![The University of Edinburgh logo]

# THE UNIVERSITY of EDINBURGH

## Edinburgh Research Explorer

# Genomic prediction with whole-genome sequence data in intensely selected pig lines

**OPEN ACCESS**

# Genomic prediction with whole-genome sequence data in intensely selected pig lines

Roger Ros-Freixedes[1,2§], Martin Johnsson[1,3], Andrew Whalen[1], Ching-Yi Chen[4], Bruno D Valente[4], William O Herring[4], Gregor Gorjanc[1], John M Hickey[1]

[1] The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of Edinburgh, Easter Bush, Midlothian, Scotland, UK

[2] Departament de Ciència Animal, Universitat de Lleida - Agrotecnio-CERCA Center, Lleida, Spain.

[3] Department of Animal Breeding and Genetics, Swedish University of Agricultural Sciences, Uppsala, Sweden.

[4] The Pig Improvement Company, Genus plc, Hendersonville, TN, USA.

[§]Corresponding author: RRF roger.ros@roslin.ed.ac.uk

# Abstract

## Background

16   Early simulations indicated that whole-genome sequence data (WGS) could improve

17   prediction accuracy and its persistence across generations and breeds. However,

18   results in real datasets have been ambiguous so far. Large data sets that capture most

19   of the genome diversity in a population must be assembled so that allele substitution

20   effects are estimated with higher accuracy. The objectives of this study were to use a

21   large pig dataset to assess the benefits of using WGS for genomic prediction

22   compared to using commercial marker arrays, to identify scenarios in which WGS

23   provides the largest advantage, and to identify potential pitfalls for its effective

24   implementation.

## Methods

25   We sequenced 6,931 individuals from seven commercial pig lines with different

26   numerical size. Genotypes of 32.8 million variants were imputed for 396,100

27   individuals (17,224 to 104,661 per line). We used BayesR to perform genomic

28   prediction for 8 real traits and 9 simulated traits with different genetic architectures.

29   Genomic predictions were performed using either data from a marker array or variants

30   preselected from WGS based on linkage disequilibrium, functional annotation, or

31   association tests. Both single and multi-line training sets were explored.

## Results

32   Using WGS improved prediction accuracy relative to the marker array, provided that

33   training sets were sufficiently large, especially for traits with high heritability and low

34   number of quantitative trait nucleotides. The performance of each set of predictor

35   variants was not robust across traits and lines. The most robust results were obtained

36 when preselected variants with statistically significant associations were added to the

37 marker array. Under this method, average improvements of prediction accuracy of 2.5

38 and 4.2 percentage points were observed in within-line and multi-line scenarios,

39 respectively, with training sets of around 80k individuals.

**Conclusions**

40 Our results evidenced the potential for WGS to improve genomic prediction accuracy

41 in intensely selected pig lines. Although the prediction accuracy improvements

42 achieved so far were modest at best, we would expect that more robust improvements

43 could be attained with a combination of larger training sets and optimised pipelines.

44

## Introduction

45       Whole-genome sequence data (WGS) has the potential to empower the

46    identification of causal variants that underlie quantitative traits or diseases [1–4],

47    increase the precision and scope of population genetic studies [5,6], and enhance

48    livestock breeding. Genomic prediction has been successfully implemented in the

49    main livestock species and it has increased the rate of genetic gain [7]. Genomic

50    prediction has provided many benefits such as greater accuracies of genetic

51    evaluations and the reduction of the generational interval in dairy cattle. However,

52    since its early implementations, genomic prediction is typically performed using

53    marker arrays that capture the effects of the (usually unknown) causal variants via

54    linkage disequilibrium. Alternatively, WGS are assumed to contain the causal variants

55    themselves. For this reason, it was hypothesized that such data could further improve

56    prediction accuracy and its persistence across generations and breeds. Early

57    simulations indicated that causal mutations from WGS could increase prediction

58    accuracy [8–13]. One simulation study indicated that the magnitude of prediction

59    accuracy improvement relative to dense marker arrays ranged from 2.5 to 3.7%, with

60    a persistence of over 10 generations [11]. Another one reached improvements of 30%

61    if causal variants with low minor allele frequency could be captured by the WGS [9].

62    However, benefits could be on the lower end of that range in standard livestock

63    populations due to small effective population sizes and long-term negative selection

64    [10].

65       During the last few years, there have been several attempts at improving the

66    accuracy of genomic prediction with the use of WGS in the main livestock species.

67    Results have been ambiguous so far. When predicting genomic breeding values within

68    breed or line, some studies found no relevant improvement of prediction accuracy for

69   WGS compared to marker arrays [14–18]. Other studies found small, and often

70   unstable, improvements (e.g., from 1 to 5% or no improvement depending on

71   prediction method [19–21], or trait-dependent results [21,22]). When predicting

72   genomic breeding values across populations, the identification of causal variants from

73   WGS can improve prediction accuracy [23–26], especially for small populations

74   where initial prediction accuracy was low or that were not included in the training

75   population [23,25–28].

76         One of the most successful strategies to exploit WGS consists in augmenting

77   available marker arrays with preselected variants from WGS based on their

78   association with the trait of interest [29–32]. In some cases, this strategy improved

79   prediction accuracy by up to 9% [31] and 11% [32]. However, it did not improve

80   prediction accuracies in other within-line scenarios [16]. Nevertheless, this shows

81   how identifying causal variants could enhance genomic prediction with WGS. Whole-

82   genome sequence data has already been applied in genome-wide association studies

83   (GWAS) to identify variants associated to a variety of traits in livestock [2,33–35],

84   including pigs [36,37]. However, the fine-mapping of causal variants remains

85   challenging due to the pervasive long-range linkage disequilibrium across extremely

86   dense variation.

87         High accuracy in estimating allele substitution effects and, ideally, the

88   identification of causal variants amongst millions of other variants are important for

89   the usefulness of WGS in research and breeding. This requires large data sets able to

90   capture most of the genome diversity in a population. Despite that low-cost

91   sequencing strategies have been developed, which typically involve sequencing a

92   subset of the individuals in a population at low coverage and then imputing WGS for

93   the remaining individuals [38–40], the cost of generating accurate WGS at this scale,

94    as well as the large computational requirements for the analyses of such datasets, have

95    limited the population sizes or number of populations tested in some of the previous

96    studies. This hinders the interpretation of results across studies, which are very

97    diverse in population structures, sequencing strategies and prediction methodologies

98    used. The largest studies on the use of WGS for genomic prediction to date have been

99    performed in cattle, for which large multi-breed reference panels are available from

100   the 1000 Bull Genomes Project [2,19,33]. This has enabled the imputation of WGS

101   for cattle populations. The lack of such available reference panels has been cited as an

102   important limiting factor for performing similar studies in other species, such as pigs

103   [36].

104         We have previously described our approach to impute WGS in large pedigreed

105   populations without the need for haplotype phased reference panels [41]. Following

106   that strategy, we generated WGS for 396,100 pigs from seven intensely selected lines

107   with diverse genetic backgrounds and numerical size. The objectives of this study

108   were to use this large pig dataset to assess the benefits of using WGS for genomic

109   prediction compared to using commercial marker arrays, to identify scenarios in

110   which WGS provides the largest advantage, and to identify potential pitfalls for its

111   effective implementation.

112

## Materials and Methods

### Populations and sequencing strategy

113         We performed whole-genome re-sequencing of 6,931 individuals from seven

114   commercial pig lines (Genus PIC, Hendersonville, TN) with a total coverage of

115   approximately 27,243x. Sequencing effort in each of the seven lines was proportional

116   to population size. Approximately 1.5% (0.9 to 2.1% in each line) of the pigs in each

117    line were sequenced. Most pigs were sequenced at low coverage, with target coverage

118    of 1 or 2x, but a subset of pigs were sequenced at higher coverage of 5, 15, or 30x.

