

1 The effects of training population design 2 on genomic prediction accuracy in wheat

3 Stefan McKinnon Edwards¹, Jaap B. Buntjer¹, Robert Jackson², Alison R. Bentley², Jacob
4 Lage³, Ed Byrne³, Chris Burt⁴, Peter Jack⁴, Simon Berry⁵, Edward Flatman⁵, Bruno Poupard⁵,
5 Stephen Smith⁶, Charlotte Hayes⁶, R. Chris Gaynor¹, Gregor Gorjanc¹, Phil Howell², Eric Ober²,
6 Ian J. Mackay⁷ and John M. Hickey^{1*}

7

8 ¹ The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of
9 Edinburgh, Easter Bush, Midlothian, Scotland, United Kingdom; ² The John Bingham
10 Laboratory, NIAB, Huntingdon Road Cambridge CB3 0LE UK; ³ KWS UK Ltd, 56 Church
11 Street Hertfordshire, SG8 7RE, UK, ⁴ RAGT UK, Grange Rd, Saffron Walden, CB10 1TA, UK,
12 ⁵ Limagrain UK Ltd, Rothwell, Market Rasen, Lincolnshire LN7 6DT, UK, ⁶ Elsoms Wheat
13 Limited, Pinchbeck Road, Spalding, Lincolnshire, PE11 1QG, UK, ⁷ IMplant Consultancy Ltd.,
14 Chelmsford, UK

15

16 *: Corresponding author (John.Hickey@roslin.ed.ac.uk)

17 Abstract

18 Genomic selection offers several routes for increasing genetic gain or efficiency of plant
19 breeding programs. In various species of livestock there is empirical evidence of increased
20 rates of genetic gain from the use of genomic selection to target different aspects of the
21 breeder's equation. Accurate predictions of genomic breeding value are central to this and the
22 design of training sets is in turn central to achieving sufficient levels of accuracy. In summary,
23 small numbers of close relatives and very large numbers of distant relatives are expected to
24 enable accurate predictions.

25 To quantify the effect of some of the properties of training sets on the accuracy of
26 genomic selection in crops we performed an extensive field-based winter wheat trial. In
27 summary, this trial involved the construction of 44 $F_{2:4}$ bi- and triparental populations, from
28 which 2992 lines were grown on four field locations and yield was measured. For each line,
29 genotype data were generated for 25,000 segregating single nucleotide polymorphism
30 markers. The overall heritability of yield was estimated to 0.65, and estimates within
31 individual families ranged between 0.10 and 0.85. Within cross genomic prediction accuracies
32 of yield BLUEs were 0.125 – 0.127 using two different cross-validation approaches, and
33 generally increased with training set size. Using related crosses in training and validation sets
34 generally resulted in higher prediction accuracies than using unrelated crosses. The results of
35 this study emphasize the importance of the training set design in relation to the genetic
36 material to which the resulting prediction model is to be applied.

37 **Keywords:** genomic selection, wheat, population design

38

39 Introduction

40 Genomic selection in plant breeding offers several routes for increasing the genetic gain
41 or efficiency of plant breeding programs (e.g., Bernardo and Yu, 2007; Hickey et al., 2014;
42 Gaynor et al., 2017). Genomic selection based strategies can achieve this by reducing breeding
43 cycle time, increasing selection accuracy and increasing selection intensity; three of the four
44 factors in the breeder's equation. Genomic prediction can reduce breeding cycle time because
45 individuals can be selected and crossed without being phenotyped. It can increase the selection
46 accuracy because genomic data enables more powerful statistical models and experimental
47 designs using more observations than can be phenotyped in a single trial round. By reducing
48 the cost of evaluating individuals via reducing the numbers phenotyped and/or reducing their
49 replication, application of genomic selection can increase selection intensity. A final advantage
50 is that the prediction models may be cumulatively updated with data of trials from previous
51 years and become more accurate, enabling individuals to be "evaluated" across a broader
52 range of environments and years.

53 In livestock there is empirical evidence of increased rates of genetic gain from the use of
54 genomic selection to target different aspects of the breeder's equation. For example the first
55 seven years of genomic selection in US dairy cattle has delivered ~50 - 100% increases in rates
56 of genetic gain (García-Ruiz et al., 2016). Much of this gain has emanated from a reduction in
57 generation interval. In commercial pig breeding, genomic selection has driven a 35% increase
58 in rate of genetic gain in the breeding program that supplies the genetics in 25% of the
59 intensively raised pigs globally. This gain came from increased accuracy of selection and a
60 better alignment of selection accuracy with the breeding goal (W. Herring, personal
61 communication).

62 Genomic selection uses genotype data to calculate the realised relationship between
63 individuals, and in a standardized statistical framework uses data from phenotyped relatives
64 to estimate genetic values of the selection candidates. The usefulness of genomic selection to
65 a breeder is a function of its accuracy. This is affected by the relatedness between the

66 phenotyped individuals in the training set and the individuals that are to be predicted (Habier
67 et al., 2007, 2010; Meuwissen, 2009; Clark et al., 2012; Hickey et al., 2014; Liu et al., 2016),
68 which may or may not be phenotyped themselves. In addition to the level of relatedness, the
69 sample size of the phenotyped individuals is an important factor in determining accuracy
70 (Zhang et al., 2017).

71 In summary, small numbers of close relatives and very large numbers of distant relatives
72 enable accurate predictions. Small or modest numbers of distant relatives do not enable
73 accurate predictions, as they share only a small proportion of genome with the selection
74 candidates, and thus provide less reliable predictions (de los Campos et al., 2013). Finally, the
75 training set should also comprise a diverse set of individuals to produce reliable predictions
76 (Calus, 2010; Pszczola et al., 2012; Pszczola and Calus, 2015), as supported by recent research
77 in both cattle (Jenko et al., 2017) and simulated barley (Neyhart et al., 2017).

78 The objective of this study was to explore the effect of level of relatedness between
79 training set and validation set on genomic prediction accuracy using data from a large set of
80 field experiments. To do this, 44 bi-parental or three-way crosses were obtained from four
81 commercial wheat breeders in the United Kingdom, as described for the GplusE Project
82 (Mackay et al., 2015). The crosses had different degrees of relatedness among each other and
83 there were many shared parents. 68 $F_{2:4}$ lines from each cross were genotyped and phenotyped
84 for yield. As this data set is of substantial size, it enabled genomic predictions while masking
85 specific fractions to assess the impact on genomic selection accuracy of training sets: (i) of
86 different sizes; and (ii) that comprise close or distant relatives, or combinations thereof.

87

88 Materials and Methods

89 Germplasm

90 Thirty-nine bi-parental and 5 triparental populations were used to develop 2992 $F_{2:4}$ lines (68
91 per cross). The parents of these populations were elite breeders' germplasm consisting of both
92 hard and soft winter wheat cultivars adapted to the United Kingdom. A total of 27 parents
93 were used, of which 5 parents were used in 6 or more crosses, 6 parents were used in 3 or 4
94 crosses, and 1 parent was used in 2 crosses. The remaining 15 parents were only used in a
95 single cross.

