstream predictors capture physiological epistasis that is transmitted from parents to their hybrid offspring. We conclude that the use of downstream “omics” data in prediction can exploit important information beyond structural genomics for leveraging the efficiency of hybrid breeding.

Communicated by Matthias Frisch.

The authors acknowledge support by the state of Baden–Württemberg through bwHPC. This project was funded by the German Federal Ministry of Education and Research (BMBF) within the projects OPTIMAL (FKZ: 0315958B, 0315958F), SYNBIOT (FKZ: 0315528D) and by the German Research Foundation (DFG, Grants No. ME 2260/5-1 and SCHO 764/6-1). Financial support for M.W. was provided by the Fiat Panis foundation, Ulm, Germany.

Electronic supplementary material The online version of this article (doi:10.1007/s00122-017-2934-0) contains supplementary material, which is available to authorized users.

✉ Stefan Scholten
stefan.scholten@uni-hamburg.de

✉ Albrecht E. Melchinger
melchinger@uni-hohenheim.de

¹ Institute of Plant Breeding, Seed Science and Population Genetics, University of Hohenheim, 70599 Stuttgart, Germany

² Institute of Animal Breeding and Husbandry, Christian-Albrechts-University Kiel, 24098 Kiel, Germany

³ Inguran LLC dba STGenetics, 22575 SH6 South, Navasota, TX 77868, USA

⁴ Biocenter Klein Flottbek, Developmental Biology, University of Hamburg, 22609 Hamburg, Germany

⁵ Max-Planck Institute of Molecular Plant Physiology, 14476 Potsdam, Germany

⁶ Plant Breeding, Technische Universität München, 85354 Freising, Germany

Introduction

Hybrid breeding, which entails crossing of lines from two genetically distant germplasm collections—called heterotic groups (Melchinger and Gumber 1998)—has emerged as a prime strategy to meet demands for a sustainable intensification of agricultural production (Duvick 2005). However, unlocking the full potential of hybrid breeding requires accurate prediction methods to efficiently identify the superior candidates out of the millions of possible hybrids that could potentially be produced in each cycle of an ordinary-sized breeding program. With the advent of the doubled haploid (DH) technology (Wedzony et al. 2009) this prediction problem has become even more challenging because, based on breeder's experience, the vast majority ($\approx 90\%$) of competing lines in each heterotic group are “new” lines without any phenotypic records on hybrid progeny from previous breeding cycles. Consequently, among all hybrid combinations possible between lines from two heterotic groups, about 81% are T0 hybrids, 18% are T1 hybrids and 1% are T2 hybrids having zero, one or two parents, respectively, that have been previously tested in other hybrid combinations. Preselection of a few hundred of the most favorable hybrids with high success rate could significantly reduce the labor-intensive and time-consuming field-testing (Kadam et al. 2016; Xu et al. 2016). This could greatly impact the efficiency of hybrid breeding and boost the annual selection gain (Longin et al. 2015).

Whereas yield and other heterotic traits of hybrids are generally poorly predicted by the performance of their parent lines (Melchinger and Gumber 1998), WGP has emerged as a major tool for tackling this challenge (Massman et al. 2013; Technow et al. 2014). Nevertheless, there is evidence that, even with complete sequence information, genomic prediction may not capture complex interactions between genes and downstream regulation, which act through the entire cascade from genotype to phenotype (Dalchau et al. 2011; Zhu et al. 2012; Rudd et al. 2015; Ritchie et al. 2015a). Most studies have evaluated predictive ability by looking at only one kind of endophenotype (intermediaries between genotype and phenotype) (Gottesman and Gould 2003; Mackay et al. 2009) such as the transcriptome (Swanson-Wagner et al. 2006; Zenke-Philippi et al. 2016; Xu et al. 2016) or the metabolome (Riedelsheimer et al. 2012; Xu et al. 2016; Dan et al. 2016). The integration of different endophenotypic and genomic data is expected to reflect more closely the variability across genotypes than genomic data alone (Mackay et al. 2009; Patti et al. 2012; Civelek and Lusis 2014). Two recent studies that integrated multiple biological strata in predicting breast cancer risk (Vazquez et al. 2016) and performance of maize inbred lines (Guo et al. 2016), respectively, demonstrated the

benefit of this strategy. However, unlike forecasting clinical or agronomic traits from endophenotypes of the same genotype, hybrid breeding requires the prediction of the genotypic values (GV) of hybrid progeny based on parental information. To achieve this objective, we used the BLUP approach—originally developed in animal breeding (Henderson 1984)—for the more complex setting of hybrids between parents from two heterotic groups (Bernardo 1996; Massman et al. 2013). Here, we measured endophenotypes of parent lines to forecast the GV of T0, T1 and T2 hybrid progeny by using prediction equations trained with “omics” information from other parent lines and phenotypic information on their hybrid offspring.

Materials and methods

Genetic material and phenotyping

The entire genetic material consisted of a set of 1536 hybrids, denoted as H_{Tot} , produced in 16 factorial mating designs between 142 Dent and 103 Flint lines from the maize breeding program at the University of Hohenheim, on which agronomic data for silage maize production, as well as pedigree and genomic data, were available. A subset of this material, albeit from different trials, has been used for genomic prediction of traits related to grain maize production (Technow et al. 2014). For hybrid prediction, we used a core set $H \subset H_{\text{Tot}}$ of 617 hybrids, produced in six factorials with hybrid sets $H_{\text{FAC}(i)} (i = 1, 2, \dots, 6; H = \bigcup_{i=1}^6 H_{\text{FAC}(i)})$ from crosses between 57 Dent and 41 Flint inbred lines, denoted as $D = \{1, 2, \dots, 57\}$ and $F = \{1, 2, \dots, 41\}$ (File S1). All hybrids were evaluated in field experiments at three or more agro-ecologically diverse locations across Germany. In the trials of each factorial, which included at least five common check genotypes, the entries were randomized in α lattice designs and planted in 2-row plots. Dry matter yield (DMY, t/ha) and dry matter content (DMC, %) of whole-plant aboveground biomass were determined by established procedures (Riedelsheimer et al. 2012). For quality traits, contents of fiber (ADF, %), fat (FAT, %), protein (PRO, %), starch (STA, %) and sugars (SUG, %) in dry matter were measured in the harvested plant material using calibrated near-infrared spectroscopy [(NIRS; Grieder et al. (2011), File S1)].

Pedigree-based relationship coefficients

Coancestry coefficients were calculated using SAS (version 9.4, SAS Institute) for all possible pairs of lines in each heterotic group according to established rules (Falconer

and Mackay 1996) under the following assumptions (Cox et al. 1986): (1) all lines in a pedigree are genetically homogeneous and homozygous, (2) pairs of genotypes with no known common parentage are unrelated and (3) a line derived from a cross or backcross obtained a proportional fraction of the genome from each parent, as expected under Mendelian inheritance in the absence of selection.

Genotyping

Genotyping of all inbred lines was performed with the Illumina SNP chip MaizeSNP50 (Ganal et al. 2011). After performing a commonly used quality check (Technow et al. 2014) and imputation of missing data (Browning and Browning 2009), a total of 21,565 polymorphic SNPs was available and used for all further analyses.

Metabolite profiling

Seedlings of all parental inbred lines were grown under controlled conditions inside climate chambers to quantify the metabolite profiles of their roots 3.5 days after sowing, as detailed by de Abreu E Lima et al. (2017). The experiment was laid out as a randomized incomplete block design with replicated germination boxes. For leaf metabolic profiles (known metabolites and unannotated chromatographic peaks), a field experiment was carried out in an α lattice design with two replications at one location in southern Germany in the spring of 2012. Excision of leaves at the third leaf stage was performed according to an established protocol (Riedelsheimer et al. 2012) 28 days after sowing, in the afternoon of a cloudy day and finalized within 45 min for the entire experiment. For both profiling procedures, all material was transferred directly into containers with dry ice and then into liquid nitrogen to quench metabolic activity.

Transcriptome profiling

For transcriptome profiling, five seeds per parent line were taken from the same seed lot as used for metabolite profiling and laid out inside a climate chamber in a randomized complete block design with five replications. Seedlings were sampled 7 days after sowing, snap-frozen in liquid nitrogen and stored at -80°C until use. Prior to mRNA extraction, roots from all replicates of a genotype were pooled and homogenized. A custom 2K-microarray (GPL22267) was assembled from a subset of the 47-K maize oligonucleotide array (GPL6438). Two-color hybridizations were carried out separately for each of the six factorials using interwoven loop designs (Kerr and Churchill 2001). The average number of shared genotypes between factorials was 4.5 and ranged from 2 to 10.

