

1 **Modelling $G \times E$ with historical weather information im-** 2 **proves genomic prediction in new environments**

3 Jussi Gillberg^{1*}, Pekka Marttinen¹, Hiroshi Mamitsuka^{1,2} & Samuel Kaski^{1*}

4 1. Helsinki Institute for Information Technology HIIT, Department of Computer Science, Aalto
5 University, PO Box 15400, 00076 Aalto, Finland.

6 2. Institute for Chemical Research, Kyoto University, Gokasho, Uji 6110011, Japan

7 **Interaction between the genotype and the environment ($G \times E$) has a strong impact on the**
8 **yield of major crop plants. Although influential, taking $G \times E$ explicitly into account in plant**
9 **breeding has remained difficult. Recently $G \times E$ has been predicted from environmental and**
10 **genomic covariates, but existing works have not shown that generalization to new environ-**
11 **ments and years without access to in-season data is possible and practical applicability re-**
12 **mains unclear. Using data from a Barley breeding program in Finland, we construct an**
13 **in-silico experiment to study the viability of $G \times E$ prediction under practical constraints. We**
14 **show that the response to the environment of a new generation of untested Barley cultivars**
15 **can be predicted in new locations and years using genomic data, machine learning and his-**
16 **torical weather observations for the new locations. Our results highlight the need for models**
17 **of $G \times E$: non-linear effects clearly dominate linear ones and the interaction between the soil**
18 **type and daily rain is identified as the main driver for $G \times E$ for Barley in Finland. Our study**
19 **implies that genomic selection can be used to capture the yield potential in $G \times E$ effects for**
20 **future growth seasons, providing a possible means to achieve yield improvements, needed for**

21 **feeding the growing population.**

22 Global yield improvements are needed to feed the growing population ¹. One possibility is to
23 breed varieties for higher environmental adaptability, known as *targeted breeding* ². By improving
24 the genetic fit of varieties in their growth environments, yield potential in the interaction between
25 the genotype and environment could be realised. While the importance of $G \times E$ for agronomic
26 performance is widely accepted, utilisation calls for methods that predict yields in new environ-
27 ments, because actual experimental data, consisting of yields of plant variety candidates from yield
28 trials, will in practice be available from only a very limited number of environments. Importantly,
29 prediction of a plant's response to a new environment cannot be based on weather data from the
30 growth season, as those will never be available at the time of prediction.

31 Methods for “cold start” prediction problems ³, where predictions are needed for completely
32 novel instances, have been developed within the machine learning community. Example appli-
33 cations include design of novel drugs for previously unseen cancers ⁴, and recommendations in
34 on-line shopping for new customers and/or products ³. These methods are based on using ex-
35 ternal covariate data that describe properties of the novel instances. We develop a new method,
36 an extension of the Kernelized Bayesian Matrix Factorization ⁵, to account for the uncertainty in
37 the covariates, which allows the use of historical records to predict weather conditions for future
38 growth seasons, and eventually makes future $G \times E$ prediction for yield possible. Therefore, our
39 new method, unlike the existing alternatives ⁶⁻⁸, does not rely on accurate weather information
40 from the growth season from the new location (Figure 1).

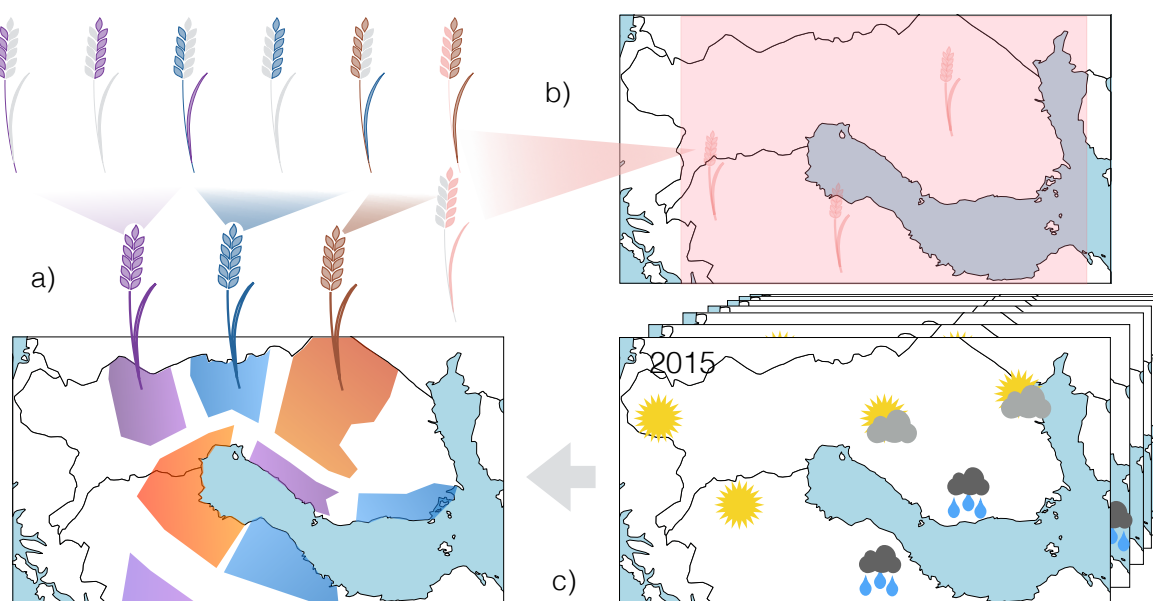


Figure 1: Outline of our approach. a) Precision breeding aims at producing varieties that are optimal for a specific environment. As compared to traditional breeding (b), targeted breeding aims at higher environmental adaptation, i.e., smaller target environments. Weather (microclimate) is a crucial driver for agronomic performance, but as it is unknown for future growth seasons, we use historical weather records (c) to predict the environmental stresses. The growth locations differ with respect to their estimated probabilities of extreme conditions and our method can be used to manage risks by trading-off yield potential for stress tolerance, when the risk in a particular environment is elevated.

41 In genomic selection (GS)⁹, field trials are replaced with genomic predictions to speed-up
42 plant breeding. We formulate an *in silico* experimental setup for GS in targeted breeding that,
43 unlike existing works^{7,8,10-12}, strictly satisfies all realistic constraints: test locations, years, and
44 genotypes are all genuinely new (not part of the training set) and yields are predicted for the off-
45 spring of the training set. In this setup, we demonstrate the feasibility of targeted breeding by in-
46 vestigating the accuracy of $G \times E$ prediction using environmental data including historical weather
47 information but without in-season data (Model M_{G+E+GE}^{hist}). We compare this with multiple com-
48 peting settings, including the non-realistic ideal situation having in-season data (M_{G+E+GE}), a

49 model without the $G \times E$ interaction (M_{G+E}), a previous implementation with $G \times E$ interactions
50 using in-season data (GE-BLUP)⁸, and the industry-standard that does not include $G \times E$ (best
51 linear unbiased prediction using genomic data¹³; GBLUP). Data from a barley breeding program
52 in Finland from Boreal Plant Breeding Ltd, including historical weather information for the target
53 environments, are divided into training, validation, and test sets, and the prediction accuracy is
54 measured as the average correlation between predicted and observed yields in the test sets⁸ (Fig-
55 ure 2c). A sensitivity analysis is done to explore the impact of model assumptions. A description
56 of the model and the setup can be found in Materials and Methods, and further details are given in
57 the Supporting Information.