96 Genotypes

97 The $F_{2:4}$ lines were genotyped using the Wheat Breeders' 35K Axiom array (Allen et al., 2016).
98 The DNA for genotyping was obtained by bulking leaves from approximately 6 F_4 plants per
99 $F_{2:4}$ line. Genotype calling was performed using the Axiom Analysis Suite 2.0 with a modified
100 version of the "best practices" workflow. Quality control threshold was reduced to 95 (97
101 normally), plate pass percent was changed to 90 (95 normally), and average call rate was
102 changed to 97 (98.5 normally). After quality control and genotype calling, a total of 35,143
103 markers were brought forward with 24,498 segregating in the 44 crosses.

104 Phenotypes

105 The $F_{2:4}$ lines and agronomic checks were evaluated in 2 by 4 meter harvested plots at 2
106 locations (Cambridge, UK and Duxford, UK) in the 2015-16 growing season, and 2 locations
107 (Hinxton, UK, and Duxford, UK) in the 2016-17 growing season. All locations were managed
108 for optimal yield by following best agronomic practice. All $F_{2:4}$ lines were evaluated in 4 plots.
109 Seed for eleven of the populations was unavailable in the 2015-16 growing season. To
110 accommodate these populations and keep the number of plots per line constant, an allocation
111 of $F_{2:4}$ lines was devised that was highly unbalanced across both years and locations as
112 described below.

113 In the 2015-16 growing season, 33 of the 44 populations were planted at two locations
114 (Table 1). The experimental design for both locations was a modified α -lattice design
115 (Patterson and Williams, 1976). The design consisted of a traditional, replicated α -lattice
116 design with un-replicated lines added to the sub-blocks. The replicated portion of the alpha-
117 lattice design was composed of the agronomic checks and half of the lines (34) from 22 of the
118 $F_{2:4}$ populations. These lines were planted in 2 blocks split into 151 sub-blocks each containing
119 5 lines. The remaining $F_{2:4}$ lines were randomly allocated to sub-blocks, bringing the total
120 number of lines per sub-block to either 9 or 10. Half of the $F_{2:4}$ lines used for the replicated
121 portion of the design differed between locations. Thus lines from 22 of the $F_{2:4}$ populations
122 were evaluated in 3 plots split across both locations and the lines from the remaining
123 populations were evaluated in 2 plots split across locations.

124 All 44 populations were planted in the 2016-17 growing season at two locations (Table
125 1); the experimental design was similar as in the previous season. The replicated portion of the
126 α -lattice design was composed of the agronomic checks and the $F_{2:4}$ lines from the 11
127 populations not planted in the 2015-16 growing season. These lines were planted in 2 blocks
128 split into 156 sub-blocks each containing 5 lines. Additional $F_{2:4}$ lines from the other
129 populations were randomly allocated to sub-blocks, bring the total number of lines per sub-
130 block to 10.

131

132 Yield Trial Analysis

133 Yield phenotypes were spatially adjusted for each trial separately. An AR1 x AR1 model
134 (Gilmour et al., 1997) was used to adjust spatial variation across both columns and rows as
135 implemented in ASREML 3.0.22 (Gilmour et al., 2009). A summary of line means after
136 adjusting for spatial effects is shown in Table 2.

137 Best linear unbiased estimates (BLUEs) for each line were estimated collectively across all
138 trials by fitting the following model:

139
$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{u} + \mathbf{e}, \quad (2)$$

140 where \mathbf{y} was the response vector of spatially adjusted yield values, \mathbf{b} site-specific means with
141 design matrix \mathbf{X} , \mathbf{u} line BLUEs to estimate, and \mathbf{e} the model residual.

142 Genomic prediction

143 This study used the genomic best linear unbiased prediction (GBLUP) model to estimate
144 heritabilities and predict line effects. The GBLUP model used was:

145
$$\mathbf{y} = \mu + \mathbf{g} + \mathbf{e}, \quad (1)$$

146 where \mathbf{y} was the response vector of yield BLUEs, μ the model intercept, \mathbf{g} the vector of genetic
147 values of genotyped $F_{2:4}$ and \mathbf{e} the model residual. We assumed that $\mathbf{g} \sim N(0, \mathbf{G}\sigma_g^2)$ with
148 genomic relationship matrix calculated as $\mathbf{G} = \mathbf{W}\mathbf{W}' / 2\sum p_i(1 - p_i)$ (VanRaden, 2008) from the
149 centred genotype matrix \mathbf{W} and allele frequencies p_i estimated in the dataset. Further, we
150 assumed that $\mathbf{e} \sim N(0, \mathbf{I}\sigma_e^2)$, which was assumed uncorrelated to \mathbf{g} .

151 The Average-Information Restricted Maximum Likelihood (AI-REML) algorithm
152 (Madsen et al., 1994; Johnson and Thompson, 1995), as implemented in DMU v. 5.1 (Madsen
153 and Jensen, 2000), was used to fit the GBLUP model to a subset of the data (training set) and
154 predict line effects ($\hat{\mathbf{g}}$) in the validation set. We defined convergence of the AI-REML algorithm
155 based on the change of variance components, $|\theta^{(t+1)} - \theta^{(t)}| < 10^{-5}$, where $\theta^{(t)}$ is the vector of
156 normalised variance components estimated at step t (Jensen et al., 1997).

157 The heritability was calculated from the trial yield data per plot as $H^2 = \frac{\sigma_g^2}{\sigma_g^2 + \frac{v}{n}}$ in which n
158 is the number of locations in which the genotype was observed and v is the residual variance
159 (Piepho and Mohring, 2007).

160 Prediction accuracies

161 We applied several cross-validation strategies for investigating prediction accuracies of
162 genomic selection with varying training set size and grouping of training sets and validation
163 sets, as described in detail in the following sections. In all strategies, the GBLUP model was

164 used as described above. The prediction accuracies were calculated as the Pearson correlation
165 (ρ) between the yield BLUEs and its prediction from the GBLUP model.

166 **Cross-validation prediction accuracy**

167 In the first approach, we used 10-fold cross-validation and leave-one-cross-out cross-
168 validation (effectively 44-fold cross-validation; refer to Figure 1). Populations were randomly
169 assigned to either training or validation set, without considering that some crosses are more
170 closely related due to sharing a parent or other ancestors. The validation sets were entire
171 populations, which means that line means of a population was confined entirely to either
172 training set or validation set. Prediction accuracies were summarised on a per cross basis and
173 encapsulate the within cross genomic prediction accuracy (sometimes referred to as the within
174 family accuracy or the accuracy of predicting the Mendelian sampling term). For the 10-fold
175 cross-validation, 10 replicates were performed where the 10 folds were re-sampled.