119    Thus, the average individual coverage was 3.9x, but the median coverage was 1.5x.

120    The number of pigs sequenced and at which coverage for each line is summarized in

121    Table 1.

122        The sequenced pigs and their coverage were selected following a three-part

123    sequencing strategy developed to represent the haplotype diversity in each line. First

124    (1), sires and dams with the highest number of genotyped progeny were sequenced at

125    2x and 1x, respectively. Sires were sequenced at a greater coverage because they

126    contributed with more progeny than dams. Then (2), the individuals with the greatest

127    genetic footprint on the population (i.e., those that carry more of the most common

128    haplotypes) and their immediate ancestors were sequenced at a coverage between 1x

129    and 30x (AlphaSeqOpt part 1; [42]). The sequencing coverage was allocated with an

130    algorithm that maximises the expected phasing accuracy of the common haplotypes

131    from the cumulated family information. Finally (3), pigs that carried haplotypes with

132    low cumulated coverage (below 10x) were sequenced at 1x (AlphaSeqOpt part 2;

133    [43]). Sets (2) and (3) were based on haplotypes inferred from marker array genotypes

134    (GGP-Porcine HD BeadChip; GeneSeek, Lincoln, NE), which were phased and

135    imputed using AlphaPhase [44] and AlphaImpute [45].

136        Most sequenced pigs and their relatives were also genotyped either at low

137    density (15k markers) using the GGP-Porcine LD BeadChip (GeneSeek) or at high

138    density (80k markers) using the GGP-Porcine HD BeadChip (GeneSeek). Quality

139    control of the marker array data was based on the individuals genotyped at high

140    density. Markers with minor allele frequency below 0.01, call rate below 0.80, or that

141   failed the Hardy-Weinberg equilibrium test were removed. After quality control,

142   38,634 to 43,966 markers remained in each line.

143

**Sequencing and data processing**

144   Tissue samples were collected from ear punches or tail clippings. Genomic

145   DNA was extracted using Qiagen DNeasy 96 Blood & Tissue kits (Qiagen Ltd.,

146   Mississauga, ON, Canada). Paired-end library preparation was conducted using the

147   TruSeq DNA PCR-free protocol (Illumina, San Diego, CA). Libraries for

148   resequencing at low coverage (1 to 5x) were produced with an average insert size of

149   350 bp and sequenced on a HiSeq 4000 instrument (Illumina). Libraries for

150   resequencing at high coverage (15 or 30x) were produced with an average insert size

151   of 550 bp and sequenced on a HiSeq X instrument (Illumina). All libraries were

152   sequenced at Edinburgh Genomics (Edinburgh Genomics, University of Edinburgh,

153   Edinburgh, UK).

154   DNA sequence reads were pre-processed using Trimmomatic [46] to remove

155   adapter sequences from the reads. The reads were then aligned to the reference

156   genome *Sscrofa11.1* (GenBank accession: GCA_000003025.6) using the BWA-MEM

157   algorithm       [47].      Duplicates     were     marked      with      Picard

158   (http://broadinstitute.github.io/picard). Single nucleotide polymorphisms (SNPs) and

159   short insertions and deletions (indels) were identified with the variant caller GATK

160   HaplotypeCaller (GATK 3.8.0) [48,49] using default settings. Variant discovery with

161   GATK HaplotypeCaller was performed separately for each individual and then a joint

162   variant set for all the individuals in each population was obtained by extracting the

163   variant positions from all the individuals.

164        We extracted the read counts supporting each allele directly from the aligned

165        reads stored in the BAM files using a pile-up function to avoid biases towards the

166        reference allele introduced by GATK when applied on low-coverage WGS [50]. That

167        pipeline uses the tool pysam (version 0.13.0; https://github.com/pysam-

168        developers/pysam), which is a wrapper around htslib and the samtools package [51].

169        We extracted the read counts for all biallelic variant positions, after filtering variants

170        in potential repetitive regions (defined as variants that had mean depth values 3 times

171        greater than the average realized coverage) with VCFtools [52]. This amounted to a

172        total of 55.6 million SNP (19.6 to 31.1 million within each line) and 10.2 million

173        indels (4.1 to 5.6 million within each line). A more complete description of the

174        variation across the lines is provided in [53].

175

### Genotype imputation

176        Genotypes were jointly called, phased and imputed for a total of 483,353

177        pedigree-related individuals using the 'hybrid peeling' method implemented in

178        AlphaPeel [54,55]. This method used all the available marker array and WGS.

179        Imputation was performed separately for each line using complete multi-generational

180        pedigrees, which encompassed from 21,129 to 122,753 individuals each (Table 1).

181        We have previously published reports on the accuracy of imputation in the same

182        populations using this method [41]. The estimated average individual-wise dosage

183        correlation was 0.94 (median: 0.97). Individuals with low predicted imputation

184        accuracy were removed before further analyses. An individual was predicted to have

185        low imputation accuracy if itself or all of its grandparents were not genotyped with a

186        marker array or if it had a low degree of connectedness to the rest of the population.

187        These criteria were based on the analysis of simulated and real data on imputation

188    accuracy [41]. A total of 396,100 individuals remained, with each line comprising

189    between 17,224 and 104,661 individuals (Table 1). The expected average individual-

190    wise dosage correlation of the remaining individuals was 0.97 (median: 0.98)

191    according to our previous estimates. We also excluded from the analyses variants with

192    a minor allele frequency lower than 0.023, as their estimated variant-wise dosage

193    correlations was lower than 0.90 [41]. After imputation, 32.8 million variants (14.5 to

194    19.9 million within each line) remained for downstream analyses, out of which 9.9

195    million segregated across all seven lines.

196

## Traits

197        We analysed data of 8 traits that are commonly included in selection

198    objectives of pig breeding programmes: average daily gain (ADG, g), backfat

199    thickness (BFT, mm), loin depth (LD, mm), average daily feed intake (ADFI, kg),

200    feed conversion ratio (FCR), total number of piglets born (TNB), litter weight at

201    weaning (LWW, kg), and return to oestrus 7 days after weaning (RET, binary trait).

202    Most pigs with records were born during the 2008–2020 period. Breeding values were

203    estimated by line with a linear mixed model that included polygenic and non-genetic

204    (as relevant for each trait) effects. Deregressed breeding values (dEBV) were obtained

205    following the method by VanRaden and Wiggans [56]. Only individuals in which the

206    trait was directly measured were retained for further analyses. The number of records

207    for each trait used in the analyses of each line is detailed in Table 2.

208

### *Simulated traits*

210        To assist in the interpretation of results, we also created 9 simulated traits with

211    different numbers of quantitative trait nucleotides (QTN; 100, 1,000 or 10,000 QTN)

212     and heritability levels ($h^2$; 0.10, 0.25 or 0.50). Positions of the QTN were sampled

213     randomly amongst all variants called across all lines. Because QTN were sampled

214     from all variants, some QTN were fixed in some of the lines while segregating in

215     others. There were only negligible differences in the number of segregating QTN per

216     line (53 to 61, 531 to 583, or 5375 to 6058, respectively). Marker effects of the QTN

217     were sampled from a gamma distribution with shape=2 and scale=5. After a polygenic

218     term was calculated for each individual using these marker effects, residual terms

219     were sampled from a normal distribution with a variance parameter adjusted to

220     produce the desired heritability level. The number of records for the simulated traits is

221     detailed in Table 2. In these simulations, we used the imputed genotypes as real

222     genotypes and, therefore, implicitly cancelled any errors that might arise from the

223     processing of the sequencing reads and genotype imputation.

224

### Training and testing sets

225     We split the individuals in each population into training and testing sets. The

226     testing sets were defined as those individuals from full-sib families from the last

227     generation of the pedigree (i.e., individuals that did not have any progeny of their

228     own). Only families with a minimum of 5 full-sibs were considered. The training set

229     was defined as all those individuals that had a pedigree coefficient of relationship

230     lower than 0.5 with any individual of the testing set. This design was chosen to mimic

231     a realistic situation in which breeding companies evaluate the selection candidates

232     available in the selection nucleus at any given time.