58 Modelling $G \times E$ with historical weather data, M_{G+E+GE}^{hist} , improves predictive accuracy as
59 compared to the industry-standard, GBLUP (Figure 2a; $p=0.011$, a two-sided paired Wilcoxon
60 signed rank test, $df=17$). The improvement is comparable to using in-season data (M_{G+E+GE} ,
61 $p=0.023$). The Bayesian methods in general show higher accuracy whereas GE-BLUP performs
62 poorly with the data available. Overall, the absolute prediction accuracy of all methods was rel-
63 atively low in this challenging setup, with M_{G+E+GE}^{hist} having the highest correlation of 0.105.
64 Nevertheless, the improvement is considerable over the industry-standard with correlation 0.077,
65 with the proposed new method explaining 85% more of the variation of the phenotype on average.

66 The sensitivity analysis demonstrates considerable variability between test environments
67 (Figure 2c). Indeed, including $G \times E$ interaction terms into the model decreased accuracy in 1/18
68 environments, had little effect in 11/18 environments, but improved the accuracy substantially in

69 6/18 environments. In the last group, increasing model complexity by adding more $G \times E$ com-
70 ponents consistently improved performance, which highlights the potential to increase accuracy
71 through complex modelling of $G \times E$. Importance of different data sources to the predictions can
72 be further analysed by investigating the weights of the different kernels, used to summarise the data
73 sources (Figure S 3). We see that the two most influential kernels were the ones that represented 1)
74 the non-linear interaction between soil type and daily rainfall, and 2) the non-linear effect of rain,
75 matching well the biological understanding of the problem.

76 Our experiments confirmed that prediction in new environments is a challenging task, as
77 reported earlier ⁸, our method reaching the highest correlation of 0.105 between predictions and
78 observations. Nevertheless, the usefulness of including multiple $G \times E$ interaction terms and non-
79 linear interactions between environmental covariates became evident from our results. We expect
80 that gains from modelling $G \times E$ will increase in the future as more data, representing further loca-
81 tions and years, will allow more accurately distinguishing the interactions from the main effects.
82 Other ways to improve the predictions include using more detailed genomic modeling, e.g. using
83 Gaussian and other kernels for summarizing the SNP data.

84 Besides targeted breeding, there are several other needs for $G \times E$ prediction models. They
85 could mitigate the problems of conventional breeding: accounting for historical weather in the
86 actual target population of environments can help prevent overfitting to the conditions in the few
87 field trials performed, as discussed in detail in SI Gains from modelling $G \times E$ for current target
88 population of environments. The assumption of the match between field trials and actual growing

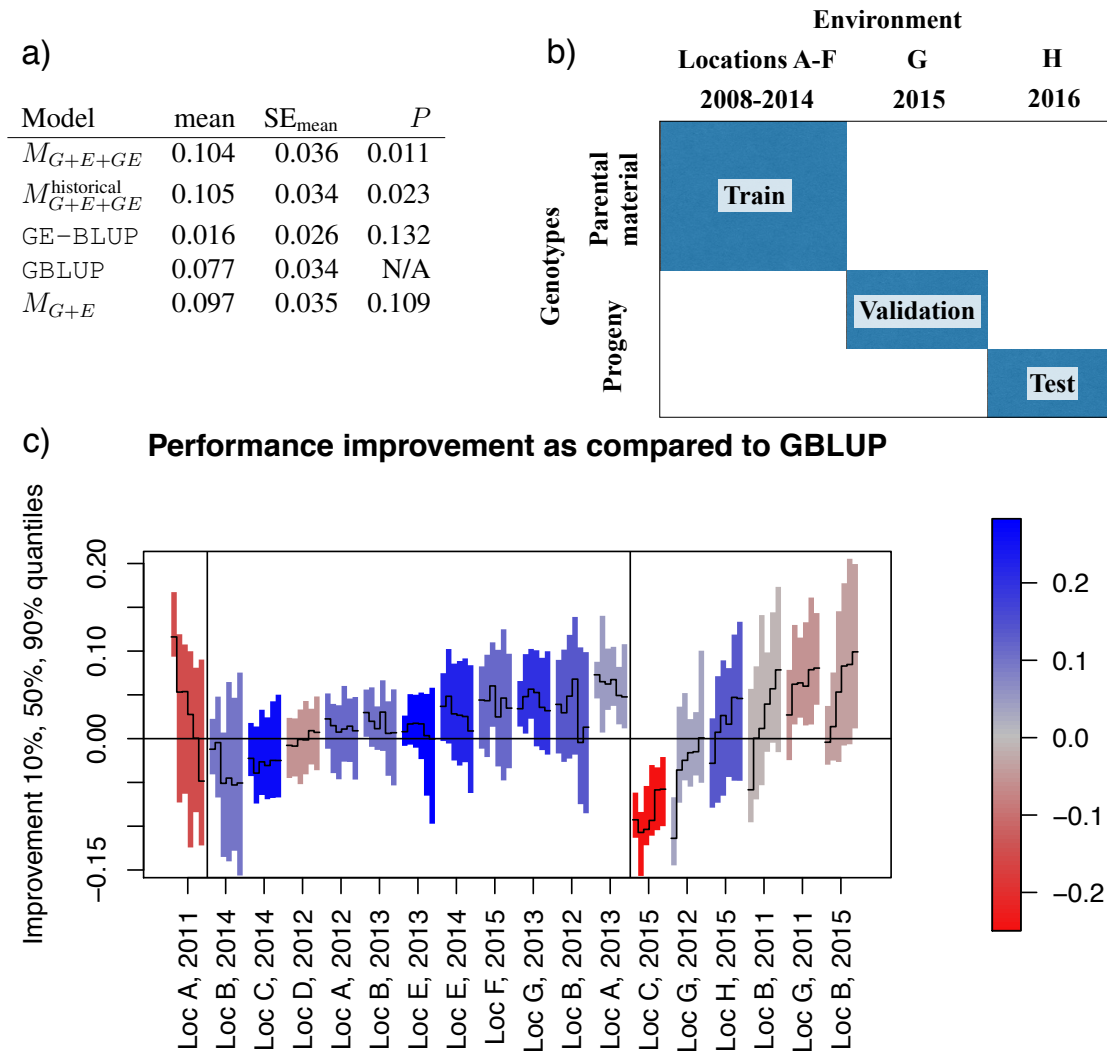


Figure 2: Predicting $G \times E$ with historical weather information improves genomic prediction accuracy in strictly new environments. **a)** Comparison of prediction accuracies; *mean*: correlation between predicted and observed yields, averaged across test environments; SE_{mean} : standard error of the mean; P : p-value compared to the industry standard (GBLUP). **b)** Outline of the *in silico* setup for comparing methods. **c)** Sensitivity analysis: the difference in prediction accuracies (y-axis) between $G \times E$ prediction with historical data (M_{G+E+GE}^{hist}) and the industry standard (GBLUP) is shown in 18 different environments (x-axis); values above the horizontal line mean that M_{G+E+GE}^{hist} is more accurate. Six vertical bars are shown for each environment, representing variability in results (median and 90 % confidence intervals). Starting from the left, they correspond to models with 0, 1, 2, 3, 4 or 5 $G \times E$ interaction terms (0 corresponds to the M_{G+E} model). The color indicates the performance of GBLUP in the environment, red meaning GBLUP performed poorly (Loc C, 2015 were omitted from the comparison as all methods performed poorly there). Vertical lines divide the environments into three groups: $G \times E$ decreased: including $G \times E$ terms to the model decreased performance; $G \times E$ neutral: 10 environments where $G \times E$ terms had neutral effect; $G \times E$ increased: 6 environments where performance increased by adding more $G \times E$ terms.