176 To evaluate the effect of training set size, the two cross-validation methods described
177 above were repeated using a subset of the total training set. For the 10-fold cross-validation,
178 10%, 20%, ..., 80%, 90% of records were randomly removed from the training set, before
179 estimating variance components and predicting line means of the validation set. For each
180 replicate and the proportion of training set masked, 10 repetitions were performed and the
181 resulting prediction accuracies encapsulate the joint across and within cross genomic
182 prediction accuracy. For the leave-one-cross-out cross-validation, 1 – 10, 15, 20, 30, 40 crosses
183 were randomly sampled to be used as training set. For each number of crosses sampled as
184 training sets, 10, 20, ..., 60, 65 records from each cross was sampled. Again, 10 repetitions
185 were performed. We emphasise that the validation sets were always entire populations (from
186 3-4 crosses in 10-fold cross-validations, from single cross in leave-one-cross-out) and no
187 records of the validated populations were included in the training set.

188 **Prediction accuracy with related or unrelated crosses**

189 In the second approach, we evaluated the prediction accuracies under different levels of
190 relatedness between validation and training sets. The 6 crosses of the 4 most frequently used

191 parents were targeted as validation crosses and tested separately. In summary, the training
192 sets consisted of varying proportions of sister-lines and half-sibs from offspring of either one
193 or both parents or unrelated crosses. Specifically, for each validation cross, training sets were
194 designed to consist of either one or several crosses of one parent, an equal number of crosses
195 from each parent, nominally unrelated crosses, or equal number of related and unrelated
196 crosses. To reduce computation time, for each training set of crosses, 5 combinations were
197 sampled from the large number of possible combinations. For each training set, the validation
198 cross contributed with 0, 1, 2, or 3 quarters of its lines. The prediction accuracies were
199 evaluated for the fourth quarter of lines that were not used in the training set. For each
200 combination of training set, 10 replicates were performed as well as cycling through all four
201 quarters of the validation cross as training set.

202 Results

203 44 bi- and tri-parental crosses from 27 parents were analysed for yield with a GBLUP model
204 (1), using BLUEs from 4 trials (2 trials in 2016, and 2 trials in 2017).

205 Trait heritability

206 The overall heritability of yield for all populations over all four trial locations was estimated at
207 0.65. Heritabilities estimated on single crosses were highly variable, ranging from as low as
208 0.1 to as high as 0.85 (Figure 2).

209 Cross-validation prediction accuracy

210 Within cross prediction accuracies were 0.125 – 0.127 using two different cross-validation
211 approaches (Table 3). In these two approaches, all lines of the crosses used for validation were
212 absent from the training set. Using a 10-fold cross-validation approach where individual lines,
213 not all lines of a cross, were selected for validation sets, the prediction accuracy was slightly
214 higher (0.142) when calculated on a per-cross basis ('10-fold, random', Table). Lastly, the
215 prediction accuracy was higher when calculated across all crosses in the validation set, due to
216 capturing variation within and between crosses (0.289 and 0.543, Table 3).

217 The prediction accuracy was found to increase with training set size. Figure 3 displays
218 the average prediction accuracy across all crosses with 10th and 90th percentile range shown as
219 the greyed area. The prediction accuracy varied greatly between the crosses (Supplemental
220 figure 1) with some accuracies as high as 0.45 (cross 7), as low as -0.20 (cross 30). For 31
221 crosses out of 44, significant positive prediction accuracies were found (Wald's test, $p < 0.05$).
222 Crosses with higher phenotypic variance generally yielded higher predictions; in
223 Supplemental figure 1, prediction accuracy plots for individual crosses are sorted with
224 decreasing phenotypic variance. Finally, the two cross-validation approaches generally
225 produced similar results (Supplemental figure 1), but when the training sets were small, the
226 accuracy of predictions from leave-one-cross-out were less stable than from 10-fold cross

227 validation. The leave-one-cross-out sampled entire crosses in contrast to the 10-fold cross-
228 validation, where lines across all crosses except the validated cross were sampled.

229 The prediction accuracy increased with an increasing number of crosses in training set
230 or increasing number of lines per cross in training set. Figure displays the average prediction
231 accuracy when sampling a number of lines from a number of crosses (x-axis). Adding an
232 additional 10 or 15 lines to a training set of 50 lines per cross generally led to a low increase in
233 prediction accuracy as compared to adding them to training sets of ≤ 40 lines per cross,
234 irrespective of the number of crosses included in the training set.

235 Prediction accuracies with related or unrelated crosses

236 Using related crosses as a training set generally resulted in higher prediction accuracies
237 compared to using unrelated crosses. This is shown in Figure 5, where the green lines (related
238 training sets) are above the purple lines (unrelated training sets). Using both related and
239 unrelated crosses in equal proportions (blue lines, Figure 5) led generally to similar
240 correlations to those for related crosses. At approximately 700 to 800 lines in the training set,
241 the prediction accuracy using both related and unrelated crosses plateaued; this was where
242 additional crosses in the training set were unrelated to the validation cross. The level of
243 prediction accuracy of the training set comprising both related and unrelated crosses (lower
244 blue line, Figure) was higher than that in Figure because results in Figure are averages over
245 just 6 crosses rather than over all crosses as in Figure 3.

246 Using only 1, 2, or 3 quarters of the validation cross as training set (grey, horizontal lines,
247 Figure 5) generally led to prediction accuracies that were higher than using a few unrelated or
248 related crosses as the training set. Adding three quarters of the validation cross to the training
249 sets of other crosses generally increased the prediction accuracy, as shown with the upper thick
250 lines in Figure . The gradual increase in prediction accuracy when adding 1, 2, or 3 quarters of
251 the validation cross to the training set is shown in the inserted plot in Figure 5.

252 Discussion

253 In this study, we have demonstrated the impact of training set size and relatedness on
254 genomic prediction in wheat, using $F_{2:4}$ lines from 44 bi- and tri-parental crosses. The results
255 were consistent with expectations from existing literature (as discussed in the next sections).
256 Specifically, we found that increasing the size of the genomic prediction training set increased
257 accuracy. We also found that training sets composed of lines more closely related to the
258 validation set produce higher prediction accuracies than equivalently sized training sets of
259 more distantly related lines.

260 It is important for genomic prediction of a complex trait that it displays a reasonable
261 heritability. Our estimate of broad sense heritability for yield (0.65) is well within range of
262 similar studies in wheat (Poland et al., 2012; Combs and Bernardo, 2013; Michel et al., 2016;
263 Schopp et al., 2017; Norman et al., 2017). We note that the heritability values within individual
264 families (Figure) cover the whole range of heritability for this trait reported in the literature.

265 The various strategies of data subset masking applied in this study has enabled us to
266 demonstrate both training set size and relatedness as parameters that influence successful
267 genomic prediction. Generally, increasing the training set size increased the prediction
268 accuracy, as expected from existing theory (Daetwyler et al, 2008, Goddard, 2009, Hickey et
269 al., 2014) and field reports (Liu et al., 2016; Zhang et al., 2017). However, we can add three
270 observations that put some nuance to this general conclusion. First (1), with a fixed training
271 set size, it is better to increase the number of populations (crosses) rather than number of lines
272 per population (cross). Second (2), the prediction accuracy plateaus when adding additional
273 crosses that are unrelated to the predicted cross (Figure). Third (3), prediction accuracies vary
274 greatly between individual crosses and this could not be explained by neither the crosses'
275 phenotypic variance nor heritability.