233

### Genome-wide association study

234       To assess whether variants from the WGS could provide a finer mapping of

235       causal variants than marker array data, and to provide an association-based criterion

236       to preselect variants for the genomic prediction tests, we performed a GWAS for each

237       trait and line. This step included only the individuals in the training set. We fitted a

238       univariate linear mixed model that accounted for the genomic relationship matrix as:

239 $$\mathbf{y} = \mathbf{x}_i\beta_i + \mathbf{u} + \mathbf{e},$$

240       where $\mathbf{y}$ is the vector of dEBV, $\mathbf{x}_i$ is the vector of genotypes for the $i$th SNP coded as

241       0 and 2 if homozygous for either allele or 1 if heterozygous, $\beta_i$ is the additive effect

242       of the $i$th SNP on the trait, $\mathbf{u} \sim N(0, \sigma_u^2\mathbf{K})$ is the vector of polygenic effects with the

243       covariance matrix equal to the product of the polygenic additive variance $\sigma_u^2$ and a

244       genomic relationship matrix $\mathbf{K}$, and $\mathbf{e}$ is a vector of uncorrelated residuals. Due to

245       computational limitations, the genomic relationship matrix $\mathbf{K}$ was calculated using

246       only imputed SNP genotypes in the marker array regardless of whether the association

247       study involves the SNPs in the marker array or the variants in WGS. We used the

248       FastLMM software [57,58] to fit the model.

249       We used the same p-value threshold ($p<10^{-6}$) for both marker array and for

250       sequence associations, because while the WGS contains many more variants, they are

251       also expected to be in higher linkage disequilibrium. This threshold was based on

252       Bonferroni's multiple test correction assuming that the markers from the marker array

253       were independent. For the simulated traits, we defined genomic regions that contained

254       significant associations and assessed whether or not they contained a QTN. These

255       regions were defined by overlapping 500-kb segments centered on the significant

256       markers.

257

## Genomic prediction in within-line scenarios

258    To test whether variants from the WGS could provide greater prediction

259    accuracy than the marker array, we tested genomic prediction using variants from the

260    marker array, from the WGS, or combining them. The marker array data (referred to

261    as 'Chip') was set as the benchmark for prediction accuracy. It contained all ~40k

262    variants in the marker array. For the sequence-based predictors, we preselected sets of

263    variants because currently available methods for genomic prediction are not yet

264    capable of handling datasets as large as the complete WGS. We tested different

265    alternative strategies for preselecting the predictor variants:

266    • *LDTags*. Tag variants retained after pruning based on linkage disequilibrium.

267        Variants were removed so that no pairs of SNPs with $r^2>0.1$ remained in any 10-

268        Mb window (windows slid by 2,000 variants) using Plink 1.9 [59]. The number

269        of predictor variants preselected by this method was on average of 30k variants

270        (range: 5k to 80k).

271    • *Top40k*. Variants preselected based on GWAS analyses. To mimic the number of

272        variants in Chip, we preselected the variants with the lowest p-value (not

273        necessarily below the significance threshold) in each of consecutive non-

274        overlapping 55-kb windows along the genome. In addition, to test the impact of

275        variant density on prediction accuracy, we preselected 10k, 25k, 75k, or 100k

276        predictor variants following the same criterion.

277    • *ChipPlusSign*. Variants preselected based on GWAS analyses as in Top40k, but

278        only significant variants ($p \leq 10^{-6}$) were preselected and merged with those in

279        Chip. When a 55-kb window contained more than one significant variant, only

280        that with the lowest p-value was selected as a proxy, in order to reduce the

281   preselection of multiple significant SNPs tagging the same causal variant. On

282   average, 309 significant variants were identified per trait and line (range: 23 to

283   1083; Table 3). These significant variants were merged with those in Chip.

284  • *Functional.* Variants that were annotated as loss-of-function or missense

285   according to Ensembl Variant Effect Predictor (Ensembl VEP; version 97, July

286   2019) [60]. The most severe predicted consequence type for each variant was

287   retrieved. The number of predictor variants preselected by this method was on

288   average of 35k variants (range: 27k to 40k).

289  • *Rand40k.* The same number of predictor variants as in Chip, chosen randomly.

290   Genomic prediction was performed by fitting a univariate model with BayesR

291 [61,62], with a mixture of normal distributions as the prior for variant effects,

292 including one distribution that sets the variant effects to zero. The model was:

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

293 where $\mathbf{y}$ is the vector of dEBV, $\mathbf{1}$ is a vector of ones, $\mu$ is the general mean, $\mathbf{X}$ is a

294 matrix of genotypes, $\boldsymbol{\beta}$ is a vector of variant effects, and $\mathbf{e}$ is a vector of uncorrelated

295 residuals. The prior variance of the variant effects in $\boldsymbol{\beta}$ had four components with

296 variances $\sigma_1^2 = 0$, $\sigma_2^2 = 0.0001\sigma_g^2$, $\sigma_3^2 = 0.001\sigma_g^2$, or $\sigma_4^2 = 0.01\sigma_g^2$, where $\sigma_g^2$ is the

297 total genetic variance. We used a uniform and almost uninformative prior for the

298 mixture distribution. We used a publicly available implementation of BayesR

299 (https://github.com/syntheke/bayesR; accessed on 30 April 2021), with default

300 settings. Prediction accuracy was calculated in the testing set as the correlation

301 between the genomic estimated breeding value and the dEBV. Bias of the prediction

302 accuracy was calculated as the regression coefficient of the dEBV on the genomic

303 estimated breeding values.

304    It has been noted that using the same reference individuals for preselecting

305    variants through GWAS and for training the predictive equation can reduce prediction

306    accuracy and bias the predicted breeding values [16,63]. To account for that, we

307    reanalysed some of the scenarios after splitting the training set into two exclusive

308    subsets, one for GWAS to preselect the predictor variants and one for training the

309    predictive equation. The GWAS subset was defined by randomly selecting either 10%

310    or 50% of the individuals in the original training set. Those individuals were excluded

311    from the subset used for training the predictive equation afterwards.

312

**Genomic prediction in multi-line scenarios**

313    We considered multi-line scenarios in which the training set consisted of

314    merging the training sets that had been defined for each line. All analyses were

315    performed as for the within-line scenarios but with a line effect. In the multi-line

316    scenarios, all SNPs from the marker array that passed quality control and were

317    imputed for at least one line were included in the baseline (referred to as 'ML-Chip').

318    For ease of computation, the strategies for preselection of predictor variants from

319    WGS were applied only to the subset of 9.9 million variants that had been called and

320    imputed in all seven lines. Thus, we defined the predictor sets 'ML-Top40k' and

321    'ML-ChipPlusSign' by preselecting variants following the same criteria as in within-

322    line scenarios, but using a multi-line GWAS analyses with line effect instead. For

323    ML-ChipPlusSign, 60 to 7247 significant variants were identified per trait (Table 3)

324    and merged with those in ML-Chip. For comparison purposes, prediction accuracy

325    was calculated for the testing set of each individual line.

326

# Results

### Prediction accuracy within line

327    Whole-genome sequence data can improve prediction accuracy of marker

328    array data when there is a sufficiently large training set and if an appropriate set of

329    predictor variants is preselected. Figure 1 shows the prediction accuracy for the case

330    with the largest training set using different sets of predictor variants. In this case, all

331    tested sets of variants from the WGS, except for LDTags, yielded increases of

332    prediction accuracy that ranged from +2.0% to +9.2%. Using WGS also reduced bias

333    relative to Chip in some scenarios. However, the performance across predictor

334    variants set was not robust for the most part, and differed for each trait and line

335    (Additional File 1), often leading to no improvements of prediction accuracy or even

336    reduced prediction accuracy relative to Chip. One stable feature of the results was

337    LDTags showing a noticeable decrease in prediction accuracy in most traits and lines.

338    The size of the training set was one of the main factors that determined the

339    capacity of predictor variants from the WGS to improve the baseline prediction

340    accuracy of Chip. Figures 2 and 3 show the difference in prediction accuracy of

341    Top40k and ChipPlusSign with respect to the baseline of Chip against the number of

342    phenotypic records available in the training set. We observed large variability for the

343    difference in prediction accuracy, especially when the training set was small. This

344    variability was larger in Top40k than in ChipPlusSign, in a way that shrinkage of

345    variation as the training set was larger was more noticeable in ChipPlusSign. Gains in

346    prediction accuracy were low-to-moderate in the most favourable cases. In the most

347    unfavourable ones we observed large losses in prediction accuracy for Top40k but

348    more restrained losses for ChipPlusSign with moderate training set sizes. For both

349    sets of predictor variants, there was a positive trend that supported the need for large

350    training sets. This trend was clearer in ChipPlusSign than in Top40k, because of the

351    apparent lower robustness of the latter. Results for the other sets of predictor variants

352    are provided in Additional File 2.