89 locations is equally crucial for the official variety trials for value of cultivation and use (VCU),
90 required in most countries to evaluate new varieties. $G \times E$ models are also needed in assessing the
91 effects of climate change and to select for varieties that react favourably to the altering conditions
92 ¹. For this purpose, the historical weather observations in M_{G+E+GE}^{hist} can be replaced with climate
93 simulations to assess the performance of varieties under various climate scenarios. To summarise,
94 we showed that $G \times E$ prediction in the setup required by targeted breeding, where the environ-
95 ments are strictly new and predictions are based on historical weather data available at the time of
96 prediction, improves prediction accuracy significantly compared to the industry standard, which is
97 needed to accelerate the implementation of targeted breeding.

98 **Methods**

99 **Data** All data used in the experiment come from a barley breeding program in Finland, which is
100 a part of a larger population of target environments for barley as varieties used in Finland are also
101 used in other Nordic countries. The phenotype consists of (z-transformed) yield measurements
102 (kg/ha) for 2,244 lines observed in trials at 11 locations across the 4 southernmost growth zones
103 in Finland from 2008 to 2015. The total number of observed *location* \times *year* combinations is
104 35. In some locations, trials have been performed on several years and several fields with varying
105 soil properties, and a total of 12,277 yield observations have been recorded. The number of ob-
106 servations per genetic line ranges from 1 to 118 (median 4). The lines were genotyped with the
107 Illumina 9k iSelect SNP Chip, SNPs with minor allele frequency (MAF) < 0.05 or with $> 5\%$
108 values missing were omitted. Also all genotypes with $> 5\%$ of SNPs missing were omitted. The

109 final proportion of missing genotype data is 0.002.

110 The soil characteristics for each field block are measured in terms of the proportions of sand,
111 silt and clay (*soil classification triangle*¹⁴) and the proportion of organic content. Meteorological
112 information consists of daily averages of temperature and rainfall, and the distances to the closest
113 meteorological station range from 1 to 40 km (average 13.5 km). The baseline approach GE-BLUP
114 ⁸ requires summarising the weather information per crop stage: vegetative (from sowing to visible
115 awns), heading time (from visible awns to the end of anthesis), and grain filling (from the end
116 of anthesis to maturity). The times of the crop stages are estimated using temperature sum accu-
117 mulation; the details are given in Section Comparison methods. In the weather observations, the
118 proportion of missing values in daily average temperature and rainfall measurements is < 0.0015
119 (max 3 missing values/environment) and < 0.0032 (max 2 missing values/environment), respec-
120 tively.

121 **Experimental setup.** To study prediction accuracy, we use a setup that strictly imposes the
122 realistic constraints related to modelling $G \times E$ in targeted breeding for new locations. Predictions
123 are required for new locations (not part of the experimental grid) and for years for which no phe-
124 notype data are available (to mimic future growth seasons). Additionally, predictions are needed
125 for the offspring of the lines in the training set, which have no phenotype data observations. More
126 details with a summary of the differences between our setup and earlier works are given in SI
127 Details of experimental setup. We measure prediction accuracy using cross-validation, where the
128 training, validation and test sets are selected to enforce the realistic constraints (Figure 2c). As the

129 performance measure for prediction accuracy, we follow the conventional approach, i.e., the Pear-
130 son correlation between the predicted and observed yields in the test set ^{8,10-12}. This correlation
131 is computed for each cross-validation fold in turn, and averaged over the test cases. Similarly to
132 Malosetti et al. ⁸, the test case -specific correlations are transformed into Fisher's z-scores before
133 averaging and back-transformed to obtain the final results. We regress the $G \times E$ interactions on
134 the average characteristics of the growing season: for each yield trial, we use the weather obser-
135 vations from the typical growing season (1st of May until end of August) regardless of the sowing
136 date. This indirect approach allows predicting with historical weather data. When predicting with
137 historical data, the prediction for each genotype is made for each year for which historical weather
138 observations are available, and the median of those is used as the final predicted value.

139 We also carry out a sensitivity analysis that allows studying the impact of modelling assump-
140 tions, such as inclusion of $G \times E$ interaction components to the model. In detail, the sensitivity
141 analysis shows variability (median and 90% interval) in the predictive performance in a given
142 environment (location, year combination) when we vary *i*) the hyperparameter values over their
143 specified ranges, *ii*) the genotype sets that we are predicting, and *iii*) the training set by removing
144 any single training environment.

145 **Model.** In the models M_{G+E} , M_{G+E+GE} and M_{G+E+GE}^{hist} we assume that *i*) the yield y_{ij} of
146 genotype i in environment j is affected by the genotype, the environmental conditions through-
147 out the growing season and the interactions between the two. We assume that *ii*) the response to
148 the environmental properties is non-linear and that *iii*) it may involve interactions between differ-

149 ent environmental properties. For instance, temperature/rainfall either too low or too high reduces
150 yield, and the response to rainfall is also affected by the soil type. We further assume that *iv*) the re-
151 sponses to the environmental conditions are highly polygenic. Assumptions *i-iv* are encoded using
152 the *kernel trick* ¹⁵, in which covariate data are represented as similarities, or kernels, between dif-
153 ferent data items. Kernel methods are a computationally effective way to model non-linearities and
154 interactions and they have been applied to breeding data ¹⁶. An additional complication in the data
155 is the low number of observed trials compared to the complexity of the problem. To handle this, we
156 constrain our model to only learn the most prominent combinations of environmental conditions
157 affecting yield, by assuming a low-rank approximation for the model parameters accounting for the
158 $G \times E$ effects. Finally, we follow the Bayesian statistical framework ¹⁷, and regularise the model
159 by placing priors on all parameters, which alleviates overfitting to the training data and improves
160 prediction accuracy in the test data.

161 Mathematically, the model for yield is formulated as

$$y_{ij} = g_i + e_j + \xi_{ij} + \epsilon_{ij}, \quad i = 1, \dots, N_g, j = 1, \dots, N_e, \quad (1)$$

162 where g_i is the genetic main effect, e_j is the environmental effect, ξ_{ij} is the effect that arises from
163 interaction between genotype i and environment j , ϵ_{ij} is noise distributed as $N(0, \sigma_j^2)$, and N_g and
164 N_e are the numbers of genotypes and environments. The genetic main effect g_i is modeled as a
165 linear function of the genomic covariates. In detail, the model for the vector of genetic main effects

166 $\mathbf{g}^* = (g_1, \dots, g_{N_g})^T$ is given in terms of a linear genomic kernel K_g by

$$\underset{N_g \times 1}{\mathbf{g}^*} = \underset{N_g \times N_g}{K_g} \cdot \underset{N_g \times 1}{\mathbf{a}_{g0}} + \underset{N_g \times 1}{\mathbf{e}_{g0}}, \quad (2)$$

167 where \mathbf{a}_{g0} are kernel regression weights and \mathbf{e}_{g0} is the noise vector with elements distributed inde-
 168 pendently as $N(0, \sigma_{g0}^2)$. The dimension of each matrix is shown in equation (2) below the corre-
 169 sponding matrix symbol. The genomic kernel K_g is computed by first concatenating the genomic
 170 covariates \mathbf{g}_i as the rows of a matrix \mathbf{G} and then using the standard linear kernel, $K_g = \mathbf{G}\mathbf{G}^T$.