Statistical analysis of agronomic traits

Agronomic data were analyzed in two stages, following Technow et al. (2014) by accounting for year, location, field replication, block and genotype effects as well as their interactions (File S1). In the first stage, and separately for each environment, best linear unbiased estimates (BLUEs) of the α designs were computed for every hybrid using REML-based linear mixed-model analyses. In the second stage, BLUEs were computed for all hybrids in H_{Tot} . The BLUEs of hybrids in the core set H served as response variables in our hybrid prediction models and cross-validation routines. For all predictions we used computationally efficient best linear unbiased predictor (BLUP) models, which have the same properties as those of a selection index because we previously accounted for fixed effects (Mrode (2014), pp. 34, 311, 312). For general and specific combining abilities (GCA and SCA) of parent lines, we used ASReml (Butler et al. 2009) to compute best linear unbiased predictors (BLUPs), variance components ($\sigma_{\text{GCA}^D}^2$, $\sigma_{\text{GCA}^F}^2$, σ_{SCA}^2) and entry-mean heritabilities (H^2) of all hybrids in H_{Tot} , treating all effects in the model as random. The covariance matrices of the GCA and SCA effects were defined by multiplying the variance components with their respective genomic relationship matrices (File S1).

Statistical analysis of endophenotypes

Raw data were normalized using established procedures for metabolites (van den Berg et al. 2006) and transcripts (Smyth and Speed 2003; Ritchie et al. 2007). From these data, we obtained BLUEs for metabolite levels and transcript abundance of each line using REML-based mixed-model analyses. The statistical models for the analysis of metabolite profiles accounted for various experimental effects as detailed by de Abreu E Lima et al. (2017). After applying quality checks and computing BLUEs, 92 leaf metabolic analytes and 283 root metabolic analytes remained for further analyses. BLUEs for transcriptomic data were computed using the R-package *limma* (Ritchie et al. 2015b) in reference to established protocols (Smyth and Speed 2003; Ritchie et al. 2007; Frisch et al. 2010). The design matrix for the linear model was based on the dye-labeling of a reference genotype. To account for possible differences between the microarrays, replicates of some genotypes across at least two factorials were included and modeled through a fixed effect term. All gene expression values were subsequently computed, based on the log-ratio relative to this common genotype (Smyth 2004). In total, 1323 gene expression profiles were available. Repeatabilities (w^2) were estimated for each endophenotype at the inbred line level using the same models as for the computation of BLUEs, but treating the genotype effect as random.

This analysis was performed jointly for the Dent and Flint lines allowing for different means and heterogeneous genotypic variances of the heterotic groups, but assuming a common error variance. Variance components were estimated by Gibbs sampling using the *R* package *MCMCglmm* (Hadfield 2010).

Prediction models and model evaluation

Predictions of hybrid performance were compared on the basis of the core set of hybrids H and the corresponding sets of parent lines D and F on which data for all five predictors (P, pedigree; G, genomic; T, transcriptomic; L, leaf metabolic; R, root metabolic data) were available with the exception of data on a few lines missing at random for R due to fungal contamination. The matrices \mathbf{W}_D and \mathbf{W}_F are matrices of standardized feature measurements for the various predictors (G, T, L, R). The matrix \mathbf{W} has dimension 'number of parent lines in the corresponding heterotic group' ($n_D = 142$, $n_F = 103$) times 'number of features' ($w_G = 21,565$, $w_T = 1,323$, $w_L = 92$, $w_R = 283$). The columns in \mathbf{W}_D and \mathbf{W}_F are centered and standardized to unit variance, respectively.

The kernels pertaining to each predictor and lines from each heterotic group—corresponding to genomic relationship matrices in the case of SNPs—can then be defined as

$$\mathbf{G}_D = \frac{1}{W} \mathbf{W}_D \mathbf{W}_D^T, \mathbf{G}_F = \frac{1}{W} \mathbf{W}_F \mathbf{W}_F^T, \quad (1)$$

where W denotes the number of features (VanRaden 2008). In the case of pedigree data (P), coancestry coefficients were used directly for \mathbf{G}_D and \mathbf{G}_F , respectively.

The universal model for GCA and SCA effects was as follows:

$$\mathbf{y} = \mu + \sum_{c=1}^C \mathbf{Z}_D \mathbf{g}_{Dc} + \sum_{c=1}^C \mathbf{Z}_F \mathbf{g}_{Fc} + \sum_{c=1}^C \mathbf{Z}_S \mathbf{s}_c + \epsilon, \quad (2)$$

where \mathbf{y} is the vector of observed hybrid performance (BLUEs), μ is the fixed model intercept, \mathbf{Z}_D is the corresponding design matrix associating the random GCA effects of the lines in D (\mathbf{g}_{Dc}) with \mathbf{y} , \mathbf{Z}_F is the corresponding design matrix associating the random GCA effects of the lines in F (\mathbf{g}_{Fc}) with \mathbf{y} and \mathbf{Z}_S is a design matrix associating the SCA effects (\mathbf{s}_c), pertaining to hybrid combinations for the c -th predictor data type with the corresponding hybrid measurements in \mathbf{y} . Thus, the model in Eq. 2 can accommodate just one ($C = 1$) or multiple ($C > 1$) predictors simultaneously. The random effects (\mathbf{g}_{Dc} and \mathbf{g}_{Fc}) have expectation zero and covariance matrices equal to $\mathbf{G}_{Dc} \sigma_{\text{GCA}_{Dc}}^2$ and $\mathbf{G}_{Fc} \sigma_{\text{GCA}_{Fc}}^2$ for the GCA effects of the Dent and Flint lines, respectively, $\mathbf{S}_c \sigma_{\text{SCA}_c}^2$ for the SCA effects and $\mathbf{I} \sigma_\epsilon^2$ for the residual error. For

each combination between crosses of lines $i \times k$ and $j \times l$, the corresponding elements in \mathbf{S}_c were obtained as the product of the respective elements f_{ij} in \mathbf{G}_{Dc} and f_{kl} in \mathbf{G}_{Fc} , respectively (Schnell 1965; Henderson 1985; Bernardo 1996; Massman et al. 2013; Technow et al. 2014; Jiang and Reif 2015) (File S1). Note that, in the majority of cases, only GCA effects were considered. In the absence of epistasis, this model is equivalent to a feature model accounting for dually defined additive effects in each heterotic group and dominance effects between them. Extensions of the single-predictor models were made by adding GCA and SCA effects for any additional predictor assuming stochastic independence of effects. In order to obtain unbiased estimates of the predictive ability and to compare different models and predictor combinations, following Technow et al. (2014), we devised a cross-validation (CV) scheme, stratified by the parent lines and using 1000 runs (CV1000, File S1). All prediction models were implemented using the *R* package *BGLR* (Pérez and de Los Campos 2014).

Comparison of predictive abilities

Predictive abilities were obtained by calculating Pearson correlations between predicted (\hat{y}) and observed phenotypes (y), separately for three test set partitions (T0, T1 and T2 hybrids). For each CV run, the training and validation sets were stored to ensure the validity of comparisons between any predictor and combinations thereof. For any two predictors, say A and B , we then have orthogonal vectors with predictive abilities r_A and r_B of length 'number of cross validation runs'.

Evaluation of a pre-selection bias in transcriptomic data

A custom 2K-microarray (GPL22267) was assembled from a subset of the 47-K maize oligonucleotide array (GPL6438), based on association of genes with hybrid performance or mid-parent heterosis for grain yield and grain dry matter content of maize. These two traits were evaluated in separate grain-yield trials with hybrids from factorial $H_{\text{FAC}(1)}$ [Frisch et al. (2010); Thiemann et al. (2010), File S1]. To ensure that no pre-selection bias was introduced in hybrid prediction using these transcriptomic data, we compared predictive abilities among the various predictors when excluding $H_{\text{FAC}(1)}$ from the entire set H .

Association mapping

For each of the seven agronomic traits, we performed a genome-wide association study (GWAS) with GCA

effects of all 142 Dent and 103 Flint parent lines as response variables using the EMMAX-method (Kang et al. 2010) as implemented in *cpgen* (Heuer 2015). To avoid using the marker data twice, GCA effects were calculated using only pedigree information. Furthermore, an eQTL analysis was carried out to examine statistically significant associations between genomic and transcriptomic data for the parent lines (*D* and *F*) of the core set *H* plus five additional lines. This was accomplished in the same way as in the GWAS for agronomic traits, but here the BLUPs of the transcriptomic data of each mRNA were used as the response variables. Associations in each GWAS were declared statistically significant at $\alpha = 0.05$ after Bonferroni correction.

Probability of success

Following Robson et al. (1967), we calculated the probability of success ($P[r, \beta]$) that a hybrid, selected at random from the upper β percent fraction of the distribution of predicted values for predictor A, has a phenotypic value contained in the upper β percent of the distribution of observed values. Denoting the predictive ability of a given predictor by r , this conditional probability was calculated assuming a bivariate normal distribution

$$\begin{bmatrix} \hat{y} \\ y \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & r \\ r & 1 \end{bmatrix} \right). \quad (3)$$

The required integrals were solved within the *R* statistical environment using the *mvtnorm* package (Genz et al. 2017).

Principal component analysis

Principal component analyses (PCA) were carried out to examine whether different predictors can distinguish between Dent and Flint parent lines and to explore whether subpopulations exist within either heterotic group. Prior to each PCA, all variables were scaled and centered. Clusters represent two component mixtures of bivariate *t*-distributions, which were estimated using Maximum Likelihood. Ellipses were drawn based on the 0.95 quantiles of the respective bivariate *t*-distributions. Unless stated otherwise, all statistical analyses were carried out inside the *R* environment for statistical computing (R Core Team 2016).