353        The genetic architecture of the traits was also related to the success of WGS

354    for improving prediction accuracy. As the true genetic architecture of real complex

355    traits is mostly unknown, we used simulated traits to show that traits with high

356    heritability and low number of QTN were more likely to show larger improvements in

357    predictive performances. With Top40k (Figure 4), heritability seemed to be the main

358    factor that affected the expected improvement with large training sets (from null

359    improvements when $h^2$=0.1 to improvements of approximately 0.05 when $h^2$=0.5,

360    regardless of number of QTN, with a training set of 92k individuals). With

361    ChipPlusSign (Figure 5), the expected improvements with the same training set (92k

362    individuals) were not only greater in magnitude but depended on both heritability and

363    number of QTN (from null improvements when $h^2$=0.1 to improvements of

364    approximately 0.03 to 0.10 when $h^2$=0.5 with a number of 100 to 10k QTN,

365    respectively). Results confirmed the trends observed for the real traits (Figures 4 and

366    5); for instance, the higher robustness of ChipPlusSign compared to Top40k.

367        We observed diminishing returns when we increased the density of the

368    predictor variants. Increasing the number of predictor variants from the 40k in

369    Top40k to 75k selected in the same way yielded small improvements in prediction

370    accuracy compared to Top40k, but increases up to 100k variants provided smaller or

371    null additional gains (Additional File 3).

372        Splitting the original training set into two exclusive subsets, one for the

373    GWAS-based preselection of the variants and one for the training of the predictive

374    equation did not improve the prediction accuracy (Additional File 4). For

375 ChipPlusSign, this strategy reduced the bias but prediction accuracy decreased too,

376 probably because of the smaller subset available for training the predictive equation.

377

### Prediction accuracy in multi-line scenarios

378 The performance of genomic predictors trained with multi-line datasets was

379 systematically lower than in the within-line scenarios (Additional File 5).

380 Nonetheless, the ML-ChipPlusSign predictor variants in general increased prediction

381 accuracy relative to ML-Chip (Figure 6). The increase in genomic prediction accuracy

382 for each line was largely dependent on the number of individuals of each line in the

383 training set. Therefore, the greatest improvements were achieved for the largest lines.

384 However, in the multi-line scenarios we observed increases of prediction accuracy for

385 some traits and lines for which no improvements were observed in the within-line

386 scenarios (Figure 7). In contrast, results for ML-Top40k were not robust (Additional

387 File 6).

388

### Association tests

389 First, we assessed the performance of GWAS using the simulated traits. Table

390 4 shows the number of regions with significant associations that were detected using

391 either Chip or WGS, and whether they contained zero, one or multiple true QTN. The

392 WGS allowed the detection of a much larger proportion of true QTN than the Chip,

393 especially for the traits with high heritability and with large population sizes. The

394 most favourable scenarios for identifying regions that contained unequivocally a

395 single QTN with WGS were those in which the trait was controlled by a low number

396 of true QTN. However, even though the genetic architecture was very simple and

397 consisted of additive effects alone, the regions with significant associations only

398  captured a small fraction of the QTN that segregated within each line. Moreover,

399  using WGS also increased the number of regions with significant associations that

400  contained no QTN, which could therefore be considered as false positives. Some of

401  the selected regions contained multiple QTN, which could indicate either a 'hit by

402  chance' or an inability to disentangle multiple causal variants. While false positives

403  also occur with Chip, their incidence was more severe with the WGS, especially for

404  traits with a large number of QTN. Large population sizes further aggravated the

405  inflation of genome-wide p-values.

406  Despite this, with the real traits we found that GWAS using WGS can

407  contribute to a better understanding of the genetic mechanisms that underlie the traits

408  of interest. To illustrate this, we examined the GWAS results for BFT in line A, for

409  which a large number of phenotypic records were available. Figure 8 shows the

410  results for chromosome 1 as an example, while Additional File 7 shows the results for

411  six genomic regions of interest. The main genomic regions and candidate genes

412  associated to BFT detected with Chip in the same genetic lines studied here were

413  reported elsewhere [64]. We will use the candidate genes reported there to refer to the

414  genomic regions with significant associations. Using Chip, we identified 6 genomic

415  regions ($p<10^{-6}$). Using WGS (with a more stringent significance threshold of $p<10^{-9}$

416  to focus on the most significant associations), we confirmed 3 of these genomic

417  regions that co-located to candidate genes *MC4R*, *DOLK*, and *DGKI* or *PTN*.

418  However, the most associated variants in each of these genomic regions located

419  outside the coding region of these putative causal genes. These signals sometimes had

420  very strong evidence of association for some variants that were relatively distant from

421  our candidate functional gene, which could cast doubts about the fine-mapping of the

422  causal mutation. The region at SSC18, 9–13 Mb, contained two candidate genes

423    *DGKI* and *PTN,* but the WGS revealed significantly associated variants within *DGKI*

424    and none within *PTN,* despite that the strongest associations were away from both

425    genes at 10.5-11 Mb. Using the WGS we also detected 24 additional genomic regions

426    that contained candidate genes such as *CYB5R4, IGF2,* and *LEPR.* These genes were

427    previously detected in other lines using the Chip but not in this one [64], sometimes

428    because there were no markers for the associated region in Chip (SSC2, 0–4 Mb). The

429    region at SSC1, 52.5–53.5 Mb, showed many significant variants that encompassed

430    not only the previously identified candidate gene *CYB5R4,* but also *MRAP2*

431    (annotated with functions on feeding behavior and energy homeostasis). In contrast,

432    candidate gene *LEPR* was located within the region at SSC6, 146.5–147.0 Mb, where

433    many significant variants were located, although the most significant variants were

434    not in the coding regions of the gene. Using the WGS we also identified additional

435    candidate genes that had not been previously detected in any of the lines, such as

436    *CYP24A1* (annotated with functions on fatty acid omega-oxidation and vitamin D

437    metabolism; not shown). For many of the other genomic regions, it was difficult to

438    pinpoint a candidate gene with the available information or there were no annotated

439    genes.

440

## Discussion

441        Our results evidenced the potential for WGS to improve genomic prediction

442    accuracy in intensely selected pig lines, provided that the training sets are large

443    enough. Improvements achieved so far were modest at best. On one hand, these

444    modest improvements indicated that the strategies that we tested were likely

445    suboptimal. On the other hand, the positive trend for the largest training sets indicated

446    that we might have not reached the critical mass of data that is needed to leverage the

447  potential of WGS, especially in scenarios where genomic prediction with marker

448  arrays is already yielding high accuracy. The results from several traits and lines with

449  different training set sizes and the use of simulated phenotypes allowed us to identify

450  the most favourable scenarios for genomic prediction with WGS. We will discuss (1)

451  the prediction accuracy that we achieved with WGS compared to commercial marker

452  array data and the scenarios in which WGS may become beneficial, and (2) the

453  potential pitfalls for its effective implementation and the need for an optimised

454  strategy.

455

### Prediction accuracy with whole-genome sequence data

456  We compared the genomic prediction accuracy of the current marker array

457  (Chip) with sets of preselected sequence variants in a way that the number of variants

458  remained similar across sets. Improvements of prediction accuracy can be limited if

459  current marker arrays are already sufficiently dense to capture a large proportion of

460  the genetic variance in intensely selected livestock populations. These populations

461  typically have small effective population size [10,19]. Nevertheless, modest

462  improvements have been achieved under certain scenarios. In our study, the most

463  robust results were obtained for the ChipPlusSign set, where variants that showed

464  statistically significant associations to the trait were preselected and added to the

465  information from the marker array. This is consistent with previous reports that

466  showed an improvement of prediction accuracy under similar approaches [29–32].