171 The environmental main effect e_j in equation (1) is modeled as a random effect,

$$e_j \sim N(0, \sigma_{e0}^2), \quad j = 1, \dots, N_e.$$

172 The $G \times E$ terms ξ_{ij} are modelled as non-linear functions of the genomic and environmental
 173 covariates, \mathbf{g}_i and \mathbf{e}_j . Each environment and genotype is first represented by R latent variables.
 174 The interactions ξ_{ij} are modelled as the inner product of the latent variable vectors corresponding
 175 to genotype i and environment j , that is,

$$\xi_{ij} = \sum_{r=1}^R h_{ir}^g \cdot h_{jr}^e, \quad i = 1, \dots, N_g, j = 1, \dots, N_e. \quad (3)$$

176 Here, h_{ik}^g is the k th latent variable for the i th genotype, and h_{jk}^e is the k th latent variable for the j th

177 environment. Using matrix notation, equation (3) can be written as

$$\Xi_{N_g \times N_e} = \begin{matrix} H_g & \cdot & H_e^T \\ N_g \times R & & R \times N_e \end{matrix}, \quad (4)$$

178 where $\Xi = [\xi_{ij}]$ is the matrix of interaction terms, and $H_g = [h_{ij}^g]$ and the $H_e = [h_{ij}^e]$ are matri-
 179 ces having as their rows the R -dimensional latent variable representations for each genotype and
 180 environment, respectively.

The latent variables H_g and H_e are obtained from genotype and environment kernels K_g and K_e :

$$\begin{aligned} H_g &= \begin{matrix} K_g & \cdot & A_g & + & E_{H_g} \\ N_g \times R & & N_g \times N_g & & N_g \times R & & N_g \times R \end{matrix} \quad \text{and} \\ H_e &= \begin{matrix} K_e & \cdot & A_e & + & E_{H_e} \\ N_e \times R & & N_e \times N_e & & N_e \times R & & N_e \times R \end{matrix}, \end{aligned}$$

181 where A_g and A_e are kernel regression weights, and E_{H_g} and E_{H_e} are matrices containing error
 182 terms distributed independently as $N(0, \sigma_g^2)$ or $N(0, \sigma_e^2)$, respectively. The environmental ker-
 183 nel K_e is obtained by combining multiple kernels K_e^1, \dots, K_e^E , computed from environmental data
 184 $\mathbf{e}_j, j = 1, \dots, N_e$, each kernel representing a different aspect of the environment (weather, soil,
 185 etc). Details about processing the raw data into kernels and about combining multiple environmen-
 186 tal kernels into a single kernel are presented in SI Data processing and kernels.

187 For inference we use variational approximation¹⁸, which is a computationally feasible way
 188 to approximate posterior distributions of parameters in complex models. The variational updates

189 required here can be derived similarly to Gönen et al. ⁵, except that we have extended their model
190 and algorithm by including the genotype and environment main effects, i.e., the terms g_i and e_j
191 in equation (1). Detailed distributions of the model parameters and the guidelines for specifying
192 hyperparameter values are given in Sections SI Detailed model specification and SI Specifying
193 hyperparameter values, respectively. Further details about the inference algorithm in SI Details of
194 the variational inference algorithm.

195 **Comparison methods.** The mixed model computations for the comparison methods GBLUP
196 and GE-BLUP are performed using the R library `rBLUP` ¹⁹. For both methods, fixed effects
197 were used to account for field block-specific effects, corresponding to the terms e_j in M_{G+E+GE}^{hist} ,
198 M_{G+E+GE} and M_{G+E} . For GBLUP, the genomic kernel (see Section Model) was used as the
199 covariance matrix Σ . For GE-BLUP, the environmental kinship model (GE-KE) ⁸, is used and the
200 full covariance matrix Σ is generated through the Kronecker product $\Sigma = \Sigma_G \otimes \Sigma_E$, where Σ_G
201 and Σ_E are the genetic and environmental covariance matrices, respectively. The environmental
202 covariance matrix Σ_E is generated from the available environmental data to describe soil properties
203 and the growth conditions during the vegetative, heading time and grain filling developmental
204 stages. All soil data and growth zone information are used as such whereas the daily average
205 temperature and rainfall measurements are summarised as the mean and the standard deviation of
206 the daily observations per crop stage. The growth periods are estimated using the sowing date
207 and temperature sum accumulation-based estimates of heading and ripening times (440.2 °C and
208 905.9 °C, respectively), which were estimated from external breeding data. The vegetative stage
209 is assumed to last 3 weeks starting from sowing, the time of heading is assumed to start 2 weeks

210 before and last 1 week after the estimated heading time and grain filling was assumed to start
211 after heading and to last 1 week longer than the estimated time of ripening. Wide estimates for the
212 growth periods were used to account for varying growth speeds. The resulting set of environmental
213 covariates is z-normalized and a linear kernel is used, which is further normalized according to
214 equation (()) in SI Data preprocessing and kernels.

215 **Data availability**

216 The data accompanied by the method code will be made available upon publication in the form of
217 kernels to allow reproducing the results.

218 **Acknowledgements**

219 We acknowledge Boreal Plant Breeding Ltd for the access to the plant breeding data and we grate-
220 fully thank Outi Manninen, Mika Isolahti and Esa Teperi for very helpful comments. This work
221 was in part funded by Tekes, the Finnish Funding Agency for Innovation (Dnro 1718/31/2014 to
222 J.G, H.M and S.K) and the Academy of Finland (Finnish Centre of Excellence in Computational
223 Inference Research COIN, and grant numbers 294238, 292334 to SK and grant numbers 286607
224 and 294015 to P.M.).

225 **Author contributions**

226 J.G. processed the data from Boreal Plant Breeding Ltd and performed the *in silico* experiments.

227 J.G. and P.M. implemented the method. All authors were involved in the conception and design of

228 the study, analyzed the results and assisted with drafting and critically revising the manuscript.