Data availability

The data and the code used to analyze the data are available upon request.

Results

Agronomic data

Mean values of the 1536 hybrids for the seven evaluated agronomic traits, relevant for animal feed and biogas production, were of the same magnitude as reported by Riedelsheimer et al. (2012) and Grieder et al. (2011). For all traits, $\sigma_{GCA^D}^2$ and $\sigma_{GCA^F}^2$ describing the main effects of the parents from each heterotic group, together explained more than 93% of the genotypic variance among hybrids (Table 1). Heritabilities were moderate to high for all agronomic traits, indicating a high precision of field experiments and data collection (Table 1).

Predictor data

Repeatabilities (w^2) for endophenotypes varied considerably in both groups of parents (Fig. S1a) with average values ranging from 0.31 to 0.41, except for transcriptomic data in Flint material where the average repeatability was only 0.18. Nevertheless, in the latter case, 291 out of 1323 transcripts still exceeded a threshold of $w^2 = 0.4$.

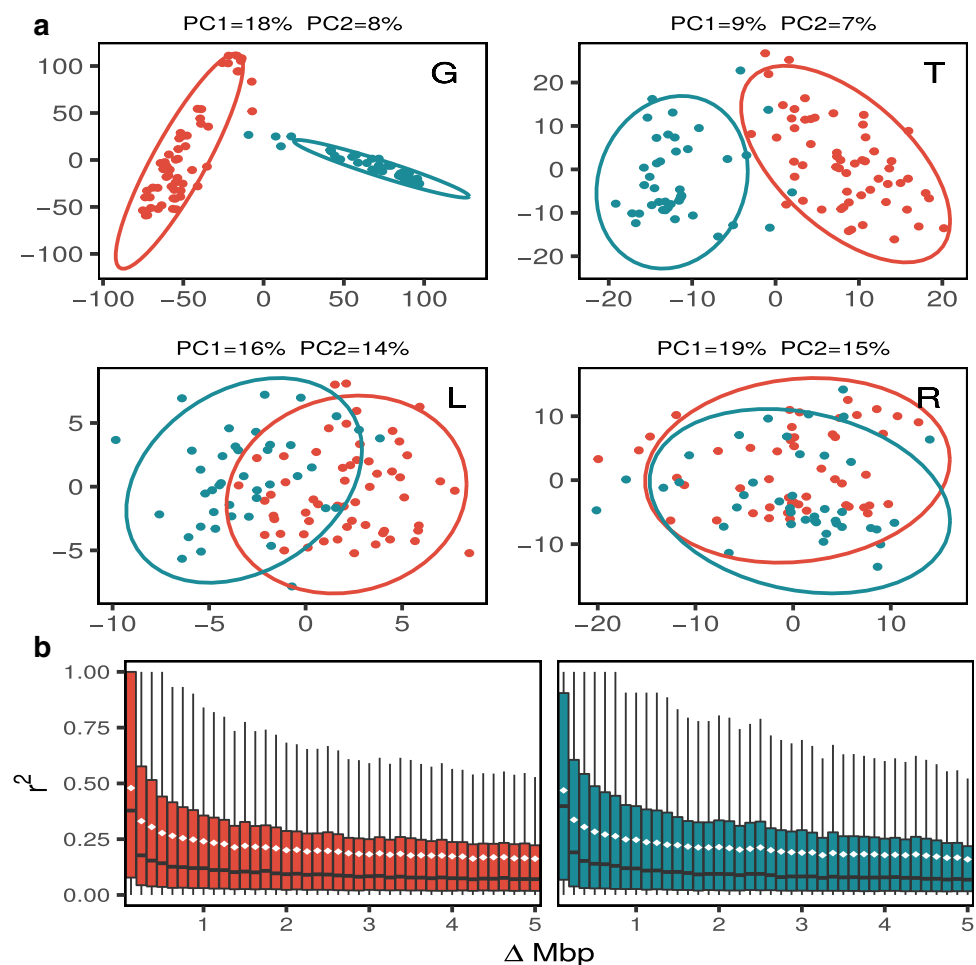
Dent and Flint lines were clearly separated in principal component analyses of genomic and transcriptomic data (Fig. 1a) without signs of subpopulations within either group. However, they overlapped for leaf metabolic and, to an even greater extent, for root metabolic data. Off-diagonal elements of the kernels \mathbf{G}_D and \mathbf{G}_F , respectively, showed moderate correlations between genomic and transcriptomic data ($\rho_D \approx 0.56$, $\rho_F \approx 0.44$, Fig. S2). Correlations between the off-diagonal elements of the \mathbf{G} -matrices were highest for the comparison between genomic and pedigree data ($\rho_D \approx 0.72$, $\rho_F \approx 0.63$). Intriguingly, the associations between the \mathbf{G} -matrices for the root and leaf metabolic data were very low ($\rho_D \approx 0.12$, $\rho_F \approx 0.06$).

Table 1 Summary of agronomic traits

Trait	μ	$\sigma_{GCA^D}^2$	$\sigma_{GCA^F}^2$	σ_{SCA}^2	H^2
DMY (t/ha)	19.00	1.51 ± 0.22	1.00 ± 0.18	0.17 ± 0.03	0.82
DMC (%)	34.13	4.17 ± 0.59	5.03 ± 0.79	0.49 ± 0.07	0.91
ADF (%)	20.93	0.27 ± 0.06	0.40 ± 0.09	0.02 ± 0.02	0.43
FAT (‰)	30.02	1.01 ± 0.19	2.10 ± 0.37	0.17 ± 0.05	0.73
PRO (‰)	69.65	3.11 ± 0.53	2.77 ± 0.51	0.29 ± 0.10	0.70
STA (%)	35.56	2.95 ± 0.51	3.33 ± 0.63	0.24 ± 0.10	0.69
SUG (‰)	38.12	31.33 ± 5.15	31.90 ± 5.73	4.12 ± 1.02	0.77

Traits are characterized by overall mean (μ), variance components of GCA effects for Dent ($\sigma_{GCA^D}^2$) and Flint lines ($\sigma_{GCA^F}^2$) and SCA effects (σ_{SCA}^2) (followed by SEM) as well as entry mean heritabilities (H^2)

Fig. 1 Properties of predictor data for Dent (*red*) and Flint (*teal*) parent lines. **(a)** Principal component analysis (PCA) of both groups for genomic (G), transcriptomic (T), leaf metabolic (L) and root metabolic (R) data. The variance explained by PC 1 (x axis) and PC 2 (y axis) are shown in the caption of each facet. **(b)** Linkage disequilibrium decay as a function of the distance between two loci using 40 bins of 0.125 Mbp width, each. The median r^2 is depicted as a horizontal bar, whereas the mean r^2 is depicted as a white diamond



We observed high median pairwise linkage disequilibrium (LD) between SNP markers ($r^2 \approx 0.39$ in Dent and $r^2 \approx 0.37$ in Flint material) at a distance of $\Delta \text{Mbp} \leq 0.125$ (Fig. 1b). After an initial drop in r^2 for $\Delta > 0.125$, substantial long-range LD remained. Large differences in allele frequencies in the two heterotic groups were present for 57% of SNPs (Fig. 2a, b)—particularly in the telomeric regions of the genome. An eQTL analysis performed with the parent lines suggests that transcript abundance integrates variegated genetic information given the fact that (1) on the same chromosome, significant associations not only occurred between adjacent but also between distant pairs of expressed genes and SNPs and (2) 50% of the significant associations ($\alpha = 0.05$, Bonferroni-corrected) occurred between expressed genes and SNPs on different chromosomes (Fig. 2).

Predictive abilities

Assuming a polygenic architecture for all traits, as suggested by results from a GWAS (Fig. S3), we chose the best linear

unbiased predictor (BLUP) method as a baseline for prediction of T0, T1 and T2 hybrids. Given that we corrected for fixed effects in advance, this method corresponds to a selection index. A cross-validation scheme with 1000 runs (CV1000), stratified by the parent lines, was devised (File S1, Fig. S4). Our main emphasis was on predicting T0 hybrids given the fact that they constitute the majority of possible hybrids in practical breeding programs (Kadam et al. 2016).

For predictive abilities (r) of T0 hybrids, transcriptomic data alone were the best predictor for the most complex and highly heterotic trait, DMY, as well as for PRO (Fig. 3a). With transcriptomic data, the predictive ability r for DMY was 14.9% higher than for genomic data, resulting in an 85% increase in the probability of successfully selecting the best hybrid candidates $P[r, \beta]$ for $\beta = 0.01\%$ (Fig. 3b). This selection intensity corresponds to picking the top 100 out of 10^6 predicted hybrids for production and intensive testing in field trials.