276 For item 1), we showed that, for example, using 10 crosses with 40 lines per cross gave
277 prediction accuracy of ≈ 0.06 , while 40 crosses with 10 lines per cross gave prediction accuracy
278 of ≈ 0.075 (Figure). We assume that in both strategies different processes increase the

279 accuracy with the addition of extra lines: In the first case, entire crosses were masked
280 simulating the future prediction of an unphenotyped cross. In comparison, increasing the
281 number of lines instead of number of crosses (while constraining the training set size) did not
282 necessarily improve the prediction accuracy. The lines capture the crosses' variance, and there
283 will be a limit to how much more variance that additional lines will capture, hence no
284 additional gain. The exception to this was adding fractions of the validation cross' lines to the
285 training set (Figure).

286 For item 2), we saw in Figure that using training sets comprised of exclusively unrelated
287 crosses resulted in lower prediction accuracies than training sets that included related crosses.
288 Using training sets comprised of either exclusively related crosses or related and unrelated
289 crosses (half-and-half) both resulted in approximately the same prediction accuracy. The
290 comparison between these three sets stops at about 800 lines in the training set, because
291 beyond this point, additional crosses were no longer distinctively related or unrelated.
292 Therefore, after this point the slope of increase in prediction accuracy is less steep, as the
293 crosses added to the training set are less related.

294 For item 3), there was no observable connection between how well the cross could be
295 predicted and the cross' heritability or the observed phenotypic variance. Likewise, these
296 values did not correspond to how well the data from the cross could be used to predict breeding
297 values in other crosses.

298

299 One of the major practical implications of this study is that increased prediction
300 accuracies can be obtained by balancing the training set for genomic selection with phenotypic
301 and genomic data of multiple related crosses, which could be taken into account in advance
302 when designing the training population, as previously proposed by Rincent et al., 2012. For
303 existing data sets, a strategy may be applied of supplementing these with phenotypic data from
304 previous trials (provided genotype-by-environment interaction is limited or can be accounted
305 for by use of trait data for control lines). Although such data might be present within the

306 context of a rolling breeding program, obtaining genomic data presents a bottleneck as this
307 requires genotyping of (old) biological material that might not be readily available, and will
308 require investment in at least low-density genotyping. In case high density genotype data sets
309 are available for the parental lines, high density genotype information for their offspring
310 populations can subsequently be obtained by imputation, as reported by Hickey et al. (2015),
311 Gorjanc et al. (2017) and others.

312 **Conclusions**

313 Genomic predictions of yield across 44 populations resulted in modest correlations between
314 observed and predicted values. The correlations did increase with training set size, but by
315 selecting training sets that comprised related crosses improved the correlation more than
316 increasing training set size. The results also showed that if the training set size is fixed, using
317 few lines from more crosses, rather than many lines from few crosses, resulted in higher
318 correlations.

319

320 **Authors' contributions**

321 Wheat crosses were made by JL, EB, CB, PJ, SB, EF, BP, SS, CH; wheat yield trials
322 were conducted by RJ, PH, EO and IJM; ARB co-ordinated genotyping; SME, RCG
323 and GG performed data analysis; SME, JBB, RCG and JMH wrote the manuscript;
324 JMH and IJM conceived the study, designed the experiment and led the project.

325

326 **Acknowledgements**

327 The authors acknowledge the financial support from BBSRC project “GplusE: Genomic
328 selection and Environment modelling for next generation wheat breeding” (grants
329 BB/L022141/1 and BB/L020467/1) and the Medical Research Council (MRC) grant
330 MR/M000370/1. The crosses, seed and genomic DNA for this study were contributed by KWS
331 UK, RAGT Seeds Ltd., Elsoms Wheat Ltd and Limagrain UK.

332

333

334 **Competing interests**

335 The authors declare that they have no competing interest.

336 References

- 337 Allen, A.M., M.O. Winfield, A.J. BurrIDGE, R.C. Downie, H.R. Benbow, G.L.A. Barker, P.A.
338 Wilkinson, J. Coghill, C. Waterfall, A. Davassi, G. Scopes, A. Pirani, T. Webster, F.
339 Brew, C. Bloor, S. Griffiths, A.R. Bentley, M. Alda, P. Jack, A.L. Phillips, and K.J.
340 Edwards. 2016. Characterization of a Wheat Breeders' Array suitable for high-
341 throughput SNP genotyping of global accessions of hexaploid bread wheat (*Triticum*
342 *aestivum*). *Plant Biotechnol. J.* Available at <http://doi.wiley.com/10.1111/pbi.12635>.
- 343 Bernardo, R., and J. Yu. 2007. Prospects for Genomewide Selection for Quantitative Traits in
344 Maize. *Crop Sci.* 47(3): 1082.
- 345 Calus, M.P.L. 2010. Genomic breeding value prediction: methods and procedures. *animal*
346 4(02): 157–164.
- 347 de los Campos, G., A.I. Vazquez, R. Fernando, Y.C. Klimentidis, and D. Sorensen. 2013.
348 Prediction of Complex Human Traits Using the Genomic Best Linear Unbiased
349 Predictor. *PLoS Genet* 9(7): e1003608.
- 350 Clark, S.A., J.M. Hickey, H.D. Daetwyler, and J.H. van der Werf. 2012. The importance of
351 information on relatives for the prediction of genomic breeding values and the
352 implications for the makeup of reference data sets in livestock breeding schemes.
353 *Genet. Sel. Evol.* 44(1): 4.
- 354 Combs, E., and R. Bernardo. 2013. Accuracy of Genomewide Selection for Different Traits with
355 Constant Population Size, Heritability, and Number of Markers. *Plant Genome* 6(1).
- 356 Daetwyler H.D., B. Villanueva, J.A. Woolliams (2008) Accuracy of Predicting the Genetic Risk
357 of Disease Using a Genome-Wide Approach. *PLoS ONE* 3(10): e3395
- 358 García-Ruiz, A., J.B. Cole, P.M. VanRaden, G.R. Wiggans, F.J. Ruiz-López, and C.P. Van
359 Tassell. 2016. Changes in genetic selection differentials and generation intervals in US
360 Holstein dairy cattle as a result of genomic selection. *Proc. Natl. Acad. Sci.* 113(28):
361 E3995–E4004.
- 362 Gaynor, R.C., G. Gorjanc, A.R. Bentley, E.S. Ober, P. Howell, R. Jackson, I.J. Mackay, and J.M.
363 Hickey. 2017. A two-part strategy for using genomic selection to develop inbred lines.
364 *Crop Sci.* 57: 1404-1420
- 365 Gilmour, A.R., B.R. Cullis, and A.P. Verbyla. 1997. Accounting for Natural and Extraneous
366 Variation in the Analysis of Field Experiments. *J. Agric. Biol. Environ. Stat.* 2(3): 269–
367 293.
- 368 Gilmour, A.R., B.J. Gogel, B.R. Cullis, and R. Thompson. 2009. ASReml User Guide Release
369 3.0. VSN International Ltd, Hemel Hempstead, UK.
- 370 Goddard, M. 2009. Genomic selection: prediction of accuracy and maximisation of long term
371 response. *Genetica* 136(2): 245-57
372
- 373 Gonen, S., R. Ros-Freixedes, M. Battagin, G. Gorjanc, and J.M. Hickey. 2017. A method for the
374 allocation of sequencing resources in genotyped livestock populations. *Genet. Sel.*
375 *Evol.* 49(1) Available at [http://gsejournal.biomedcentral.com/articles/10.1186/s12711-](http://gsejournal.biomedcentral.com/articles/10.1186/s12711-017-0322-5)
376 017-0322-5 (verified 22 May 2017).