229 **Competing interests**

230 The authors declare no competing financial interests.

- 232 1. Tester M, Langridge P (2010) Breeding Technologies to Increase Crop Production in a Chang-
233 ing World. *Science* 327(5967):818–822.
- 234 2. Braun HJ, Rajaram S, Ginkel M (1996) Cimmyt’s approach to breeding for wide adaptation.
235 *Euphytica* 92(1):175–183.
- 236 3. Schein AI, Popescul A, Ungar LH, Pennock DM (2002) Methods and Metrics for Cold-start
237 Recommendations in *Proceedings of the 25th Annual International ACM SIGIR Conference*
238 *on Research and Development in Information Retrieval, SIGIR ’02*. (ACM, New York, NY,
239 USA), pp. 253–260.
- 240 4. Costello JC, et al. (2014) A community effort to assess and improve drug sensitivity prediction
241 algorithms. *Nature biotechnology* 32(12):1202–1212.
- 242 5. Gönen M, Kaski S (2014) Kernelized Bayesian Matrix Factorization. *IEEE Transactions on*
243 *Pattern Analysis and Machine Intelligence* 36(10):2047–2060.
- 244 6. Jarquín D, et al. (2014) A reaction norm model for genomic selection using high-dimensional
245 genomic and environmental data. *Theoretical and Applied Genetics* 127(3):595–607.
- 246 7. Heslot N, Akdemir D, Sorrells ME, Jannink JL (2014) Integrating environmental covariates
247 and crop modeling into the genomic selection framework to predict genotype by environment
248 interactions. *Theoretical and Applied Genetics* 127(2):463–480.
- 249 8. Malosetti M, Bustos-Korts D, Boer MP, van Eeuwijk FA (2016) Predicting responses in mul-
250 tiple environments: issues in relation to genotype × environment interactions. *Crop Science*
251 56(5):2210–2222.

- 252 9. Meuwissen THE, Hayes BJ, Goddard ME (2001) Prediction of Total Genetic Value Using
253 Genome-Wide Dense Marker Maps. *Genetics* 157(4):1819–1829.
- 254 10. Burgueño J, de los Campos G, Weigel K, Crossa J (2012) Genomic prediction of breeding val-
255 ues when modeling genotype× environment interaction using pedigree and dense molecular
256 markers. *Crop Science* 52(2):707–719.
- 257 11. Albrecht T, et al. (2014) Genome-based prediction of maize hybrid performance across genetic
258 groups, testers, locations, and years. *Theoretical and Applied Genetics* 127(6):1375–1386.
- 259 12. Saint Pierre C, et al. (2016) Genomic prediction models for grain yield of spring bread wheat
260 in diverse agro-ecological zones. *Scientific Reports* 6:27312.
- 261 13. de los Campos G, Hickey JM, Pong-Wong R, Daetwyler HD, Calus MPL (2013) Whole-
262 Genome Regression and Prediction Methods Applied to Plant and Animal Breeding. *Genetics*
263 193(2):327–345.
- 264 14. Shepard FP (1954) Nomenclature based on sand-silt-clay ratios. *Journal of Sedimentary*
265 *Petrology* 24(3):151–158.
- 266 15. Shawe-Taylor J, Cristianini N (2004) *Kernel methods for pattern analysis*. (Cambridge Uni-
267 versity Press).
- 268 16. Gianola D, Morota G, Crossa J (2014) Genome-enabled prediction of complex traits with
269 kernel methods: What have we learned? in *Proceedings, 10th World Congress of Genetics*
270 *Applied to Livestock Production*. p. 6.

- 271 17. Gelman A, et al. (2013) Bayesian data analysis, 3rd edition.
- 272 18. Beal MJ (2003) *Variational algorithms for approximate Bayesian inference*. (Gatsby Compu-
273 tational Neuroscience Unit, University College London).
- 274 19. Endelman JB (2011) Ridge Regression and Other Kernels for Genomic Selection with R Pack-
275 age rrBLUP. *The Plant Genome* 4(3):250–255.
- 276 20. Gönen M (2012) Bayesian Efficient Multiple Kernel Learning. *Proc. 29th International Con-
277 ference on Machine Learning, ICML 2012* pp. 1–8.

278 **Supplementary Information (SI)**

279 **Data preprocessing and kernels.** A summary of different kernels, including transformations
280 specific to each data source, preprocessing and kernel transformations used, is given in Table S 1.
281 The bandwidth parameter of all the Gaussian kernels is set to the conventional default value equal
282 to the number of covariates used to compute the kernel. All kernels K are normalized to make
283 them unit diagonal:

$$\tilde{K} = (\mathbf{d}^{-1/2} \times \mathbf{d}^{-1/2}) \cdot K \quad (1)$$

284 where \mathbf{d} is a vector of the diagonal values of kernel K , \times denotes the outer product, and the
285 $\mathbf{d}^{-1/2}$ denotes a vector with all elements of \mathbf{d} raised to the power of $-1/2$. The interaction kernel
286 between the soil type and rainfall is computed from other kernels as

$$K_{\text{soil} \times \text{rain}} = \tilde{K}_{\text{soil, Gaussian}} \odot \tilde{K}_{\text{rain, Gaussian}}, \quad (2)$$

287 where \odot denotes the Hadamard (elementwise) product. Finally, all kernels are normalized with
288 respect to their summed total variance by multiplication with a constant c

$$\tilde{\tilde{K}} = c \cdot \tilde{K} \quad (3)$$

289 where $c = \left[\sum_{i=1}^{N_z} \text{Var}(\tilde{\mathbf{k}}_i) \right]^{-1/2}$ and $\tilde{\mathbf{k}}_i$ is the i th column of \tilde{K} . The motivation for this normalisa-
290 tion comes from the expectation that *a priori* each kernel explains the same amount of variance,
291 and, when combining the kernels as described below, this prior expectation is realised by the nor-

292 malisation.

293 **Combining environmental kernels.** The final environmental kernel K_e is obtained as a
294 weighted sum of the different normalized (\tilde{K}) kernels in Table S 1. The weights are learnt from the
295 training data by fitting BEMKL²⁰, a multiple kernel learning regression method, using experiment-
296 specific yield means as the target variable. For BEMKL, shape (α) and scale (β) parameter values
297 for the prior Gamma distributions are set to 1 except for the λ parameter, for which the scale is
298 fixed to 10, providing stronger regularization. Regression bias term b is set to 0. For further details
299 of BEMKL, see Gönen et al.²⁰. Before combining the kernels, the learnt weights are normalized
300 such that their sum of squares is equal to 1, and the largest weight (in absolute value) is positive.
301 The distributions of the normalised kernel weights from the sensitivity analysis are presented in
302 Figure S 3. The composite kernel K_e is again normalized according to equation (3).

303 **Detailed model specification.** The distributional assumptions of the model are

$$y_{ij}|H_g, H_e, g_i, e_j, \sigma_j^2 \sim \mathcal{N}(g_i + e_j + (\mathbf{h}_i^g)^T \mathbf{h}_j^e, \sigma_j^2), \quad \forall(i, j)$$

$$\sigma_j^{-2} \sim \mathcal{G}(\alpha_j, \beta_j), \quad \forall(j)$$

$$a_i^{g0}|\lambda_{g0} \sim \mathcal{N}(0, \lambda_{g0}^{-1}), \quad \forall(i)$$

$$g_i|\mathbf{a}_{g0}, K_g, \sigma_{g0}^2 \sim \mathcal{N}(\mathbf{a}_{g0}^T \mathbf{k}_i^g, \sigma_{g0}^2), \quad \forall(i)$$

$$a_{ij}^g|\lambda_g \sim \mathcal{N}(0, \lambda_g^{-1}), \quad \forall(i, j)$$

Variable (unit)	transformation	preprocessing parameters	missing value imputation	kernel transformation(s)
Soil content (%, $N_{covs} = 3$)	log transformation	z-normalization	(none)	linear and Gaussian
Soil organic content (%, $N_{covs} = 3$)	log transformation	z-normalization	(none)	Gaussian
daily rainfall (mm, $N_{covs} = 123$)	7-day moving average (6 previous days)	z-normalization with 3 rd order polynomial smoothing	0-imputation	linear and Gaussian
daily average temperature (C°, $N_{covs} = 123$)		z-normalization with 3 rd order polynomial smoothing of daily mean/scale parameters	0-imputation	linear and Gaussian
growth zone (1-4, $N_{covs} = 1$)		z-normalization	(none)	Gaussian
genotype markers (SNPs, $N_{covs} = 5696$)		Minor allele frequency scaling for SNP A: $\frac{A - 2 \cdot MAF_A}{\sqrt{2 \cdot MAF_A \cdot (1 - MAF_A)}}$	mean imputation	linear kernel

Table S 1: Preprocessings and kernel functions applied to covariates.