Compared to other individual predictors, r obtained with genomic data alone were higher for FAT and SUG. Root

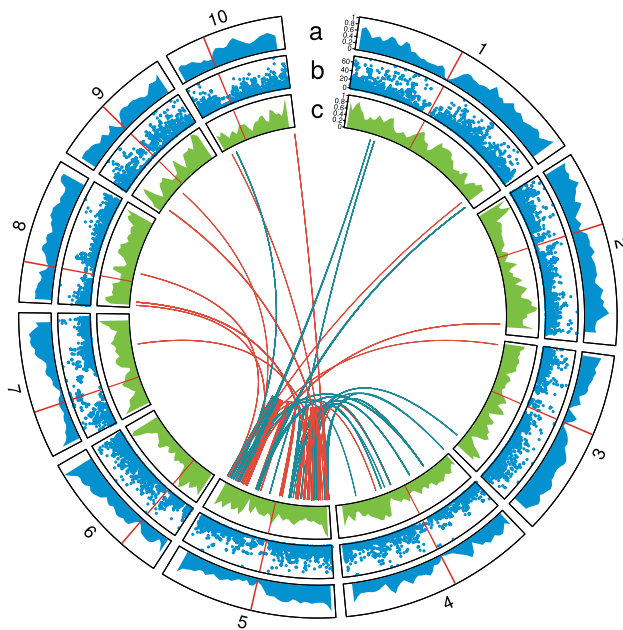


Fig. 2 Distribution and relationship of genomic and transcriptomic data for the ten maize chromosomes. Centromeres for each chromosome are depicted as vertical red lines. (a, c) Density of SNPs and mRNAs, respectively, across the ten maize chromosomes. (b) Statistical significance ($-\log_{10} p$ value) for differences in SNP allele frequencies between Dent and Flint lines. Center links between any statistically significant ($\alpha = 0.05$ after Bonferroni correction) association between SNPs and mRNAs. Associations are displayed as links for SNPs on chromosome 5, for which the distribution of associations is representative for the entire genome, using red color for Dent parent lines and teal color for Flint parent lines

metabolites displayed moderate to high predictive abilities for DMY and FAT, but did not perform well otherwise. Leaf metabolites performed relatively poorly for all traits. Regardless of the trait, combinations of genomic and transcriptomic information displayed robust and consistently high predictive abilities. Except for PRO, incorporating additional endophenotypes as predictors into our models did not yield notable improvements but remained at the same level compared to combining genomic and transcriptomic data. Incorporating SCA effects into our models did not further improve predictive abilities (Fig. S5). Results for the combination of other predictors with metabolic data are not presented because no improvement of predictive abilities over the combination of genomic data with transcriptomic data and pedigree data could be achieved. Finally, we assessed the influence of the number of SNPs and mRNAs on predictive abilities. For genomic data, a subset of 5000 SNPs already yielded the same predictive ability as when using the entire available set. For transcriptomic data, the predictive ability improved only marginally with subsets larger than 50% of the available transcripts (Fig. S6).

Discussion

A paradigm shift in hybrid breeding

Hybrid breeding programs are generally based on genetically divergent heterotic groups. Their use enables a better exploitation of heterosis when conducting crosses between them (Melchinger and Gumber 1998) and is expected to reduce the ratio of specific to general combining ability variance ($\sigma_{SCA}^2 : \sigma_{GCA}^2$) in the crosses, thereby allowing for the selection of hybrids largely on the basis of GCA of their parent lines (Reif et al. 2007). However, obtaining accurate estimates of GCA requires the evaluation of new lines in combinations with testers from the opposite heterotic group in multi-environment field trials. The promise of hybrid prediction is to accelerate breeding programs by skipping a large share of these tests in favor of selecting the most promising hybrids before they are even produced (Technow et al. 2014). This approach involves the prediction of an impressive number of putative hybrid candidates (n^2) using predictor data collected on only $2n$ parent lines. Crucial for hybrid prediction are predictors, which not only reflect the relationship between parental inbred lines but also the interaction of the two parental genomes in their hybrid progeny.

Heterotic groups Because of genetic drift and selection for hybrid performance, allele frequencies are expected to diverge in the two heterotic groups, thereby enlarging their genetic distance (Falconer and Mackay 1996; Reif et al. 2007; Lari pe et al. 2017). Consistent with this hypothesis and two pilot studies with U.S. maize lines (Gerke et al. 2015; Hall et al. 2016), Dent and Flint lines in our study were clearly separated in principal component analyses of genomic and transcriptomic data. With large differences in allele frequencies p^I and p^{II} in the two heterotic groups, as observed for 57% of SNPs, dominance variance σ_D^2 becomes very small because it is a function of the product $p^I(1 - p^I)p^{II}(1 - p^{II})$ [Stuber and Cockerham (1966), File S1].

Dominance variance (σ_D^2) is the main component contributing to the variance of the specific combining ability effects (σ_{SCA}^2), describing all types of interactions among the parental genomes in hybrid combinations. It was, therefore, not surprising that the variances of the general combining ability (GCA) effects ($\sigma_{GCA^D}^2$ and $\sigma_{GCA^F}^2$), describing the main effects of the parents from each heterotic group, together explained more than 93% of the genotypic variance among hybrids for agronomic traits, which is consistent with earlier studies on silage maize of the Dent \times Flint heterotic pattern (Geiger et al. 1986; Argillier et al. 2000). While the magnitude of SCA effects was trait-specific, it was low for all observed traits, which is in agreement

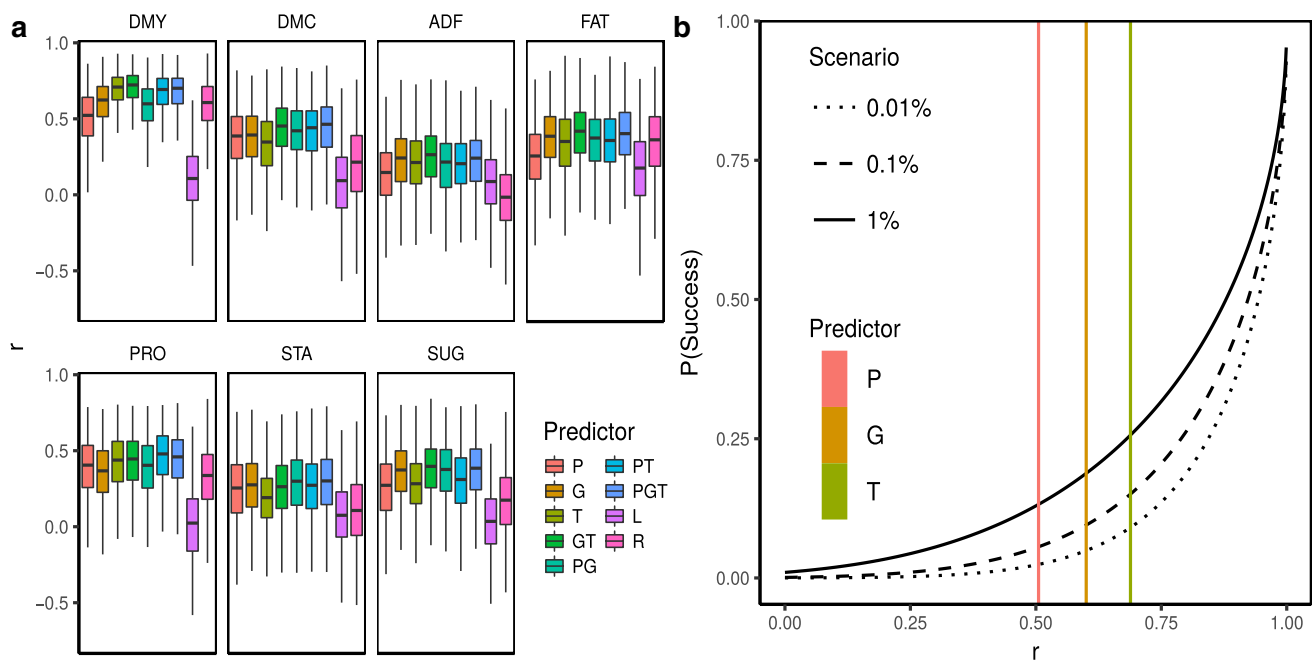


Fig. 3 Predictive abilities (r) from BLUP models using a CV scheme with sampling of $|H_{\text{TRN}}| = 200$ hybrids, $|D_{\text{TRN}}| = 40$ Dent and $|F_{\text{TRN}}| = 33$ Flint parent lines for various predictors and combinations thereof (P, pedigree; G, genome; T, transcriptome; L, leaf metabolome; R, root metabolome). (a) Comparison of r values from 1000 CV runs for T0 hybrids and seven agronomic traits. (b) Success rate of selecting superior hybrids ($P[r, \beta]$). $P[r, \beta]$ is a function of the

predictive ability $r = r(y, \hat{y})$ and refers to the conditional probability of a hybrid, selected at random from the upper $\beta\%$ fraction of the distribution of predicted values (\hat{y}), having a phenotypic value contained in the upper $\beta\%$ of the distribution of phenotypic values y . Observed predictive abilities (r) for T0 hybrids and the trait DMY are displayed as vertical, colored lines for three predictors

with previously reported values for yield and quality traits in silage maize (Grieder et al. 2012). The importance of GCA in our material was further corroborated by merely marginal differences in predictive abilities between models using only GCA effects and those that additionally incorporated SCA effects (Fig. S5). Nevertheless, in crops such as wheat, with yet no clearly defined heterotic groups (Zhao et al. 2015) and greater importance of SCA, inclusion of SCA effects in the model should improve predictive abilities.