- 377 Gorjanc, G., M. Battagin, J.-F. Dumasy, R. Antolin, R.C. Gaynor, and J.M. Hickey. 2017.
378 Prospects for Cost-Effective Genomic Selection via Accurate Within-Family
379 Imputation. *Crop Sci.* 57(1): 216.
- 380 Gorjanc, G., M.A. Cleveland, R.D. Houston, and J.M. Hickey. 2015. Potential of genotyping-
381 by-sequencing for genomic selection in livestock populations. *Genet. Sel. Evol.* 47(1):
382 12.
- 383 Habier, D., R.L. Fernando, and J.C.M. Dekkers. 2007. The impact of genetic relationship
384 information on genome-assisted breeding values. *Genetics* 177(4): 2389–2397.
- 385 Habier, D., J. Tetens, F.-R. Seefried, P. Lichtner, and G. Thaller. 2010. The impact of genetic
386 relationship information on genomic breeding values in German Holstein cattle.
387 *Genet. Sel. Evol.* 42(1): 5.
- 388 Hickey, J.M., S. Dreisigacker, J. Crossa, S. Hearne, R. Babu, B.M. Prasanna, M. Grondona, A.
389 Zambelli, V.S. Windhausen, K. Mathews, and G. Gorjanc. 2014. Evaluation of genomic
390 selection training population designs and genotyping strategies in plant breeding
391 programs using simulation. *Crop Sci.* 54: 1476–1488.
- 392 Hickey JM, G. Gorjanc, T.K. Varshney and C. Nettelblad (2015) Imputation of Single
393 Nucleotide Polymorphism Genotypes in Biparental, Backcross, and Topcross
394 Populations with a Hidden Markov Model *Crop Science* 55: 1934-1946
395
- 396 Jenko, J., G.R. Wiggans, T.A. Cooper, S. a. E. Eaglen, W.G. de L. Luff, M. Bichard, R. Pong-
397 Wong, and J.A. Woolliams. 2017. Cow genotyping strategies for genomic selection in a
398 small dairy cattle population. *J. Dairy Sci.* 100(1): 439–452.
- 399 Jensen, J., E.A. Mantysaari, P. Madsen, and R. Thompson. 1997. Residual Maximum
400 Likelihood Estimation of (Co) Variance Components in Multivariate Mixed Linear
401 Models using Average Information. *J. Indian Soc. Agric. Stat.* 49: 215–236.
- 402 Johnson, D.L., and R. Thompson. 1995. Restricted Maximum Likelihood Estimation of
403 Variance Components for Univariate Animal Models Using Sparse Matrix Techniques
404 and Average Information. *J. Dairy Sci.* 78(2): 449–456.
- 405 Liu, G., Y. Zhao, M. Gowda, C.F.H. Longin, J.C. Reif, and M.F. Mette. 2016. Predicting Hybrid
406 Performances for Quality Traits through Genomic-Assisted Approaches in Central
407 European Wheat (L Lukens, Ed.). *PLOS ONE* 11(7): e0158635.
- 408 Mackay, I., E. Ober, and J. Hickey. 2015. GplusE: beyond genomic selection. *Food Energy*
409 *Secur.* 4(1): 25–35.
- 410 Madsen, P., and J. Jensen. 2000. A User's Guide to DMU. A Package for Analysing
411 Multivariate Mixed Models. Version 6, release 5.1. : 32.
- 412 Madsen, P., J. Jensen, and R. Thompson. 1994. Estimation of (co)variance components by
413 REML in multivariate mixed linear models using average of observed and expected
414 information. p. 455–462. *In* 5th WCGALP. Guelph, Canada.
- 415 Meuwissen, T.H. 2009. Accuracy of breeding values of “unrelated” individuals predicted by
416 dense SNP genotyping. *Genet. Sel. Evol.* 41(1): 35.
- 417 Michel, S., C. Ametz, H. Gungor, D. Epure, H. Grausgruber, F. Löschenberger, and H.
418 Buerstmayr. 2016. Genomic selection across multiple breeding cycles in applied bread
419 wheat breeding. *Theor. Appl. Genet.* 129(6): 1179–1189.

- 420 Neyhart, J.L., T. Tiede, A.J. Lorenz, and K.P. Smith. 2017. Evaluating Methods of Updating
421 Training Data in Long-Term Genomewide Selection. *G3* 5(8): 1499–1510.
422 GenesGenomesGenetics 7(5): 1499–1510.
- 423 Norman, A., J. Taylor, E. Tanaka, P. Telfer, J. Edwards, J.-P. Martinant, and H. Kuchel. 2017.
424 Increased genomic prediction accuracy in wheat breeding using a large Australian
425 panel. *Theor. Appl. Genet.* 130(12): 2543–2555.
- 426 Patterson, H.D., and E.R. Williams. 1976. A new class of resolvable incomplete block designs.
427 *Biometrika* 63: 83–92.
- 428 Piepho, H. and J. Mohring. 2007. Computing heritability and selection response from
429 unbalanced plant breeding trials. *Genetics* 177: 1881–1888.
- 430 Poland, J.A., J. Endelman, J. Rutkoski, S. Wu, Y. Manes, S. Dreisigacker, J. Crossa, H.
431 Sánchez-Villeda, M. Sorrells, and J.-L. Jannink. 2012. Genomic Selection in Wheat
432 Breeding using Genotyping-by-Sequencing. *Plant Genome J.* 5(3): 103.
- 433 Pszczola, M., and M.P.L. Calus. 2015. Updating the reference population to achieve constant
434 genomic prediction reliability across generations. *animal*: 1–7.
- 435 Pszczola, M., T. Strabel, H.A. Mulder, and M.P.L. Calus. 2012. Reliability of direct genomic
436 values for animals with different relationships within and to the reference population.
437 *J. Dairy Sci.* 95(1): 389–400.
- 438 R. Rincent, D. Laloë, et al. 2012. Maximizing the reliability of genomic selection by
439 optimizing the calibration set of reference individuals: Comparison of methods in two
440 diverse groups of maize inbreds (*Zea mays* L.). *Genetics* 192: 715–728
441
- 442 Ros-Freixedes, R., S. Gonen, G. Gorjanc, and J.M. Hickey. 2017. A method for allocating low-
443 coverage sequencing resources by targeting haplotypes rather than individuals. *Genet.*
444 *Sel. Evol.* 49(1): 78.
- 445 Schopp, P., D. Müller, Y.C.J. Wientjes, and A.E. Melchinger. 2017. Genomic Prediction Within
446 and Across Biparental Families: Means and Variances of Prediction Accuracy and
447 Usefulness of Deterministic Equations. *G3* 5(8): 1499–1510.
448 g3.300076.2017.
- 449 VanRaden, P.M. 2008. Efficient methods to compute genomic predictions. *J. Dairy Sci.* 91(11):
450 4414–23.
- 451 Zhang, A., H. Wang, Y. Beyene, K. Semagn, Y. Liu, S. Cao, Z. Cui, Y. Ruan, J. Burgueño, F. San
452 Vicente, M. Olsen, B.M. Prasanna, J. Crossa, H. Yu, and X. Zhang. 2017. Effect of Trait
453 Heritability, Training Population Size and Marker Density on Genomic Prediction
454 Accuracy Estimation in 22 bi-parental Tropical Maize Populations. *Front. Plant Sci.*
455 8 Available at <http://journal.frontiersin.org/article/10.3389/fpls.2017.01916/full>
456 (verified 22 November 2017).