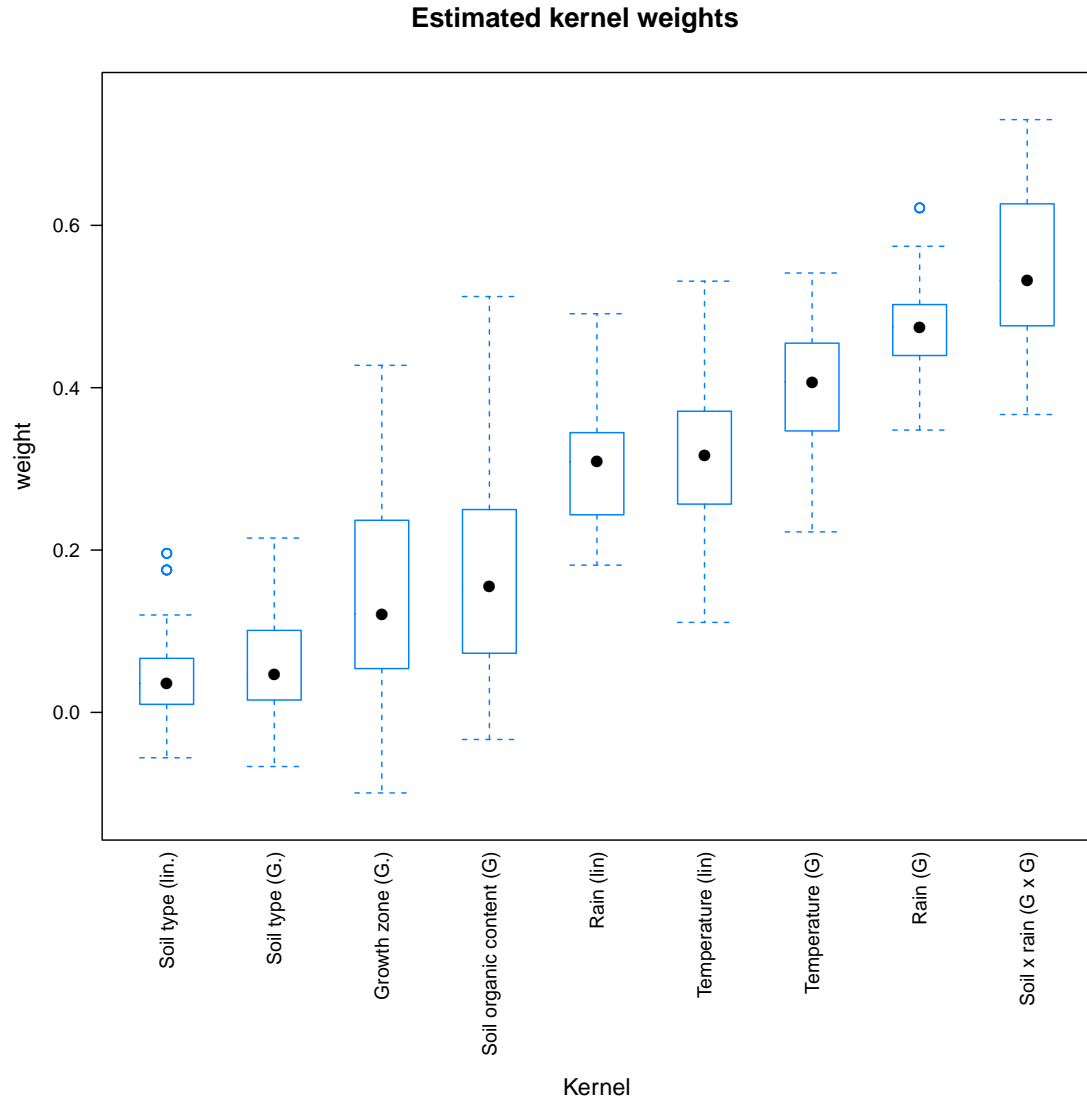


Figure S 3: Estimated normalized kernel weights in the sensitivity analysis.

$$h_{ij}^g | A_g, K_g, \sigma_g^2 \sim \mathcal{N}((\mathbf{k}_i^g)^T \mathbf{a}_j^g, \sigma_g^2), \quad \forall(i, j)$$

$$e_j | \sigma_{e0}^2 \sim \mathcal{N}(0, \sigma_{e0}^2), \quad \forall(j)$$

$$a_{ij}^e | \lambda_e \sim \mathcal{N}(0, \lambda_e^{-1}), \quad \forall(i, j)$$

$$h_{ij}^e | A_e, K_e, \sigma_e^2 \sim \mathcal{N}((\mathbf{k}_i^e)^T \mathbf{a}_j^e, \sigma_e^2), \quad \forall(i, j),$$

304 where $\mathbf{k}_i^g, \mathbf{k}_j^e, \mathbf{a}_j^g, \mathbf{a}_j^e$, denote columns of matrices K_g, K_e, A_g, A_e , with subscripts i and j specifying
305 the column index; \mathbf{h}_i^g and \mathbf{h}_j^e denote i th and j th rows of H_g and H_e , represented as column vectors;
306 a_i^{g0} is the i th element of vector \mathbf{a}_{g0} ; a_{ij}^g and a_{ij}^e are the (i, j) th elements in matrices A_g and A_e . \mathcal{N}
307 and \mathcal{G} denote the Gaussian and Gamma distributions, respectively.

308 **Specifying hyperparameter values.** Prior knowledge about the approximate weights of
309 different sources of variance, e.g. the relative weight of genetic and environmental main effects, is
310 used to specify hyperparameter values. We determine for each hyperparameter either a single fixed
311 value or a grid of values to be selected from by cross-validation. Parameters (α_j, β_j) of the Gamma
312 distribution for environment-specific residual noise variances σ_j^2 are set to $(10, 1)$, corresponding
313 to an expected value of approximately 0.1 for σ_j^2 . The variance of environment mean effects σ_{e0}^2 is
314 fixed to 0.25. To set the parameters λ_{g0} and σ_{g0}^2 that determine the amount of signal and noise in
315 the genetic main effects, we find values for them such that two conditions are satisfied. First, 95%
316 of the variance of the genetic effects \mathbf{g}^* is assumed to be signal, that is,

$$\frac{\text{Var}(K_g \cdot \mathbf{a}_{g0})}{\text{Var}(K_g \cdot \mathbf{a}_{g0}) + \sigma_{g0}^2} = 0.95.$$

317 The second condition is that the variance of the genetic main effects, $\text{Var}(K_g \cdot \mathbf{a}_{g0}) + \sigma_{g0}^2$, is either
318 0.2, 0.4, or 0.6. In practice we find these values for σ_{g0}^2 and λ_{g0} by simulating multiple realisations
319 from the model with specific values for the parameters, and select values that on average satisfy
320 the two conditions.