Properties of well-established predictors While pedigree data reflect the expected relationship between genotypes, they do not necessarily depict their realized relationship. Genomic data and downstream endophenotypes offer to improve upon this pedigree-based approximation by more closely mirroring the transmission of genes between genotypes and their interactions. Genomic data have the advantage of reliably capturing Mendelian sampling, thereby improving pedigree-based prediction for many traits. However, genomic data alone may not be the final answer for the prediction of complex traits for two major reasons: First, the number of samples in most studies is considerably smaller than the number of genetic markers or even nucleotides of a genome. This implies that just modeling

additive effects already necessitates shrinkage of effects. More importantly, however, interactions between loci throughout the genome can be frequent (Brem et al. 2005; Brown et al. 2014), but attempts to incorporate this epistasis for the prediction of heterotic traits using genomic data have been disappointing when the prediction and training set did not share the same or closely related parents (Jiang and Reif 2015). This was true even when using recently developed, efficient models (Jarquín et al. 2014; Martini et al. 2016) and suggests that genomic data capture only statistical epistasis, referring to genetic variation at the population level (Sackton and Hartl 2016), which is generally of negligible magnitude (Hill et al. 2008; Mackay 2014; Guo et al. 2016; Vazquez et al. 2016).

Complementation of predictors

Flow of biological information It is well-known that genetic effects on the phenotype are mediated through multiple layers of endophenotypes (Civelek and Lusis 2014; Ritchie et al. 2015a) with information mainly flowing from the genome toward the phenotype via the transcriptome, the proteome and the metabolome with metabolite fluxes ultimately governing energy production and growth (Fiévet

et al. 2010). For most traits in our material, metabolite- and pedigree-based predictive abilities were lower than those obtained with either transcriptomic or genomic information. However, consistently high predictive abilities across multiple traits could be realized when combining multiple predictors, as has been reported previously in humans (Vazquez et al. 2016) and maize inbred lines (Guo et al. 2016). This suggests complementary properties of the different predictors resulting in better proxies for the complex interplay in gene networks than genomic information alone. Such an advantage is particularly important for hybrid prediction when parents of prediction set hybrids are not closely related to parents of training set hybrids (File S1) as was shown by the relative excellence of transcriptomic data and the use of multiple predictors for the prediction of traits in T0 hybrids compared to T1 and T2 hybrids.

Tapping new sources of information Whereas pedigree and genomic information are static, subsequent endophenotypes are characterized by pervasive interactions among and between each other (Dalchau et al. 2011; Zhu et al. 2012) and are, to varying degrees, influenced by biotic (Rudd et al. 2015; Tzin et al. 2015) and abiotic perturbations (Caldana et al. 2011; Witt et al. 2012). So while endophenotypes report do not exclusively on physiological epistasis but also on non-heritable effects, they seem to capture important information not represented by the genome given their intermediate position in the genotype-phenotype cascade. We get support for this hypothesis from (1) merely low to moderate correlations between off-diagonal elements of the kernels of different predictors in our study, (2) mounting evidence for further improvements of predictive abilities when complementing genomic prediction with other endophenotypes despite sufficient marker densities [Fig. S6, Guo et al. (2016)] and (3) the integration of SNP information from close and distant eQTL in the transcripts analyzed in our study. However, we concede that the number of parental genotypes in our mRNA assays was too small to warrant a reasonable statistical power for detecting epistasis in the expression of transcripts. In breeding programs, predictive abilities are largely driven by relationships—including Mendelian sampling—among genotypes compared to LD between SNP markers and causal QTL (Schopp et al. 2017). Increasing marker densities therefore have limited utility for improving genomic predictions as observed in our material, where SNP-based predictive abilities reached a plateau after using 5000 equally spaced SNPs (Fig. S6). While two other studies also attempted to model interactions between different predictors, we refrained from this approach given that their reported predictive abilities based on interactions were not different from those in additive models despite using much larger sample sizes (Vazquez et al. 2016; Guo et al. 2016).

Transcriptomic data

Utility of transcriptomic data for trait predictions Of particular note was the excellent performance of transcriptomic data in predicting dry matter yield and protein. Evidence that parental gene expression patterns might be predictive of hybrid performance is given by (1) prevailing additive expression patterns in maize hybrids (Springer and Stupar 2007a; Stupar et al. 2008), (2) a positive correlation of the proportion of additive gene expression with the yield of hybrids (Guo et al. 2006) and (3) co-localization of additively expressed genes with heterotic QTL (Thiemann et al. 2014). According to metabolic flux theory, gene expression in hybrids at the mid-parent level can generate hybrid vigor by counterbalancing opposing detrimental expression levels in their parent lines on a genome-wide scale (Kacser and Burns 1981; Springer and Stupar 2007b). The same concept is expected to apply to other quantitative endophenotypes (Lisec et al. 2011).

Pre-selection bias As pointed out earlier, our transcripts were pre-selected based on associations with grain dry matter yield and grain dry matter content in hybrids, using a subset of the data included in our study ($H_{FAC(1)}$). Hence, genotypes used for the pre-selection (i.e. $H_{FAC(1)}$) could be regarded as a training set. By combining this “training set” and genotypes from the remaining five factorials, we might have introduced a bias by using predictors that have already seen the response variable in the $H_{FAC(1)}$ genotypes. To rule out the existence of such a bias, we have compared the predictive abilities of different predictors for the complement of $H_{FAC(1)}$. Two findings indicate that no bias in the comparison of predictive abilities was introduced: (1) Relative differences in predictive abilities between transcriptomic and pedigree or genomic data did not change when excluding genotypes from $H_{FAC(1)}$ from the data (Fig. S7) and (2) transcriptomic data performed rather poorly in predicting dry matter content although this trait was also among the criteria for the pre-selection procedure. Finally, an independent study using RNA-Seq data for the prediction of traits in maize inbred lines also reported exceptionally good performance of transcriptomic data in the prediction of multiple yield-related traits (Guo et al. 2016).

Relative excellence of predictors for different traits

Tissue and sampling time Despite the great prospects of using endophenotypes for trait predictions, some aspects require careful consideration when using this approach. A particular challenge in endophenotype-based prediction efforts is the choice of a suitable tissue and sampling time. Tissue-related effects regarding gene expression were found in studies on humans (Yang et al. 2015; Mele et al. 2015; Searle et al. 2016) and *A. thaliana* (Schmid et al. 2005)

and in maize hybrids with respect to metabolome composition and metabolite abundance (Witt et al. 2012). Moreover, the age of an organism can selectively influence the expression of genes as observed in studies on humans (Mele et al. 2015; Yang et al. 2015) and *C. elegans* (Vinuela et al. 2010; Francesconi and Lehner 2014). The low correlations between the off-diagonal elements of the kernels calculated from root and leaf metabolites might, therefore, be a reflection of highly dynamic processes differing between tissues and during different developmental stages. Whereas root metabolic data and transcriptomic data were obtained from seedlings germinated in standard controlled conditions, leaf metabolic data were derived from field-grown plants at a much later developmental stage, thereby increasing the possibility of environmentally induced modifications. One might hypothesize that the choice of sampling time and tissue could influence the chances of successful trait prediction if such age- or tissue-dependent transcripts and metabolites are associated with a phenotypic or clinical trait.

Feature selection Another explanation for trait-dependent excellence of any predictor might lie in the sampling of features. In this study, only a small subset of metabolites was sampled and even very recent technologies (Xu et al. 2016; Dan et al. 2016) capture only a fraction of the estimated set of metabolites (Ferne 2007). Moreover, the smaller differences in metabolite levels between both heterotic groups (Fig. 1) were most likely not conducive to capturing basic components underlying complex heterotic traits. It is also possible that transcriptomic data are associated with more biological processes than metabolite data and better capture the genetic effects relevant for the prediction of T0 hybrids.

Prospects for metabolites Previously observed moderate metabolite-based predictive abilities for T1 hybrids (Riedelsheimer et al. 2012) were confirmed in our study (Fig. S8), but for the majority of traits, root metabolites reached only medium and leaf metabolites even lower predictive abilities when predicting T0 hybrids. Despite the aforementioned shortcomings of metabolites, they have shown to be intriguing predictors due to their physiological proximity to the phenotype, which provides information that is impossible to infer from DNA or proteins (Ferne and Stitt 2012), as well as encouraging results from other studies (Guo et al. 2016; Dan et al. 2016). A recently introduced technology, allowing for live-measurements of small molecules in the blood of living and awake animals (Arroyo-Currás et al. 2017), might overcome the problem of poorly time-resolved snap-shots of some metabolites with extremely fast turnover rates (Arrivault et al. 2009) if modified to properly work in plants.

Predictor requirements Besides improving upon predictions based on pedigree relationships by capturing Mendelian sampling, the widespread use of genomic information

in trait prediction has been driven by the ease of its application. In order to compete with genomic data, other 'omics' data, therefore, require the use of standardized sampling conditions to obtain large repeatabilities and the possibility of season-independent sample extraction from seeds, seedlings or young roots to achieve high throughput.