457

458

459

460 **Figures**

461 **Table 1: Trial design summary showing number of plots per tested line per location**

# Lines	2015/2016		2016/2017	
	Cambridge	Duxford	Duxford	Hinxton
367	2	1	1	0
381	2	1	0	1
381	1	2	1	0
367	1	2	0	1
748	1	1	1	1
748	0	0	2	2
Total plots	2992	2992	2992	2992

462

463 **Table 2: Summary of line means per location after adjusting for spatial effects.**

	No. lines	Avg. value	Coef. Variation	Correlation [†]
2016 Cambridge	2,247	8.58	6.1%	0.63
2016 Duxford	2,248	10.82	6.3%	0.81
2017 Hinxton	2,249	4.64	10.3%	0.71
2017 Duxford	2,235	8.24	6.6%	0.62

464 [†]: Correlation between moisture corrected yield values and spatially adjusted values.

465 **Table 3: Prediction accuracies using the largest training sets by cross-validation**
 466 **approach.**

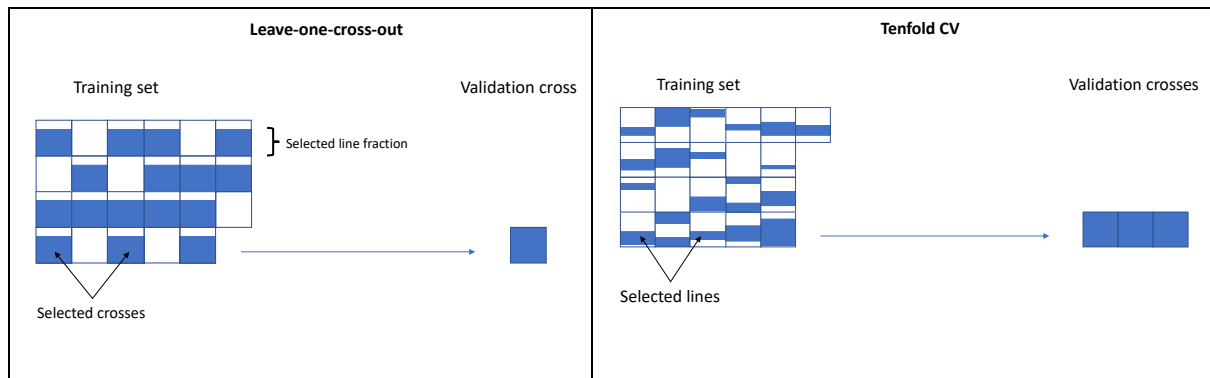
	Correlation metric	Training set size	Correlation [†]
Leave-one-cross-out	By cross	2,787	0.127 _{0.222}
10-fold, crosses	By cross	2,563	0.125 _{0.193}
10-fold, random [‡]	By cross	2,567	0.142 _{0.195}
10-fold, crosses	Across all [‡]	2,567	0.289 _{0.259}
10-fold, random [‡]	Across all [‡]	2,567	0.543 _{0.009}

467 [†]: Average across all replicates. Small font displays inter-quantile range for correlations.

468 [‡]: 10-fold cross-validation where validation and training sets were grouped by lines instead of crosses.

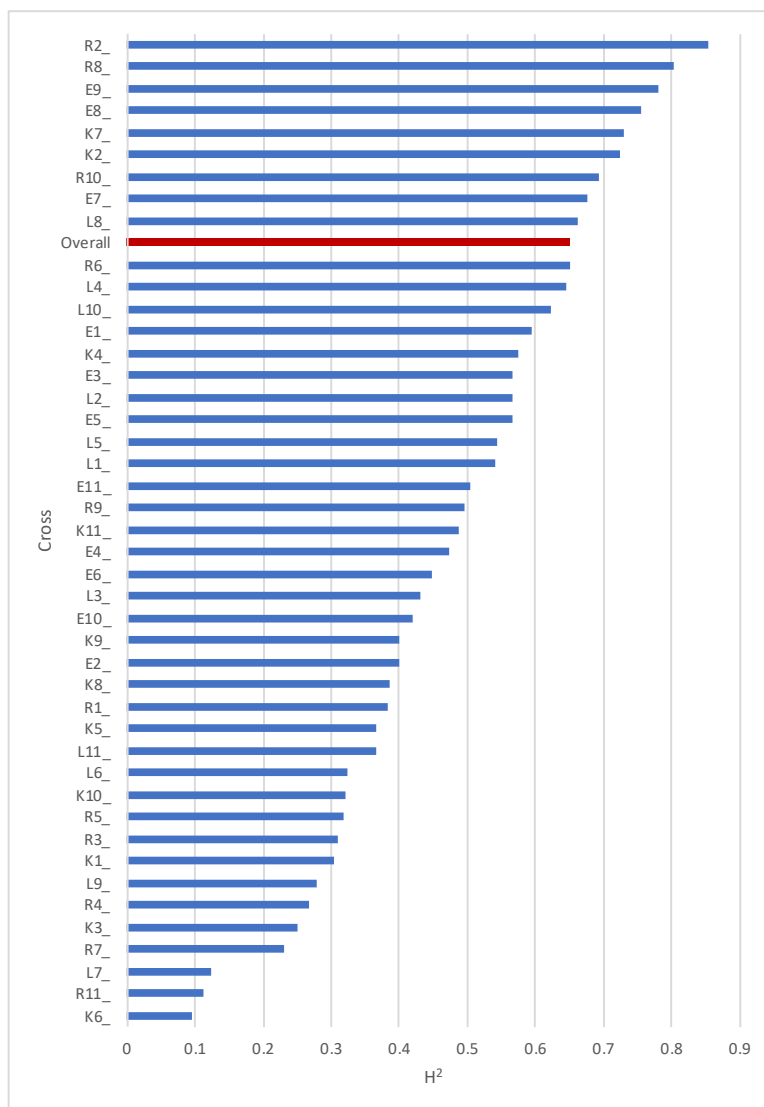
469 [‡]: Correlations were calculated across multiple crosses in validation set.