The parameters $\lambda_g, \sigma_g^2, \lambda_e$, and σ_e^2 , controlling the proportion of signal and noise in the la-

tent components H_g and H_e that model the $G \times E$ interactions, are selected according to similar principles: by inspecting the proportion of signal of the total variance of the latent factors and the relative contribution of the interaction terms compared to the genetic main effects. In detail, we assume first that

$$\frac{\text{Tr}(\text{Var}(K_g \cdot A_g))}{\text{Tr}(\text{Var}(K_g \cdot A_g)) + R\sigma_g^2} = 0.95, \text{ and}$$

$$\frac{\text{Tr}(\text{Var}(K_e \cdot A_e))}{\text{Tr}(\text{Var}(K_e \cdot A_e)) + R\sigma_e^2} = 0.95,$$

321 where $\text{Tr}()$ denotes the trace of a matrix. Second, we assume that the total variance of the interac-
 322 tions is either the same or half of the total variance from the genetic main effects, i.e.

$$\text{Tr}(\text{Var}(H_g \cdot H_e^T)) = \Phi \times R \times [\text{Var}(K_g \cdot \mathbf{a}_{g0}) + \sigma_{g0}^2],$$

323 where Φ is either 0.5 or 1, to be selected with cross-validation.

324 **Details of the variational inference algorithm.** For short-hand, the hyper-parameters in the
 325 model are denoted jointly by

$$\zeta = \{\alpha_j, \beta_j, \sigma_{g0}^2, \sigma_{e0}^2, \sigma_g^2, \sigma_e^2, \lambda_{g0}, \lambda_g, \lambda_e\},$$

326 and the parameters by

$$\Theta = \{\mathbf{a}_{g0}, A_g, A_e, H_g, H_e, \mathbf{g}^*, \mathbf{e}^*, \sigma_*^2\},$$

327 where $\sigma_*^2 = (\sigma_1^2, \dots, \sigma_{N_e}^2)$. In the following the dependence on ζ is omitted for clarity. We assume
 328 the factorized variational approximation

$$p(\Theta|K_g, K_e, Y) \approx q(\Theta) = q(\mathbf{a}_{g0})q(A_g)q(A_e)q(H_g)q(H_e)q(\mathbf{g}^*)q(\mathbf{e}^*)q(\sigma_*^2)$$

and define each factor in the ensemble just like its full conditional:

$$\begin{aligned} q(\mathbf{a}_{g0}) &= \mathcal{N}(\mathbf{a}_{g0}; \mu(\mathbf{a}_{g0}), \Sigma(\mathbf{a}_{g0})) \\ q(A_g) &= \prod_{r=1}^R \mathcal{N}(\mathbf{a}_r^g; \mu(\mathbf{a}_r^g), \Sigma(\mathbf{a}_r^g)) \\ q(A_e) &= \prod_{r=1}^R \mathcal{N}(\mathbf{a}_r^e; \mu(\mathbf{a}_r^e), \Sigma(\mathbf{a}_r^e)) \\ q(H_g) &= \prod_{i=1}^{N_g} \mathcal{N}(\mathbf{h}_i^g; \mu(\mathbf{h}_i^g), \Sigma(\mathbf{h}_i^g)) \\ q(H_e) &= \prod_{j=1}^{N_e} \mathcal{N}(\mathbf{h}_j^e; \mu(\mathbf{h}_j^e), \Sigma(\mathbf{h}_j^e)) \\ q(\mathbf{g}^*) &= \prod_{i=1}^{N_g} \mathcal{N}(g_i; \mu(g_i), \Sigma(g_i)) \\ q(\mathbf{e}^*) &= \prod_{j=1}^{N_e} \mathcal{N}(e_j; \mu(e_j), \Sigma(e_j)) \\ q(\sigma_*^2) &= \prod_{j=1}^{N_e} \mathcal{G}(\sigma_j^{-2}; \alpha(\sigma_j^{-2}), \beta(\sigma_j^{-2})). \end{aligned}$$

329 The parameters in the factor distributions can be derived as by Gönen et al ⁵, and they are therefore
 330 omitted from here.

331 **Initialisation of the variational algorithm.** The parameter \mathbf{g}^* was initialised to the main

332 genetic effects learnt by GBLUP, and \mathbf{e}^* was initialised to the average yields in the different en-
333 vironments. Parameters H_g and H_e were initialised by applying the regularized Singular Value
334 Decomposition (SVD) implemented in R library `softImpute` to the yield matrix Y after re-
335 gressing out the initialised main effects \mathbf{g}^* and \mathbf{e}^* . Parameters \mathbf{a}_{g0} , A_g and A_e were initialised to
336 0. Environment-specific residual variance parameters σ_*^2 were initialised to environment-specific
337 sample variances.

338 **Details of experimental setup.** Different prediction tasks, distinguished by the availability
339 of different data types, are presented in Figure S4. Setups 1-4 correspond to those studied by
340 Malosetti et al. ⁸: in setup 1, phenotype measurements are available for the genotypes and envi-
341 ronments to be predicted, and both genotypes and environmental covariates are fully observed. In
342 setups 2 and 3, phenotype measurements are still available but only for the genotypes or the envi-
343 ronments to be predicted, but not both, and covariates are fully observed. In setup 4, no phenotype
344 data are available for environments/genotypes to be predicted, but both genetic and environmental
345 covariates are still fully observed.

346 Two additional setups can be considered. In setups 5 and 6 environmental covariates from the
347 environments of interest are only partially available: location and soil characteristics are known but
348 the in-season weather measurements are not available for the year of interest. However, historical
349 observations for the same locations are available and they are used to estimate the performance of
350 each genotype. Setups 5 and 6 differ depending on whether phenotype measurements are available
351 from some other environment for the genotypes (5) or not at all (6). The results in this paper are for

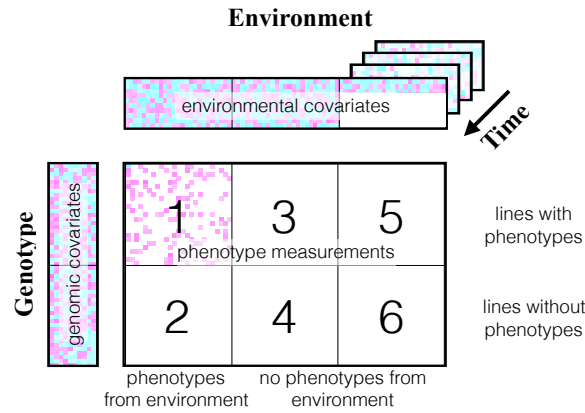


Figure S4: Comparison of prediction setups with respect to the availability of phenotype data and the genomic and environmental covariates as presented by Malosetti et al ⁸. White colour indicates missing value. In setups 1, 3 and 5, "lines with phenotypes", the lines to be predicted have phenotype observations (from some environments). In setups 1 and 2, "phenotypes from environment", phenotypes have been measured from the prediction target environments (for some lines). In setups 1-4 presented by Malosetti et al ⁸, environmental covariates are available for all environments, whereas in the new setups 5 and 6, environmental covariates from the trials of interest are missing and they are replaced by using several years of historical data.