Conclusions

The use of whole-genome information has considerably advanced trait prediction over traditional pedigree-based BLUP by incorporating previously unobservable Mendelian sampling. Combining variegated sources of information promises to capture complex interactions between genes and endophenotypes, leading to stable predictions across traits. Especially if an extremely small fraction of the candidates is selected from the millions of possible new hybrids from each breeding cycle, the success of forecasts is a strongly convex function of predictive ability (Fig. 3b). Therefore, considering endophenotypes could have a substantial effect on the success and economics of hybrid breeding. Given the anticipated technological improvements in RNA-Seq and metabolite profiling, as well as the forthcoming adoption of the DH-technology for many crops (Kelliher et al. 2017), a paradigm shift from exclusively genomic prediction models to more inclusive approaches seems imminent.

Author contribution statement WS and AEM developed the lines and hybrids, WS and AEM designed the field experiments, TAS analyzed the agronomic and pedigree data, LW, MS, AS, AEM and MW designed the metabolic experiments, AS conducted the metabolic experiments, SS designed the transcriptomic experiments, FS and AT conducted the transcriptomic experiments, MW analyzed the metabolic and transcriptomic data, MW, TAS, GT, CH and AEM devised the prediction models, MW, CH and TAS implemented the prediction models and developed software, HFU and AE contributed to the statistical analysis. MW, AEM, SS, GT, ZN and CCS wrote the manuscript.

Acknowledgements We thank the staff of the Agricultural Experimental Research station, University of Hohenheim, for excellent technical assistance in conducting the field experiments. We are indebted to the group of R. Fries from Technische Universität München for the SNP genotyping of the parent inbred lines, to X. Mi for his assistance in preparing auxiliary figures based on the Mathematica software, to C. Zenke for advice on the computation of transcriptomic BLUEs and to P. Schopp for advice on prediction models.

Compliance with ethical standards

Conflict of interest The authors declare that they have no conflict of interest.

References

- Argillier O, Méchin V, Barrière Y (2000) Inbred line evaluation and breeding for digestibility-related traits in forage maize. *Crop Sci* 40(6):1596–1600. doi:[10.2135/cropsci2000.4061596x](https://doi.org/10.2135/cropsci2000.4061596x)
- Arrivault S, Guenther M, Ivakov A, Feil R, Vosloh D, Van Dongen JT, Sulpice R, Stitt M (2009) Use of reverse-phase liquid chromatography, linked to tandem mass spectrometry, to profile the Calvin cycle and other metabolic intermediates in *Arabidopsis* rosettes at different carbon dioxide concentrations. *Plant J* 59(5):824–839. doi:[10.1111/j.1365-3113X.2009.03902.x](https://doi.org/10.1111/j.1365-3113X.2009.03902.x)
- Arroyo-Currás N, Somerson J, Vieira PA, Ploense KL, Kippin TE, Plaxco KW (2017) Real-time measurement of small molecules directly in awake, ambulatory animals. *Proc Natl Acad Sci USA* 114(4):645–650. doi:[10.1073/pnas.1613458114](https://doi.org/10.1073/pnas.1613458114)
- van den Berg RA, Hoefsloot HJ, Westerhuis JA, Smilde AK, van der Werf MJ (2006) Centering, scaling, and transformations: improving the biological information content of metabolomics data. *BMC Genom* 7:1–15. doi:[10.1186/1471-2164-7-142](https://doi.org/10.1186/1471-2164-7-142)
- Bernardo R (1996) Best linear unbiased prediction of maize single-cross performance. *Crop Sci* 36:50–56
- Brem RB, Storey JD, Whittle J, Kruglyak L (2005) Genetic interactions between polymorphisms that affect gene expression in yeast. *Nature* 436(7051):701–3. doi:[10.1038/nature03865](https://doi.org/10.1038/nature03865)
- Brown AA, Buil A, Vinuela A, Lappalainen T, Zheng HF, Richards JB, Small KS, Spector TD, Dermitzakis ET, Durbin R (2014) Genetic interactions affecting human gene expression identified by variance association mapping. *eLife* 3:1–16. doi:[10.7554/eLife.01381](https://doi.org/10.7554/eLife.01381)
- Browning BL, Browning SR (2009) A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am J Hum Genet* 84(2):210–223. doi:[10.1016/j.ajhg.2009.01.005](https://doi.org/10.1016/j.ajhg.2009.01.005)
- Butler DG, Cullis BR, Gilmour AR, Gogel BJ (2009) Mixed models for S language environments, ASReml-R reference manual. Training and development series, No QE02001. QLD Department of Primary Industries and Fisheries, Brisbane
- Caldana C, Degenkolbe T, Cuadros-Inostroza A, Klie S, Sulpice R, Leisse A, Steinhauser D, Fernie AR, Willmitzer L, Ma Hannah (2011) High-density kinetic analysis of the metabolomic and transcriptomic response of *Arabidopsis* to eight environmental conditions. *Plant J* 67(5):869–884. doi:[10.1111/j.1365-3113X.2011.04640.x](https://doi.org/10.1111/j.1365-3113X.2011.04640.x)
- Civelek M, Lusk AJ (2014) Systems genetics approaches to understand complex traits. *Nat Rev Genet* 15(1):34–48. doi:[10.1038/nrg3575](https://doi.org/10.1038/nrg3575)
- Cox TS, Murphy JP, Rodgers DM (1986) Changes in genetic diversity in the red winter wheat regions of the United States. *Proc Natl Acad Sci USA* 83(15):5583–5586. doi:[10.1073/pnas.83.15.5583](https://doi.org/10.1073/pnas.83.15.5583)
- Dalchau N, Baek SJ, Briggs HM, Robertson FC, Dodd AN, Gardner MJ, Stancombe MA, Haydon MJ, Stan GB, Gonçalves JM, Webb AAR (2011) The circadian oscillator gene *GIGANTEA* mediates a long-term response of the *Arabidopsis thaliana* circadian clock to sucrose. *Proc Natl Acad Sci USA* 108(12):5104–5109. doi:[10.1073/pnas.1015452108](https://doi.org/10.1073/pnas.1015452108). arXiv:1408.1149
- Dan Z, Hu J, Zhou W, Yao G, Zhu R, Zhu Y, Huang W (2016) Metabolic prediction of important agronomic traits in hybrid rice (*Oryza sativa* L.). *Nat Sci Rep* 6(October 2015):1–9. doi:[10.1038/srep21732](https://doi.org/10.1038/srep21732)
- de Abreu E, Lima F, Westhues M, Willmitzer L, Melchinger AE, Nikołoski Z (2017) Metabolic robustness in young roots underpins a predictive model of maize hybrid performance in the field. *Plant J* 90(2):319–329. doi:[10.1111/tpj.13495](https://doi.org/10.1111/tpj.13495)
- Duvick DN (2005) Genetic progress in yield of United States maize (*Zea mays* L.). *Maydica* 50:193–202
- Falconer DS, Mackay TFC (1996) Introduction to quantitative genetics, 4th edn. Pearson Prentice, Harlow, UK
- Fernie AR (2007) The future of metabolic phytochemistry: larger numbers of metabolites, higher resolution, greater understanding. *Phytochemistry* 68(22–24):2861–2880. doi:[10.1016/j.phytochem.2007.07.010](https://doi.org/10.1016/j.phytochem.2007.07.010)
- Fernie AR, Stitt M (2012) On the discordance of metabolomics with proteomics and transcriptomics: coping with increasing complexity in logic, chemistry, and network interactions scientific correspondence. *Plant Physiol* 158(3):1139–45. doi:[10.1104/pp.112.193235](https://doi.org/10.1104/pp.112.193235)
- Fiévet JB, Dillmann C, de Vienne D (2010) Systemic properties of metabolic networks lead to an epistasis-based model for heterosis. *Theor Appl Genet* 120(2):463–73. doi:[10.1007/s00122-009-1203-2](https://doi.org/10.1007/s00122-009-1203-2)
- Francesconi M, Lehner B (2014) The effects of genetic variation on gene expression dynamics during development. *Nature* 505(7482):208–11. doi:[10.1038/nature12772](https://doi.org/10.1038/nature12772)
- Frisch M, Thiemann A, Fu J, Ta Schrag, Scholten S, Melchinger AE (2010) Transcriptome-based distance measures for grouping of germplasm and prediction of hybrid performance in maize. *Theor Appl Genet* 120(2):441–450. doi:[10.1007/s00122-009-1204-1](https://doi.org/10.1007/s00122-009-1204-1)
- Ganal MW, Durstewitz G, Polley A, Bérard A, Buckler ES, Charcoset A, Clarke JD, Graner EM, Hansen M, Joets J, Le Paslier MC, McMullen MD, Montalent P, Rose M, Schön CC, Sun Q, Walter H, Martin OC, Falque M (2011) A large maize (*Zea mays* L.) SNP genotyping array: development and germplasm genotyping, and genetic mapping to compare with the B73 reference genome. *PloS One* 6(12):e28,334. doi:[10.1371/journal.pone.0028334](https://doi.org/10.1371/journal.pone.0028334)
- Geiger HH, Melchinger AE, Schmidt G (1986) Analysis of factorial crosses between flint and dent maize inbred lines for forage performance and quality traits. In: Dolstra O, Miedema P (eds) Breeding of silage maize. Pudoc, Wageningen, pp 147–154
- Genz A, Bretz F, Miwa T, Mi X, Leisch F, Scheipl F, Hothorn T (2017) mvtnorm: multivariate normal and t distributions. <http://cran.r-project.org/package=mvtnorm>
- Gerke JP, Edwards JW, Guill KE, Ross-Ibarra J, McMullen MD (2015) The genomic impacts of drift and selection for hybrid performance in maize. *Genetics* 201(3):1201–1211. doi:[10.1534/genetics.115.182410](https://doi.org/10.1534/genetics.115.182410). arXiv:1307.7313
- Gottesman II, Gould TD (2003) The endophenotype concept in psychiatry: etymology and strategic intentions. *Am J Psychiatry* 160(4):636–645. doi:[10.1176/appi.ajp.160.4.636](https://doi.org/10.1176/appi.ajp.160.4.636)
- Grieder C, Mittweg G, Dhillon B, Montes J, Orsini E, Melchinger AE (2011) Determination of methane fermentation yield and its kinetics by near infrared spectroscopy and chemical composition in maize. *J Near Infrared Spectrosc* 19(6):463–477
- Grieder C, Dhillon BS, Schipprack W, Melchinger AE (2012) Breeding maize as biogas substrate in Central Europe: II. Quantitative-genetic parameters for inbred lines and correlations with testcross performance. *Theor Appl Genet* 124(6):981–988. doi:[10.1007/s00122-011-1762-x](https://doi.org/10.1007/s00122-011-1762-x)
- Guo M, Ma Rupe, Yang X, Crasta O, Zinselmeier C, Smith OS, Bowen B (2006) Genome-wide transcript analysis of maize hybrids: allelic additive gene expression and yield heterosis. *Theor Appl Genet* 113(5):831–845. doi:[10.1007/s00122-006-0335-x](https://doi.org/10.1007/s00122-006-0335-x)
- Guo Z, Magwire MM, Basten CJ, Xu Z, Wang D (2016) Evaluation of the utility of gene expression and metabolic information for genomic prediction in maize. *Theor Appl Genet* 129(12):2413–2427. doi:[10.1007/s00122-016-2780-5](https://doi.org/10.1007/s00122-016-2780-5)
- Hadfield JD (2010) MCMC methods for multi-response generalized linear mixed models: the MCMCglmm R package. *J Stat Softw* 33(2):1–22. doi:[10.1002/ana.22635](https://doi.org/10.1002/ana.22635). arXiv:1501.0228
- Hall BD, Fox R, Zhang Q, Baumgarten A, Nelson B, Cummings J, Drake B, Phillips D, Hayes K, Beatty M, Zastrow-Hayes G, Zeka B, Hazebrook J, Smith S (2016) Comparison of