467  We added 23 to 1083 significant variants to those in Chip in different scenarios. In the

468  most successful ones, at least around 200 significant variants were added and average

469  improvements of prediction accuracy of 2.5 percentage points were observed with

470  training sets of around 80k individuals. In other instances, however, additions of a

471 larger number of variants have been proposed. The addition of 1623 variants

472 (preselected as the combination of 3-5 variants for each of the top QTL per trait and

473 breed) to a 50k array increased prediction reliability (accuracy squared) by up to 5

474 percentage points in Nordic cattle [29]. Adding the 16k SNPs with largest estimated

475 effects to a 60k array increased prediction reliability on average by 2.7 (up to 4.8)

476 percentage points in Holstein cattle [30]. For the custom 50k array for Hanwoo cattle,

477 it has been reported that adding at least around 12k SNPs (3k for each of four traits)

478 improved prediction accuracy by up to ~6 percentage points [32]. The addition of

479 ~400 variants preselected by GWAS with regional heritability mapping to a 50k array

480 increased prediction accuracy by 9 percentage points in sheep [31]. In other cases in

481 Nordic cattle, however, the addition of ~1500 variants preselected by GWAS to a 54k

482 panel produced negligible improvements in the prediction of traits with low

483 heritability [65].

484 Preselecting an entirely new set of predictor variants from WGS, as in

485 Top40k, proved more challenging than ChipPlusSign. In Top40k, we preselected the

486 variants with the lowest p-value in each of consecutive non-overlapping 55-kb

487 windows along the genome. This strategy did not perform much differently from just

488 taking random variants from these windows, as in Rand40k. One possible reason for

489 these results is that at this variant density, random variants effectively tag the same

490 associations as Top40k thanks to linkage disequilibrium. Denser sets of predictor

491 variants provided only small further improvements of prediction accuracy with

492 diminishing returns.

493 The modest performance of ChipPlusSign and Top40k could also be a

494 consequence of the difficulty for fine-mapping causal variants through GWAS with

495 WGS. Theoretically, the identification of all causal variants associated with a trait

496 should enable the improvement of prediction accuracy [12]. Even though WGS allows

497 the detection of a very large number of associations, problems such as false positives

498 or p-value inflation also become more severe in a way that added noise might offset

499 the detected signal. For instance, results in cattle showed that GWAS with WGS did

500 not detect clearer associated regions relative to marker arrays and failed to capture

501 QTL for genomic prediction [14], as the effect of potential QTL were spread across

502 multiple variants. Therefore, WGS performed better with simple genetic architectures

503 (i.e., traits with low number of QTN). This is consistent with expectations and

504 simulation results [8] that indicated that the benefit of using WGS for genomic

505 prediction would be limited by the number and size of QTN. When there are many

506 QTN with small effects it becomes much more difficult to properly estimate their

507 effects accurately. Therefore, for largely polygenic traits (as most traits of interest in

508 livestock production), training sets need to be very large before WGS can increase

509 prediction accuracy [8].

510 The advantage of using WGS might be limited by the current training set

511 sizes, especially in scenarios where marker arrays are already yielding high prediction

512 accuracy [14,20]. Multi-line training sets could be particularly beneficial with the use

513 of WGS because they allow a larger training set with low pairwise relationship degree

514 among individuals. Previous simulations suggested that WGS might be the most

515 beneficial with multi-breed reference panels [66], especially for numerically small

516 populations. Our results with a multi-line training set indicated that WGS can improve

517 prediction accuracy in scenarios that are less optimised than within-line genomic

518 prediction. The average improvements of prediction accuracy of 4.2 percentage points

519 were observed for the populations that contributed around 80k individuals to the

520 training set. However, in general those predictions were still less accurate than using

521    variants preselected under within-line training sets. In our multi-line scenarios we

522    only used variation that segregated across all seven lines. We observed that

523    population-specific variation accounted only for small fractions of genetic variance

524    [53] and it seems unlikely that they would contribute much to prediction accuracy

525    across breeds. Another possible obstacle is the differences in the allele substitution

526    effects of the causal mutations across breeds. This can be caused by differences in

527    allele frequency, contributions of non-additive effects and different genetic

528    backgrounds, or even gene-by-environment interactions among others [24,67].

529        We observed low robustness of genomic prediction with WGS across traits

530    and lines, and drops in prediction accuracy in those scenarios where genomic

531    prediction with WGS failed. Regarding bias, we did not observe a systematic increase

532    for ChipPlusSign despite using the same individuals for variant preselection and for

533    training the predictors [16,63]. When we split the training set into two subsets, one for

534    GWAS-based variant preselection and the other for training of the predictive

535    equations, we did not observe any improvement in accuracy or bias. One hypothesis is

536    that both subsets belong to the same population and therefore retained similar inter-

537    relationship degrees (i.e., they are not strictly independent sets of individuals).

538    Moreover, the reduction in individuals available for training the predictors negatively

539    affected prediction accuracy.

540        We did not directly test persistence of prediction accuracy, but previous

541    studies with real data found no higher persistence of prediction accuracy for WGS,

542    not even with low degree of relationship between training and testing sets [14]. We

543    would expect such obstacles to persistence of accuracy until causal variants can be

544    successfully identified.

545

## Suboptimal strategy and pitfalls

546        The use of WGS for genomic prediction can only be reached after many other

547    steps are completed to produce the genotypes at whole-genome level. Each of these

548    steps has its pitfalls. It is unavoidable that the success of using of WGS is sensitive

549    not only to the prediction methodology itself but also to the strategy followed until

550    genotyping. This strategy includes the choice of which individuals to sequence, the

551    bioinformatics pipeline to call variants, the imputation of the WGS and choice of

552    variant filters. When combined with the multiplicity of prediction methods and the

553    preselection of predictor variants (which is unavoidable with current datasets,

554    predictive methodologies and computational capacities), there are many options and

555    variables in the whole process that can affect the final result and that are not yet well

556    understood. Therefore, a much greater effort for optimising such strategies is required.

557    Here we tested relatively simple approaches to see how they performed with large

558    WGS datasets. We have discussed what in our opinion are the main pitfalls of our

559    approach for selection of the individuals to sequence [55] and the biases that may

560    appear during processing of sequencing reads [50] elsewhere, and therefore here we

561    will focus on imputation of WGS and its use for genomic prediction.

562

563    *Imputation accuracy*

564        It is widely recognized that imputation from marker arrays to WGS from very

565    few sequenced individuals can introduce genotyping errors and that genotype

566    uncertainty can be high [17,21,68,69]. The accuracy of the imputed WGS is one of the

567    main factors that may limit its performance for genomic prediction. In a simulation

568    study, van den Berg et al. [17] quantified the impact of imputation errors on

569    prediction accuracy and showed that prediction accuracy decreases as errors

570    accumulate, especially in the testing set.

571         Imputation of WGS is particularly challenging because typically we have to

572    impute a very large number of variants for a very large number of individuals from

573    very few sequenced individuals. We assessed the imputation accuracy of our approach

574    elsewhere [41,55] and recommended that ~2% of the population should be sequenced.

575    In our study, line D was the line where prediction accuracy with Top40k performed

576    the worst, mostly performing below Chip predictors. In this line, only 0.9% of the

577    individuals in the population had been sequenced and therefore lower imputation

578    accuracy could be expected. Although there was not enough evidence for establishing

579    a link between these two features (sequencing effort and prediction accuracy), we

580    recommend cautious design of a sequencing strategy that is suited to the intended

581    imputation method [55].

582         Prediction accuracy could be improved by accounting for genotype uncertainty

583    of the imputed WGS. For that, it could be advantageous to use allele dosages rather

584    than best-guess genotypes [69], although most current implementations cannot handle

585    such information.