470



471 **Figure 1:** Resampling strategies applied to assess the impact of training set design. Leave-
472 one-cross out strategy (left) tests the impact of inclusion of the amount of crosses as well as
473 training set size, while the ten-fold cross validation (right) tests training set size only.

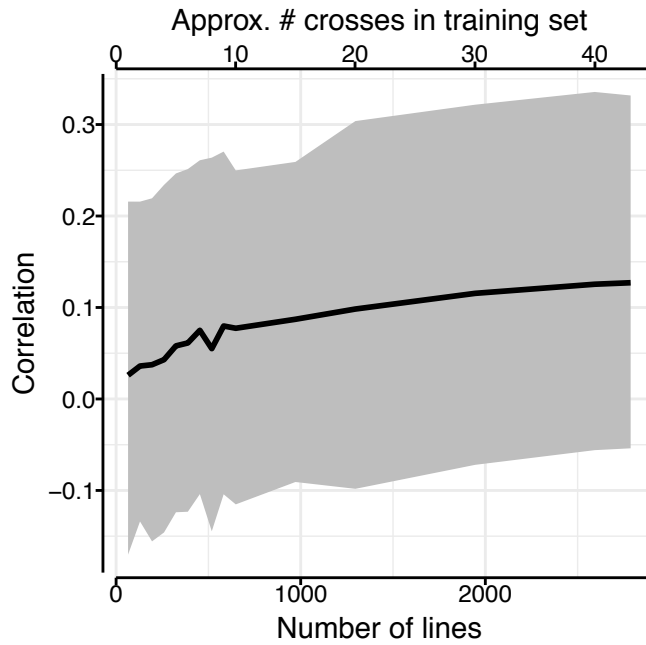
474



475

476 **Figure 2: Yield heritabilities when estimated per cross.** Crosses (blue bars) are ordered by

477 heritability value, overall heritability for this trait is shown in red.



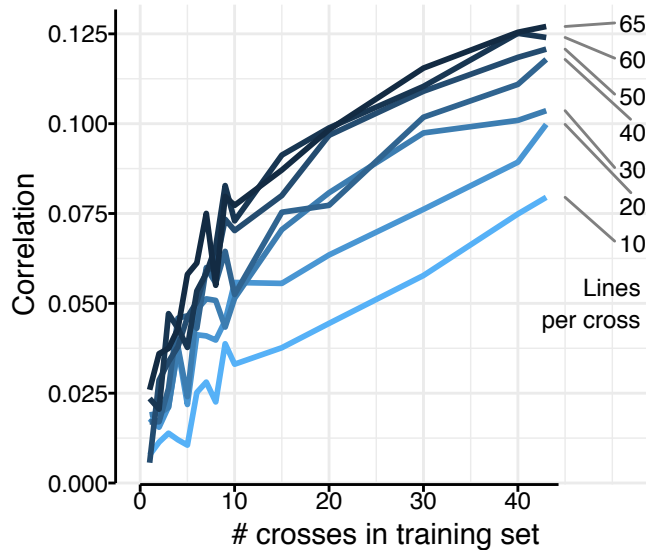
478

479 **Figure 3: Increasing training set size increased prediction accuracy (correlation).** Solid

480 line shows average of all leave-one-cross-out cross-validations with 10th and 90th percentile range shown

481 by greyed area.

482



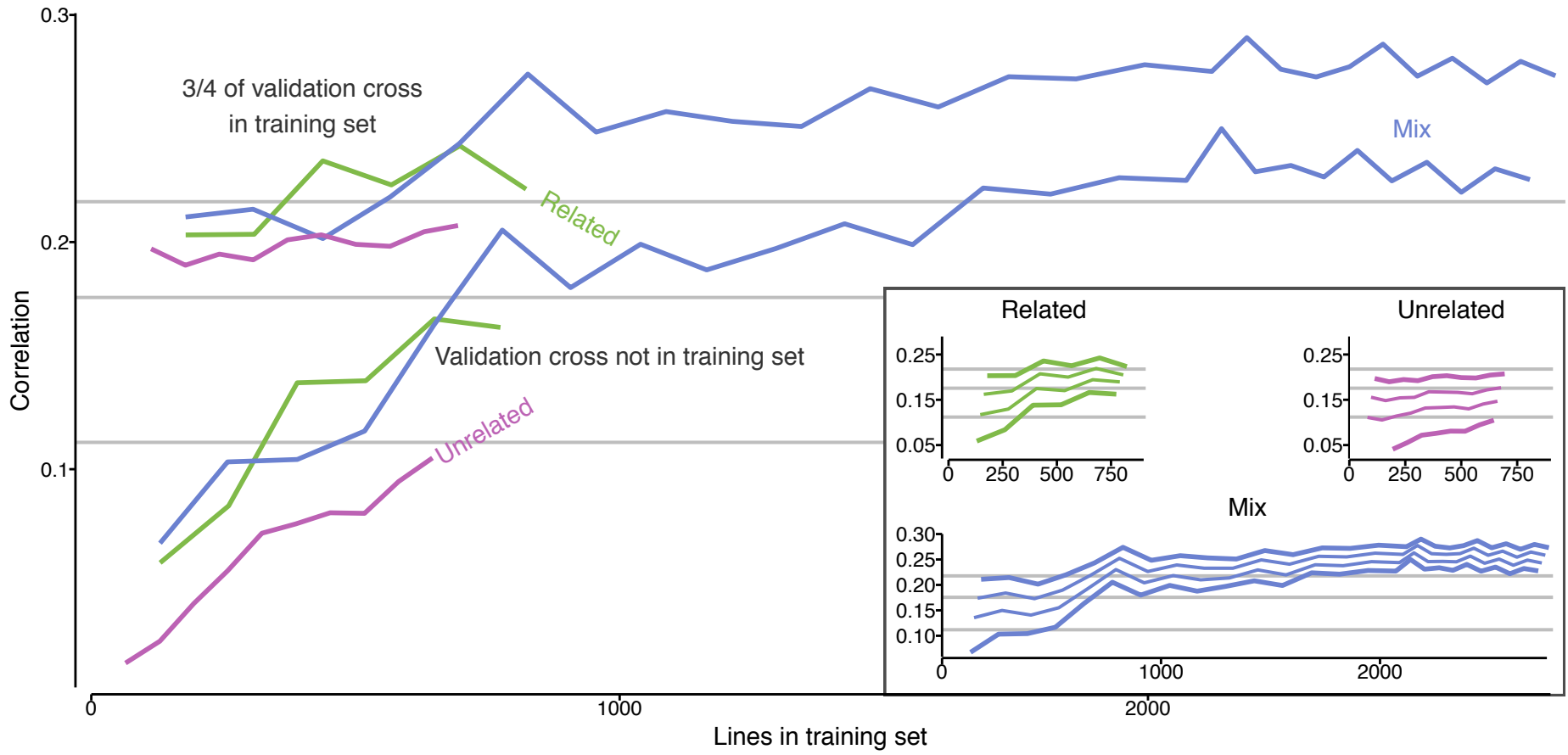
483

484 **Figure 4: Prediction accuracies increased with the increasing number of crosses or the**