- genotypic and expression data to determine distinctness among inbred lines of maize for granting of plant variety protection. *Crop Sci* 56(4):1443–1459
- Henderson C (1985) Best linear unbiased prediction of nonadditive genetic merits in noninbred populations. *J Anim Sci* 60:111–117
- Henderson CR (1984) Applications of linear models in animal breeding models. University of Guelph, Guelph
- Heuer C (2015) cpgen: parallelized genomic prediction and GWAS. <https://cran.r-project.org/package=cpgen>
- Hill WG, Goddard ME, Visscher PM (2008) Data and theory point to mainly additive genetic variance for complex traits. *PLoS Genet* 4(2):1–10. doi:10.1371/journal.pgen.1000008
- Jarquín D, Crossa J, Lacaze X, Du Cheyron P, Daucourt J, Lorgeou J, Piraux F, Guerreiro L, Pérez P, Calus M, Burgueño J, de los Campos G (2014) A reaction norm model for genomic selection using high-dimensional genomic and environmental data. *Theor Appl Genet* 127(3):595–607. doi:10.1007/s00122-013-2243-1
- Jiang Y, Reif JC (2015) Modelling epistasis in genomic selection. *Genetics* 201(2):759–768. doi:10.1534/genetics.115.177907
- Kacser H, Burns JA (1981) The molecular basis of dominance. *Genetics* 97:639–666
- Kadam D, Potts S, Bohn MO, Lipka AE, Lorenz A (2016) Genomic prediction of hybrid combinations in the early stages of a maize hybrid breeding pipeline. *G3*(6):3443–3453. doi:10.1101/054015
- Kang HM, Sul JH, Service SK, Zaitlen Na, Kong SY, Freimer NB, Sabatti C, Eskin E (2010) Variance component model to account for sample structure in genome-wide association studies. *Nat Genet* 42(4):348–354. doi:10.1038/ng.548
- Kelliher T, Starr D, Richbourg L, Chintamanani S, Delzer B, Nuccio ML, Green J, Chen Z, McCuiston J, Wang W, Liebler T, Bullock P, Martin B (2017) MATRILINEAL, a sperm-specific phospholipase, triggers maize haploid induction. *Nature* 542(7639):105–109. doi:10.1038/nature20827
- Kerr MK, Churchill GA (2001) Experimental design for gene expression microarrays. *Biostatistics* 2(2):183–201. doi:10.1093/biostatistics/2.2.183
- Larièpe A, Moreau L, Laborde J, Bauland C, Mezouk S, Décousset L, Mary-Huard T, Fiévet JB, Gallais A, Dubreuil P, Charcosset A (2017) General and specific combining abilities in a maize (*Zea mays* L.) test-cross hybrid panel: relative importance of population structure and genetic divergence between parents. *Theor Appl Genet* 130(2):403–417. doi:10.1007/s00122-016-2822-z
- Lisec J, Römisch-Margl L, Nikoloski Z, Piepho HP, Giavalisco P, Selbig J, Gierl A, Willmitzer L (2011) Corn hybrids display lower metabolite variability and complex metabolite inheritance patterns. *Plant J* 68(2):326–336. doi:10.1111/j.1365-3113X.2011.04689.x
- Longin CFH, Mi X, Würschum T (2015) Genomic selection in wheat: optimum allocation of test resources and comparison of breeding strategies for line and hybrid breeding. *Theor Appl Genet* 128(7):1297–1306. doi:10.1007/s00122-015-2505-1
- Mackay TFC (2014) Epistasis and quantitative traits: using model organisms to study gene-gene interactions. *Nat Rev Genet* 15(1):22–33. doi:10.1038/nrg3627
- Mackay TFC, Stone EA, Ayroles JF (2009) The genetics of quantitative traits: challenges and prospects. *Nat Rev Genet* 10(8):565–77. doi:10.1038/nrg2612
- Martini JWR, Wimmer V, Erbe M, Simianer H (2016) Epistasis and covariance: how gene interaction translates into genomic relationship. *Theor Appl Genet* 129(5):963–976. doi:10.1007/s00122-016-2675-5
- Massman JM, Gordillo A, Lorenzana RE, Bernardo R (2013) Genomewide predictions from maize single-cross data. *Theor Appl Genet* 126(1):13–22. doi:10.1007/s00122-012-1955-y
- Melchinger AE, Gumber RK (1998) Overview of heterosis and heterotic groups in agronomic crops. In: Lamkey K, Staub J (eds) Concepts and breeding of heterosis in crop plants. CSSA, Madison, p 16
- Mele M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, Young TR, Goldmann JM, Pervouchine DD, Sullivan TJ, Johnson R, Segre AV, Djebali S, Niarchou A, Consortium TG, Wright FA, Lappalainen T, Calvo M, Getz G, Dermitzakis ET, Ardlie KG, Guigo R (2015) The human transcriptome across tissues and individuals. *Science* 348(6235):660–665. doi:10.1126/science.aaa0355
- Mrode RA (2014) Linear models for the prediction of animal breeding values, 3rd edn. CABI, Oxfordshire. doi:10.1017/CBO9781107415324.004. arXiv:1011.1669v3
- Patti GJ, Yanes O, Siuzdak G (2012) Metabolomics: the apogee of the omics trilogy. *Nat Rev Mol Cell Biol* 13(4):263–9. doi:10.1038/nrm3314
- Pérez P, de Los Campos G (2014) Genome-wide regression & prediction with the BGLR statistical package. *Genetics* 198(October):483–495. doi:10.1534/genetics.114.164442
- R Core Team (2016) R: a language and environment for statistical computing. <https://www.r-project.org/>
- Reif JC, Gumpert F, Fischer S, Melchinger AE (2007) Impact of interpopulation divergence on additive and dominance variance in hybrid populations. *Genetics* 176(3):1931–1934. doi:10.1534/genetics.107.074146
- Riedelsheimer C, Czedik-Eysenberg A, Grieder C, Lisec J, Technow F, Sulpice R, Altmann T, Stitt M, Willmitzer L, Melchinger AE (2012) Genomic and metabolic prediction of complex heterotic traits in hybrid maize. *Nat Genet* 44(2):217–20. doi:10.1038/ng.1033
- Ritchie MD, Holzinger ER, Li R, Pendergrass SA, Kim D (2015a) Methods of integrating data to uncover genotype-phenotype interactions. *Nat Rev Genet* 16:85–97. doi:10.1038/nrg3868
- Ritchie ME, Silver J, Oshlack A, Holmes M, Diyagama D, Holloway A, Smyth GK (2007) A comparison of background correction methods for two-colour microarrays. *Bioinformatics* 23(20):2700–2707. doi:10.1093/bioinformatics/btm412
- Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK (2015b) Limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* 43(7):e47. doi:10.1093/nar/gkv007
- Robson DS, Powers L, Urquhart NS (1967) The proportion of genetic deviates in the tails of a normal population. *Der Züchter Genet Breed Res* 37(4):205–216. doi:10.1007/BF00329530
- Rudd JJ, Kanyuka K, Hassani-Pak K, Derbyshire M, Andongabo A, Devonshire J, Lysenko A, Saqi M, Desai NM, Powers SJ, Hooper J, Ambroso L, Bharti A, Farmer A, Hammond-Kosack KE, Dietrich RA, Courbot M (2015) Transcriptome and metabolite profiling of the infection cycle of *Zymoseptoria tritici* on wheat reveals a biphasic interaction with plant immunity involving differential pathogen chromosomal contributions and a variation on the hemibiotrophic lifestyle def. *Plant Physiol* 167(3):1158–1185. doi:10.1104/pp.114.255927
- Sackton TB, Hartl DL (2016) Perspective genotypic context and epistasis in individuals and populations. *Cell* 166:279–287. doi:10.1016/j.cell.2016.06.047
- Schmid M, Davison TS, Henz SR, Pape UJ, Demar M, Vingron M, Schölkopf B, Weigel D, Lohmann JU (2005) A gene expression map of *Arabidopsis thaliana* development. *Nat Genet* 37(5):501–506. doi:10.1038/ng1543
- Schnell F (1965) Die Covarianz zwischen Verwandten in einer gen-orthogonalen population. I. Allgemeine Theorie. *Biom Z* 7(1):2–49