586

587    ***Preselection of predictor variants***

588         Simply using WGS to increase the number of markers does not improve

589    prediction accuracy [18,21,24]. Due to the large dimensionality of WGS, there is a

590    need to remove uninformative variants [24,31,66,68,70]. Predictor variants must be

591    causal or at least informative of the causal variants, which depends on the distance

592    between the markers and the causal variants [13]. For this reason, variants that are in

593    weak linkage disequilibrium with causal mutations have a 'dilution' effect, i.e., they

594   add noise and limit prediction accuracy [24,31,70]. However, if too stringent filters

595   are applied during preselection of predictor variants, there is a risk of removing true

596   causal variants, and that would debilitate persistence of accuracy across generations

597   and across populations [66,71]. For instance, the impact of removing predictor

598   variants with low minor allele frequency can vary depending on the minor allele

599   frequency of the causal variants as well as the distance between predictor and causal

600   variants [13]. Losing causal or informative variants would negatively affect multi-line

601   or multi-breed prediction.

602        A popular strategy to preselect the predictor variants is based on association

603   tests. Genome-wide association studies on WGS are expected to confirm associations

604   that were already detected with marker arrays and identify novel associations (e.g.,

605   [36,72]). However, preliminary inspection of our GWAS results for the real traits

606   showed that the added noise could easily offset the added information and fine-

607   mapping remains challenging. Multi-breed GWAS [4] and meta-analyses [73] are

608   suitable alternatives for GWAS to accommodate much larger population sizes and for

609   combining results of populations with diverse genetic backgrounds. Multi-breed

610   GWAS can be more efficient to identify informative variants than single-breed

611   GWAS, which may benefit even prediction within lines [74]. Because the signal of

612   some variants may go undetected for some traits but not for other correlated traits,

613   combining GWAS information of several traits can also help identifying weak or

614   moderate associations [25]. We did not test whether combining the significant

615   markers from the different single-trait GWAS yielded greater improvements in

616   prediction accuracy [29,32]. Multi-trait GWAS models could be more suited for that

617   purpose [72,75]. To improve fine-mapping, other GWAS models that incorporate

618     biological information have been proposed (e.g., functional annotation [76] or

619     metabolomics [77]).

620         There have been other suggested methods that may improve variant

621     preselection for genomic prediction. VanRaden et al. [30] suggested that preselecting

622     variants based on the genetic variance that they contribute rather than the significance

623     of the association could be more advantageous, as the former would indirectly

624     preselect variance with higher minor allele frequency. Other authors proposed

625     preselection of variants using statistics that do not depend on GWAS, such as the

626     fixation index ($F_{ST}$) score between groups of individuals with high and low phenotype

627     values [70], as an alternative to avoid the negative impact of spurious associations.

628         Preselecting predictor variants based on functional annotation was not useful,

629     as it reduced prediction accuracy in several traits and lines. Previous studies showed

630     that subsets of variants based on functionality either did not improve or reduced

631     prediction accuracy [20] and that adding preselected variants from coding regions to

632     marker arrays produced lower prediction accuracy than just adding the same number

633     of variants without considering functional classification [32]. A plausible explanation

634     is that functional variants are enriched for lower minor allele frequency, which can be

635     less informative for prediction [13]. Furthermore, functional annotation does not

636     necessarily capture true effects, and the method we used is biased towards protein-

637     coding variants, which may lead to an underrepresentation of functional non-coding

638     variants that may explain a large fraction of quantitative trait variance. Xiang et al.

639     [78] found that expression QTL and non-coding variants explained more variation in

640     quantitative traits in cattle than protein-coding functional variants. When functional

641     annotation is not considered, intergenic variants are more likely to be preselected by

642     chance. Such variants tend to be more common and widespread across populations,

643    and therefore can act as tag variants and capture much larger fractions of trait variance

644    [53].

645         Another popular strategy to reduce the number of variants is to prune variants

646    based on linkage disequilibrium (LDTags). This strategy performed very poorly in our

647    populations. Other studies reported different outcomes, where pruning for $r^2>0.9$

648    provided positive results [18,21]. It is possible that this was in part due to the stringent

649    threshold ($r^2>0.1$) that we used in order to retain only a small number of variants.

650

651    ***New models and methods***

652         It is also likely that models, methods, and their implementations need to be

653    improved to handle the complexity of WGS and to efficiently estimate marker effects

654    of so many variants with high accuracy, among other features. This is a very active

655    area of research and multiple novel methodologies have been proposed over the last

656    years. Some examples are a combination of subsampling and Gibbs sampling [79],

657    and a model that simultaneously fits a GBLUP term for a polygenic effect and a

658    BayesC term for variants with large effects selected by the model (BayesGC) [26].

659    Testing alternative models and methods for genomic prediction was out of the scope

660    of this report. However, together with refinements in the preselection of predictor

661    variants, it remains an interesting avenue for further optimisation of the analysis

662    pipeline.

663         Some of the most promising methods are designed to incorporate prior

664    biological information into the models. One of such methods is BayesRC [23], which

665    extends BayesR by assigning flatter prior distributions to classes of variants that are

666    more likely to be causal [19,22]. Similarly, GFBLUP [80] could be used to

667    incorporate prior biological information from either QTL databases or GWAS as

668  genomic features [21,35,68]. The model MBMG [27], which fits two genomic

669  relationship matrices according to prior biological information, has also been

670  proposed for multi-breed scenarios to improve genomic prediction in small

671  populations. Haplotype-based prediction methods could provide greater prediction

672  accuracy with WGS than SNP-based methods in pigs [81] and cattle [82]. These

673  methods reduce the number of model dimensions. However, the uptake of such

674  methods has been limited so far due to their greater complexity, for example, to define

675  haplotype blocks.

676

## Conclusion

677      Our results evidenced the potential for WGS to improve genomic prediction

678  accuracy in intensely selected pig lines. The performance of each set of predictor

679  variants was not robust across traits and lines and the improvements that we achieved

680  so far were modest at best. The most robust results were obtained when variants that

681  showed statistically significant associations to the trait were preselected and added to

682  the marker array. With this method, average improvements of prediction accuracy of

683  2.5 and 4.2 percentage points were observed in within-line and multi-line scenarios,

684  respectively, with training sets of around 80k individuals. We would expect that a

685  combination of larger training sets and improved pipelines could help achieve greater

686  improvements of prediction accuracy. The robustness of the whole strategy for

687  generating WGS at the population level must be carefully stress-tested and further

688  optimised.

689

## Ethics approval and consent to participate

690  The samples used in this study were derived from the routine breeding activities of

691  PIC.

## Consent for publication

692  Not applicable.

## Availability of data and material

693  The software packages AlphaSeqOpt, AlphaPhase, AlphaImpute and AlphaPeel are

694  available from the AlphaGenes website (http://www.alphagenes.roslin.ed.ac.uk). The

695  datasets generated and analysed in this study are derived from the PIC breeding

696  programme and not publicly available.

## Competing interests

697  The authors declare that they have no competing interests. BDV, CYC, and WOH are

698  employees of Genus PIC.

## Funding

## Authors' contributions

704  RRF, GG and JMH designed the study; CYC assisted in preparing the datasets; RRF,

705  AW and MJ performed the analyses; RRF wrote the first draft; AW, CYC, BDV,

706    WHO, GG and JMH assisted in the interpretation of the results and provided

707    comments on the manuscript. All authors read and approved the final manuscript.

## Acknowledgements

710

## References

711    1. Pasaniuc B, Rohland N, McLaren PJ, Garimella K, Zaitlen N, Li H, et al.
712    Extremely low-coverage sequencing and imputation increases power for genome-
713    wide association studies. Nat Genet. 2012;44:631–5.

714    2. Daetwyler HD, Capitan A, Pausch H, Stothard P, van Binsbergen R, Brondum RF,
715    et al. Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and
716    complex traits in cattle. Nat Genet. 2014;46:858–65.

717    3. Nicod J, Davies RW, Cai N, Hassett C, Goodstadt L, Cosgrove C, et al. Genome-
718    wide association of multiple complex traits in outbred mice by ultra-low-coverage
719    sequencing. Nat Genet. 2016;48:912–8.

720    4. Sanchez M-P, Govignon-Gion A, Croiseau P, Fritz S, Hozé C, Miranda G, et al.
721    Within-breed and multi-breed GWAS on imputed whole-genome sequence variants
722    reveal candidate mutations affecting milk protein composition in dairy cattle. Genet
723    Sel Evol. 2017;49:68.