485 **increasing number of lines per cross in training set.** Right-hand numbers show number of lines

486 per cross in training set.

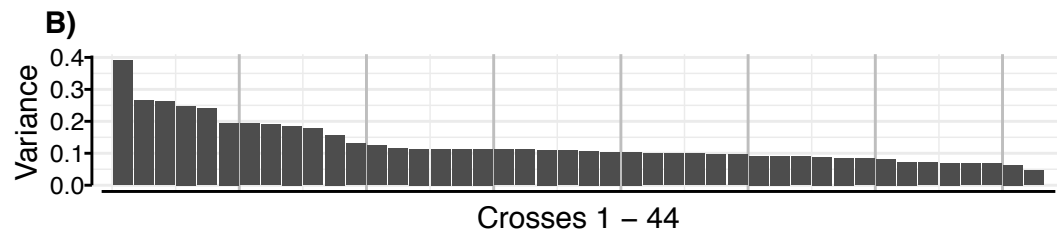
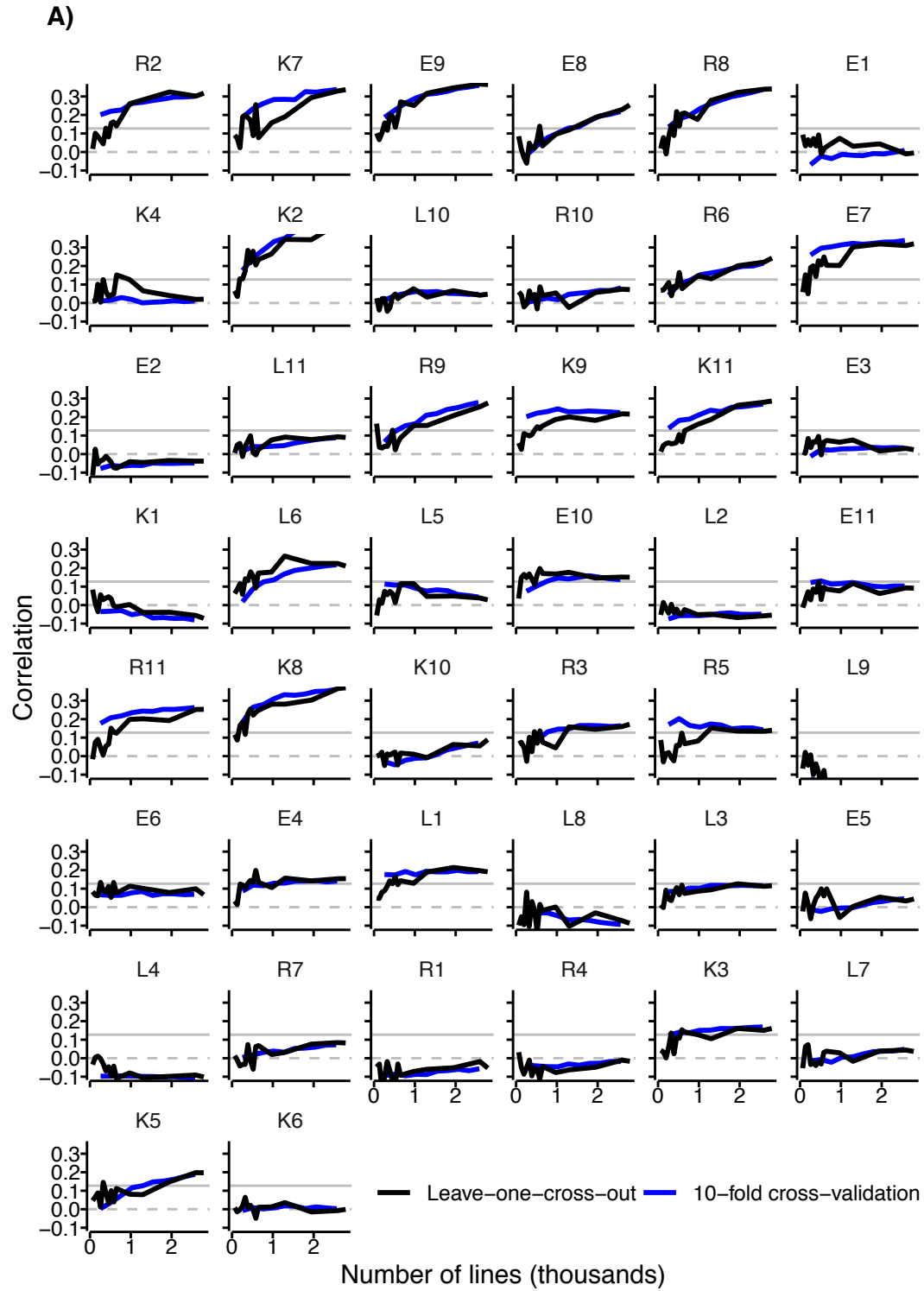
487



489 **Figure 5: Prediction accuracies increased when the validation cross was partly in training set or had its related crosses in training set.**

490 Results show average prediction accuracies for 6 validation crosses. Lines show prediction accuracies when training set comprised of related crosses (green solid
491 line), unrelated crosses (purple dashed line), or a mix of both (blue dotted line). Lower set of lines show prediction accuracies when validation crosses were not
492 included on the training set; upper set of lines show prediction accuracies when validation crosses were included in the training set with 3/4 of lines. Grey
493 horizontal lines show average prediction accuracy using *only* 1/4, 2/4, or 3/4 of validation cross as training set. Inserted figure shows the increase in accuracy
494 when adding 1/4, 2/4, and 3/4 of the validation group to the training set. The thick lines in the inserted figure denote the lines of the main figure.

495 Supplementary materials



497 **Supplemental figure 1: Per-cross correlation under two approaches (A), ordered by**
498 **decreasing variance of crosses' BLUEs (B).** Grey, horizontal lines are guides for zero correlation
499 (dashed) and overall average correlation of 0.127 (solid). Crosses in A) are ordered with decreasing
500 variance of their BLUEs, same order as in B).