- Schopp P, Müller D, Technow F, Melchinger AE (2017) Accuracy of genomic prediction in synthetic populations depending on the number of parents, relatedness and ancestral linkage disequilibrium. *Genetics* 205:441–454. doi:[10.1534/genetics.116.193243](https://doi.org/10.1534/genetics.116.193243)
- Searle BC, Gittelman RM, Manor O, Akey JM (2016) Detecting sources of transcriptional heterogeneity in large-scale RNA-Seq data sets. *Genetics* 204(December):1391–1396. doi:[10.1534/genetics.116.193714](https://doi.org/10.1534/genetics.116.193714)
- Smyth G (2004) Linear models and empirical bayes methods for assessing differential expression in microarrays experiments. *Stat Appl Genet Mol Biol* 3(1):1–26
- Smyth GK, Speed T (2003) Normalization of cDNA microarray data. *Methods* 31(4):265–273. doi:[10.1016/S1046-2023\(03\)00155-5](https://doi.org/10.1016/S1046-2023(03)00155-5)
- Springer NM, Stupar RM (2007a) Allele-specific expression patterns reveal biases and embryo-specific parent-of-origin effects in hybrid maize. *Plant Cell* 19(8):2391–2402. doi:[10.1105/tpc.107.052258](https://doi.org/10.1105/tpc.107.052258)
- Springer NM, Stupar RM (2007b) Allelic variation and heterosis in maize: how do two halves make more than a whole? *Genome Res* 17(3):264–275. doi:[10.1101/gr.5347007](https://doi.org/10.1101/gr.5347007)
- Stuber CW, Cockerham CC (1966) Gene effects and variances in hybrid populations. *Genetics* 54(6):1279–1286
- Stupar RM, Gardiner JM, Oldre AG, Haun WJ, Chandler VL, Springer NM (2008) Gene expression analyses in maize inbreds and hybrids with varying levels of heterosis. *BMC Plant Biol* 8(33):1–19. doi:[10.1186/1471-2229-8-33](https://doi.org/10.1186/1471-2229-8-33)
- Swanson-Wagner RA, Jia Y, DeCook R, Borsuk LA, Nettleton D, Schnable PS (2006) All possible modes of gene action are observed in a global comparison of gene expression in a maize F1 hybrid and its inbred parents. *Proc Natl Acad Sci USA* 103(18):6805–6810. doi:[10.1073/pnas.0510430103](https://doi.org/10.1073/pnas.0510430103)
- Technow F, Schrag TA, Schipprack W, Bauer E, Simianer H, Melchinger AE (2014) Genome properties and prospects of genomic prediction of hybrid performance in a breeding program of maize. *Genetics* 197:1343–1355. doi:[10.1534/genetics.114.165860](https://doi.org/10.1534/genetics.114.165860)
- Thiemann A, Fu J, Ta Schrag, Melchinger AE, Frisch M, Scholten S (2010) Correlation between parental transcriptome and field data for the characterization of heterosis in *Zea mays* L. *Theor Appl Genet* 120(2):401–413. doi:[10.1007/s00122-009-1189-9](https://doi.org/10.1007/s00122-009-1189-9)
- Thiemann A, Fu J, Seifert F, Grant-Downton RT, Ta Schrag, Pospisil H, Frisch M, Melchinger AE, Scholten S (2014) Genome-wide meta-analysis of maize heterosis reveals the potential role of additive gene expression at pericentromeric loci. *BMC Plant Biol* 14(88):1–14. doi:[10.1186/1471-2229-14-88](https://doi.org/10.1186/1471-2229-14-88)
- Tzin V, Fernandez-Pozo N, Richter A, Schmelz EA, Schoettner M, Schäfer M, Ahern KR, Meihls LN, Kaur H, Huffaker A, Mori N, Degenhardt J, Mueller LA, Jander G (2015) Dynamic maize responses to aphid feeding are revealed by a time series of transcriptomic and metabolomic assays. *Plant Physiol* 169(November):1727–1743. doi:[10.1104/pp.15.01039](https://doi.org/10.1104/pp.15.01039)
- VanRaden PM (2008) Efficient methods to compute genomic predictions. *J Dairy Sci* 91(11):4414–4423. doi:[10.3168/jds.2007-0980](https://doi.org/10.3168/jds.2007-0980)
- Vazquez AI, Veturi YC, Behring M, Shrestha S, Kirst M, Resende MF Jr, de los Campos G (2016) Increased proportion of variance explained and prediction accuracy of survival of breast cancer patients with use of whole-genome multi-omic profiles. *Genetics* 203(3):1425–1438. doi:[10.1534/genetics.115.185181](https://doi.org/10.1534/genetics.115.185181)
- Vinuela A, Snoek LB, Riksen JAG, Kammenga JE (2010) Genome-wide gene expression regulation as a function of genotype and age in *C. elegans*. *Genome Res* 20:929–937. doi:[10.1101/gr.102160.109](https://doi.org/10.1101/gr.102160.109)
- Wedzony M, Forster B, Zur I, Golemic E, Szechynska-Hebda M, Dubas E, Gotebiowska G (2009) Progress in doubled haploid technology in higher plants. In: Touarev A, Forster BP, Mohan JS (eds) *Advances in haploid production in higher plants*. Springer, Berlin, pp 1–33
- Witt S, Galicia L, Lisek J, Cairns J, Tiessen A, Araus JL, Palacios-Rojas N, Fernie AR (2012) Metabolic and phenotypic responses of greenhouse-grown maize hybrids to experimentally controlled drought stress. *Mol Plant* 5(2):401–17. doi:[10.1093/mp/ssr102](https://doi.org/10.1093/mp/ssr102)
- Xu S, Xu Y, Gong L, Zhang Q (2016) Metabolomic prediction of yield in hybrid rice. *Plant J* 88(2):219–227. doi:[10.1111/tpj.13242](https://doi.org/10.1111/tpj.13242)
- Yang J, Huang T, Petralia F, Long Q, Zhang B, Argmann C, Zhao Y, Mobbs CV, Schadt EE, Zhu J, Tu Z (2015) Synchronized age-related gene expression changes across multiple tissues in human and the link to complex diseases. *Nat Sci Rep* 5(15):145. doi:[10.1038/srep15145](https://doi.org/10.1038/srep15145)
- Zenke-Philippi C, Thiemann A, Seifert F, Schrag T, Melchinger AE, Scholten S, Frisch M (2016) Prediction of hybrid performance in maize with a ridge regression model employed to DNA markers and mRNA transcription profiles. *BMC Genom* 17(1):262. doi:[10.1186/s12864-016-2580-y](https://doi.org/10.1186/s12864-016-2580-y)
- Zhao Y, Li Z, Liu G, Jiang Y, Maurer HP, Würschum T, Mock HP, Matros A, Ebmeyer E, Schachschneider R, Kazman E, Schacht J, Gowda M, Longin CFH, Reif JC (2015) Genome-based establishment of a high-yielding heterotic pattern for hybrid wheat breeding. *Proc Natl Acad Sci USA* 112(51):15,624–15,629. doi:[10.1073/pnas.1514547112](https://doi.org/10.1073/pnas.1514547112)
- Zhu J, Sova P, Xu Q, Dombek KM, Xu EY, Vu H, Tu Z, Brem RB, Bumgarner RE, Schadt EE (2012) Stitching together multiple data dimensions reveals interacting metabolomic and transcriptomic networks that modulate cell regulation. *PLoS Biol* 10(4):e1001301. doi:[10.1371/journal.pbio.1001301](https://doi.org/10.1371/journal.pbio.1001301)