724    5. Das A, Panitz F, Gregersen VR, Bendixen C, Holm L-E. Deep sequencing of
725    Danish Holstein dairy cattle for variant detection and insight into potential loss-of-
726    function variants in protein coding genes. BMC Genomics. 2015;16:1043.

727    6. Gudbjartsson DF, Helgason H, Gudjonsson SA, Zink F, Oddson A, Gylfason A, et
728    al. Large-scale whole-genome sequencing of the Icelandic population. Nat Genet.
729    2015;47:435–44.

730    7. VanRaden PM. Symposium review: How to implement genomic selection. J Dairy
731    Sci. 2020;103:5291–301.

732    8. Clark SA, Hickey JM, van der Werf JH. Different models of genetic variation and
733    their effect on genomic evaluation. Genet Sel Evol. 2011;43:18.

734    9. Druet T, Macleod IM, Hayes BJ. Toward genomic prediction from whole-genome
735    sequence data: impact of sequencing design on genotype imputation and accuracy of
736    predictions. Heredity. 2014;112:39–47.

737    10. MacLeod IM, Hayes BJ, Goddard ME. The Effects of Demography and Long-
738    Term Selection on the Accuracy of Genomic Prediction with Sequence Data.
739    Genetics. 2014;198:1671–84.

740    11. Meuwissen T, Goddard M. Accurate Prediction of Genetic Values for Complex
741    Traits by Whole-Genome Resequencing. Genetics. 2010;185:623–31.

742    12. Pérez-Enciso M, Rincón JC, Legarra A. Sequence- vs. chip-assisted genomic
743    selection: accurate biological information is advised. Genet Sel Evol. 2015;47:43.

744    13. van den Berg I, Boichard D, Guldbrandtsen B, Lund MS. Using Sequence
745    Variants in Linkage Disequilibrium with Causative Mutations to Improve Across-
746    Breed Prediction in Dairy Cattle: A Simulation Study. G3 GenesGenomesGenetics.
747    2016;6:2553–61.

748    14. van Binsbergen R, Calus MPL, Bink MCAM, van Eeuwijk FA, Schrooten C,
749    Veerkamp RF. Genomic prediction using imputed whole-genome sequence data in
750    Holstein Friesian cattle. Genet Sel Evol. 2015;47:71.

751    15. Calus MPL, Bouwman AC, Schrooten C, Veerkamp RF. Efficient genomic
752    prediction based on whole-genome sequence data using split-and-merge Bayesian
753    variable selection. Genet Sel Evol. 2016;48:49.

754    16. Veerkamp RF, Bouwman AC, Schrooten C, Calus MPL. Genomic prediction
755    using preselected DNA variants from a GWAS with whole-genome sequence data in
756    Holstein–Friesian cattle. Genet Sel Evol. 2016;48:95.

757    17. van den Berg I, Bowman PJ, MacLeod IM, Hayes BJ, Wang T, Bolormaa S, et al.
758    Multi-breed genomic prediction using Bayes R with sequence data and dropping
759    variants with a small effect. Genet Sel Evol. 2017;49:70.

760    18. Frischknecht M, Meuwissen THE, Bapst B, Seefried FR, Flury C, Garrick D, et al.
761    Short communication: Genomic prediction using imputed whole-genome sequence
762    variants in Brown Swiss Cattle. J Dairy Sci. 2018;101:1292–6.

763    19. Hayes BJ, MacLeod IM, Daetwyler HD, Bowman PJ, Chamberlain AJ, Vander
764    Jagt CJ, et al. Genomic prediction from whole genome sequence in livestock: the
765    1000 Bull Genomes Project. Proc 10th World Congr Genet Appl Livest Prod
766    WCGALP. Vancouver, BC, Canada; 2014. p. 183.

767    20. Heidaritabar M, Calus MPL, Megens H-J, Vereijken A, Groenen MAM,
768    Bastiaansen JWM. Accuracy of genomic prediction using imputed whole-genome
769    sequence data in white layers. J Anim Breed Genet. 2016;133:167–79.

770    21. Song H, Ye S, Jiang Y, Zhang Z, Zhang Q, Ding X. Using imputation-based
771    whole-genome sequencing data to improve the accuracy of genomic prediction for
772    combined populations in pigs. Genet Sel Evol. 2019;51:58.

773    22. Zhang C, Kemp RA, Stothard P, Wang Z, Boddicker N, Krivushin K, et al.
774    Genomic evaluation of feed efficiency component traits in Duroc pigs using 80K,
775    650K and whole-genome sequence variants. Genet Sel Evol. 2018;50:14.

776   23. MacLeod IM, Bowman PJ, Vander Jagt CJ, Haile-Mariam M, Kemper KE,
777   Chamberlain AJ, et al. Exploiting biological priors and sequence variants enhances
778   QTL discovery and genomic prediction of complex traits. BMC Genomics.
779   2016;17:144.

780   24. Raymond B, Bouwman AC, Schrooten C, Houwing-Duistermaat J, Veerkamp RF.
781   Utility of whole-genome sequence data for across-breed genomic prediction. Genet
782   Sel Evol. 2018;50:27.

783   25. Xiang R, MacLeod IM, Daetwyler HD, de Jong G, O'Connor E, Schrooten C, et
784   al. Genome-wide fine-mapping identifies pleiotropic and functional variants that
785   predict many traits across global cattle populations. Nat Commun. 2021;12:860.

786   26. Meuwissen T, van den Berg I, Goddard M. On the use of whole-genome sequence
787   data for across-breed genomic prediction and fine-scale mapping of QTL. Genet Sel
788   Evol. 2021;53:19.

789   27. Raymond B, Bouwman AC, Wientjes YCJ, Schrooten C, Houwing-Duistermaat J,
790   Veerkamp RF. Genomic prediction for numerically small breeds, using models with
791   pre-selected and differentially weighted markers. Genet Sel Evol. 2018;50:49.

792   28. Moghaddar N, Brown DJ, Swan AA, Gurman PM, Li L, Werf JH. Genomic
793   prediction in a numerically small breed population using prioritized genetic markers
794   from whole-genome sequence data. J Anim Breed Genet. 2021;

795   29. Brøndum RF, Su G, Janss L, Sahana G, Guldbrandtsen B, Boichard D, et al.
796   Quantitative trait loci markers derived from whole genome sequence data increases
797   the reliability of genomic prediction. J Dairy Sci. 2015;98:4107–16.

798   30. VanRaden PM, Tooker ME, O'Connell JR, Cole JB, Bickhart DM. Selecting
799   sequence variants to improve genomic predictions for dairy cattle. Genet Sel Evol.
800   2017;49:32.

801   31. Al Kalaldeh M, Gibson J, Duijvesteijn N, Daetwyler HD, MacLeod I, Moghaddar
802   N, et al. Using imputed whole-genome sequence data to improve the accuracy of
803   genomic prediction for parasite resistance in Australian sheep. Genet Sel Evol.
804   2019;51:32.

805   32. Lopez BIM, An N, Srikanth K, Lee S, Oh J-D, Shin D-H, et al. Genomic
806   Prediction Based on SNP Functional Annotation Using Imputed Whole-Genome
807   Sequence Data in Korean Hanwoo Cattle. Front Genet. 2021;11:603822.

808   33. Hayes BJ, Daetwyler HD. 1000 Bull Genomes Project to Map Simple and
809   Complex Genetic Traits in Cattle: Applications and Outcomes. Annu Rev Anim
810   Biosci. 2019;7:89–102.

811   34. Sanchez M-P, Guatteo R, Davergne A, Saout J, Grohs C, Deloche M-C, et al.
812   Identification of the ABCC4, IER3, and CBFA2T2 candidate genes for resistance to
813   paratuberculosis from sequence-based GWAS in Holstein and Normande dairy cattle.
814   Genet Sel Evol. 2020;52:14.

815  35. Yang R, Xu Z, Wang Q, Zhu D, Bian C, Ren J, et al. Genome-wide association
816  study and genomic prediction for growth traits in yellow-plumage chicken using
817  genotyping-by-sequencing. Genet Sel Evol. 2021;53:82.