352 setup 6 where no phenotype data are available for any of the lines of interest. We emphasize that
 353 a further difference to earlier work ⁸ is that we strictly require the test environments to be simul-
 354 taneously both from a location and from a year not included among the training environments and
 355 that the genotypes in the test/validation sets are from the progeny of the training set. A summary
 356 of the differences between our setup to those presented by earlier works is given in Table S 2

357 **Gains from modelling $G \times E$ for current target population of environments.** Our results
 358 indicate targeted breeding could improve yields by dividing a single target population of environ-
 359 ments (TPE) into several parts, but the same methodology could be used even when developing
 360 only 1 variety for a larger population of target environments as in traditional breeding. Traditional
 361 breeding makes the implicit assumption that varieties' observed yields $g \in 1, \dots, G$ in trial exper-
 362 iments in environments (location \times year) $e \in 1, \dots, E$, are representative of the yield in the TPE,

Publication	New environment	New genotypes
Burgueño et al. 2012 (CV1/CV2) ¹⁰	CV1/CV2: test locations and years are present in the location-year combinations in the training data	new lines in CV1: not restricted to the offspring generation. In CV2 the test lines have phenotype observations
Heslot et al. 2013 ⁷	Random split, balanced wrt years and locs → year and locations not new	only 544/2195 genotypes have no phenotype observations, test set not restricted to the offspring generation
Albrecht et al. 2014 ¹¹	the year-location combination is new but the test locations and years are present in other location-year combinations in the training data	genotypes are new and from the offspring
Malosetti et al. 2016 ⁸	time-structured DTD: 2/6 test locations new according to strict criteria; physically structured DTD: none of the environments are strictly new (as the year is not new)	all genotypes within the same family, not from the next generation.
Saint Pierre et al. 2016 ¹² (leave-one-side-out)	location new but year part of the training set	test lines have phenotype observations

Table S 2: Comparison of the proposed *in silico* setup to the existing setups.

363 in other words

$$p(\text{yield}_g|\text{TPE}) \approx \frac{1}{E} \sum_e p(\text{yield}_g|\text{environment}_e) \quad (4)$$

However, with geographic field use information and weather data widely available, this strong assumption can be replaced with an estimate for the yield in the TPE given the actual fields and their microclimates:

$$p(\text{yield}_g|\text{TPE}) \approx \sum_f^F P_f \times p(\text{yield}_g|f) \quad (5)$$

$$= \sum_f^F P_f \times \int_{\theta_f} p(\text{yield}_g|\theta_f) \times p(\theta_f) d\theta_f, \quad (6)$$

364 where $f \in 1, \dots, F$, are fields in the TPE used for cultivation of the new variety, θ_f are parameters
365 (e.g. weather conditions) related to a certain field f , $p(\theta_f)$ is the uncertainty related to these
366 conditions, estimated from historical records, $p(\text{yield}_g|\theta_f)$ is the predictive distribution for the
367 yield under conditions θ_f , obtained from the model, and P_f is the proportion of the total volume
368 cultivated in field f .

	environment	GBLUP	M_{G+E}	M_{G+E+GE}	M_{G+E+GE}^{hist}	GE-BLUP	N_{test}
1	Loc B, 2011	-0.147	-0.146	-0.109	-0.07	0.068	59
2	Loc A, 2011	-0.237	-0.086	-0.206	-0.2	-0.015	59
3	Loc G, 2011	-0.02	-0.028	0.045	0.031	0.183	58
4	Loc B, 2013	-0.016	0.057	0.051	0.046	0.182	182
5	Loc A, 2012	0.167	0.198	0.167	0.163	0.057	106
6	Loc G, 2012	0.023	-0.079	0.017	0.031	-0.166	106
7	Loc D, 2012	0.01	-0.031	-0.01	-0.026	-0.075	105
8	Loc E, 2013	0.124	0.2	0.265	0.245	-0.354	91
9	Loc B, 2012	0.089	0.166	0.17	0.195	-0.047	106
10	Loc G, 2013	0.228	0.21	0.278	0.282	-0.256	91
11	Loc B, 2012	0.171	0.223	0.179	0.192	0.127	260
12	Loc A, 2012	0.046	0.113	0.087	0.043	0.1	243
13	Loc E, 2013	0.467	0.488	0.503	0.518	-0.122	153
14	Loc G, 2013	0.2	0.308	0.353	0.344	-0.107	152
15	Loc C, 2014	0.208	0.197	0.068	0.045	0.054	79
16	Loc B, 2013	0.267	0.309	0.291	0.225	0.147	153
17	Loc A, 2013	0.05	0.128	0.116	0.114	0.089	153
18	Loc E, 2014	-0.02	-0.007	0.027	-0.013	0.354	79
19	Loc B, 2014	-0.037	-0.021	-0.046	-0.044	-0.133	79
20	Loc B, 2014	0.258	0.272	0.012	0.054	0.122	106
21	Loc C, 2014	0.344	0.324	0.375	0.345	0.065	106
22	Loc E, 2014	0.325	0.42	0.371	0.404	0.275	105
23	Loc H, 2015	0.199	0.12	0.166	0.122	0.251	64
24	Loc F, 2015	0.246	0.245	0.197	0.295	0.097	64
25	Loc B, 2015	-0.069	-0.093	0.227	0.257	-0.289	64
26	Loc B, 2013	0.091	0.084	0.078	0.083	0.058	488
27	Loc G, 2013	0.181	0.236	0.151	0.138	-0.015	244
28	Loc E, 2013	0.243	0.222	0.234	0.236	-0.052	244
29	Loc C, 2014	0.292	0.36	0.385	0.382	-0.069	120
30	Loc F, 2015	0.053	0.169	0.186	0.161	0.002	39
31	Loc E, 2014	0.232	0.217	0.25	0.208	-0.062	120
32	Loc B, 2015	-0.056	-0.083	-0.062	-0.071	-0.292	39
33	Loc B, 2014	0.01	-0.016	-0.145	-0.137	0.086	91
34	Loc E, 2014	0.221	0.306	0.382	0.367	-0.157	91
35	Loc C, 2014	0.273	0.25	0.232	0.258	0.089	91
36	Loc H, 2015	0.044	0.006	0.159	0.123	-0.206	42
37	Loc B, 2015	-0.041	0.035	0.145	0.128	0.219	42
38	Loc B, 2015	0.095	-0.039	0.008	0.012	0.039	64
39	Loc F, 2015	-0.037	0	-0.039	0.03	0.078	64
40	Loc H, 2015	0.25	0.259	0.261	0.262	0.253	63
41	Loc C, 2015	-0.204	-0.341	-0.353	-0.359	-0.143	60

Table S 3: Results for individual test folds.