818  36. Yan G, Liu X, Xiao S, Xin W, Xu W, Li Y, et al. An imputed whole-genome
819  sequence-based GWAS approach pinpoints causal mutations for complex traits in a
820  specific swine population. Sci China Life Sci. 2021;

821  37. Yang R, Guo X, Zhu D, Tan C, Bian C, Ren J, et al. Accelerated deciphering of
822  the genetic architecture of agricultural economic traits in pigs using a low-coverage
823  whole-genome sequencing strategy. GigaScience. 2021;10:giab048.

824  38. Li Y, Sidore C, Kang HM, Boehnke M, Abecasis GR. Low-coverage sequencing:
825  Implications for design of complex trait association studies. Genome Res.
826  2011;21:940–51.

827  39. Hickey JM. Sequencing millions of animals for genomic selection 2.0. J Anim
828  Breed Genet. 2013;130:331–2.

829  40. Hickey JM, Gorjanc G, Cleveland MA, Kranis A, Jenko J, Mésázros G, et al.
830  Sequencing Millions of Animals for Genomic Selection 2.0. Proc 10th World Congr
831  Genet Appl Livest Prod WCGALP. Vancouver, BC, Canada; 2014. p. 377.

832  41. Ros-Freixedes R, Whalen A, Chen C-Y, Gorjanc G, Herring WO, Mileham AJ, et
833  al. Accuracy of whole-genome sequence imputation using hybrid peeling in large
834  pedigreed livestock populations. Genet Sel Evol. 2020;52:17.

835  42. Gonen S, Ros-Freixedes R, Battagin M, Gorjanc G, Hickey JM. A method for the
836  allocation of sequencing resources in genotyped livestock populations. Genet Sel
837  Evol. 2017;49:47.

838  43. Ros-Freixedes R, Gonen S, Gorjanc G, Hickey JM. A method for allocating low-
839  coverage sequencing resources by targeting haplotypes rather than individuals. Genet
840  Sel Evol. 2017;49:78.

841  44. Hickey JM, Kinghorn BP, Tier B, Wilson JF, Dunstan N, van der Werf JH. A
842  combined long-range phasing and long haplotype imputation method to impute phase
843  for SNP genotypes. Genet Sel Evol. 2011;43:12.

844  45. Hickey JM, Kinghorn BP, Tier B, van der Werf JH, Cleveland MA. A phasing
845  and imputation method for pedigreed populations that results in a single-stage
846  genomic evaluation. Genet Sel Evol. 2012;44:9.

847  46. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina
848  sequence data. Bioinformatics. 2014;30:2114–20.

849  47. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-
850  MEM. arXiv. 2013;1303.3997v1 [q – bio.GN].

48. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat Genet. 2011;43:491–8.

49. Poplin R, Ruano-Rubio V, DePristo MA, Fennell TJ, Carneiro MO, Van der Auwera GA, et al. Scaling accurate genetic variant discovery to tens of thousands of samples. bioRxiv. 2018;10.1101/201178.

50. Ros-Freixedes R, Battagin M, Johnsson M, Gorjanc G, Mileham AJ, Rounsley SD, et al. Impact of index hopping and bias towards the reference allele on accuracy of genotype calls from low-coverage sequencing. Genet Sel Evol. 2018;50:64.

51. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics. 2009;25:2078–9.

52. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. Bioinformatics. 2011;27:2156–8.

53. Ros-Freixedes R, Valente B, Chen C-Y, Herring WO, Gorjanc G, Hickey JM, et al. Rare and population-specific functional variants across pig lines. bioRxiv [Internet]. 2022; Available from: https://doi.org/10.1101/2022.02.01.478603

54. Whalen A, Ros-Freixedes R, Wilson DL, Gorjanc G, Hickey JM. Hybrid peeling for fast and accurate calling, phasing, and imputation with sequence data of any coverage in pedigrees. Genet Sel Evol. 2018;50:67.

55. Ros-Freixedes R, Whalen A, Gorjanc G, Mileham AJ, Hickey JM. Evaluation of sequencing strategies for whole-genome imputation with hybrid peeling. Genet Sel Evol. 2020;52:18.

56. VanRaden PM, Wiggans GR. Derivation, Calculation, and Use of National Animal Model Information. J Dairy Sci. 1991;74:2737–46.

57. Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST linear mixed models for genome-wide association studies. Nat Methods. 2011;8:833–5.

58. Widmer C, Lippert C, Weissbrod O, Fusi N, Kadie C, Davidson R, et al. Further Improvements to Linear Mixed Models for Genome-Wide Association Studies. Sci Rep. 2015;4:6874.

59. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the challenge of larger and richer datasets. GigaScience. 2015;4.

60. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, et al. The Ensembl Variant Effect Predictor. Genome Biol. 2016;17:122.

61. Erbe M, Hayes BJ, Matukumalli LK, Goswami S, Bowman PJ, Reich CM, et al. Improving accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density single nucleotide polymorphism panels. J Dairy Sci. 2012;95:4114–29.

890  62. Moser G, Lee SH, Hayes BJ, Goddard ME, Wray NR, Visscher PM. Simultaneous
891  Discovery, Estimation and Prediction Analysis of Complex Traits Using a Bayesian
892  Mixture Model. Haley C, editor. PLOS Genet. 2015;11:e1004969.

893  63. MacLeod IM, Bolormaa S, Schrooten C, Goddard ME, Daetwyler H. Pitfalls of
894  pre-selecting subsets of sequence variants for genomic prediction. Proc 22nd Conf
895  Assoc Adv Anim Breed Genet AAABG. Townsville, Queensland, Australia; 2017. p.
896  141–4.

897  64. Gozalo-Marcilla M, Buntjer J, Johnsson M, Batista L, Diez F, Werner CR, et al.
898  Genetic architecture and major genes for backfat thickness in pig lines of diverse
899  genetic backgrounds. Genet Sel Evol. 2021;53:76.

900  65. Gebreyesus G, Lund MS, Sahana G, Su G. Reliabilities of Genomic Prediction for
901  Young Stock Survival Traits Using 54K SNP Chip Augmented With Additional
902  Single-Nucleotide Polymorphisms Selected From Imputed Whole-Genome
903  Sequencing Data. Front Genet. 2021;12:667300.

904  66. Iheshiulor OOM, Woolliams JA, Yu X, Wellmann R, Meuwissen THE. Within-
905  and across-breed genomic prediction using whole-genome sequence and single
906  nucleotide polymorphism panels. Genet Sel Evol. 2016;48:15.

907  67. Legarra A, Garcia-Baccino CA, Wientjes YCJ, Vitezica ZG. The correlation of
908  substitution effects across populations and generations in the presence of non-additive
909  functional gene action. PREPRINT. 2021;

910  68. Sarup P, Jensen J, Ostersen T, Henryon M, Sørensen P. Increased prediction
911  accuracy using a genomic feature model including prior information on quantitative
912  trait locus regions in purebred Danish Duroc pigs. BMC Genet. 2016;17:11.

913  69. Pausch H, MacLeod IM, Fries R, Emmerling R, Bowman PJ, Daetwyler HD, et al.
914  Evaluation of the accuracy of imputed sequence variant genotypes and their utility for
915  causal variant detection in cattle. Genet Sel Evol. 2017;49:24.

916  70. Ling AS, Hay EH, Aggrey SE, Rekaya R. Dissection of the impact of prioritized
917  QTL-linked and -unlinked SNP markers on the accuracy of genomic selection. BMC
918  Genomic Data. 2021;22:26.

919  71. Fragomeni BO, Lourenco DAL, Masuda Y, Legarra A, Misztal I. Incorporation of
920  causative quantitative trait nucleotides in single-step GBLUP. Genet Sel Evol.
921  2017;49:59.

922  72. Bolormaa S, Swan AA, Stothard P, Khansefid M, Moghaddar N, Duijvesteijn N,
923  et al. A conditional multi-trait sequence GWAS discovers pleiotropic candidate genes
924  and variants for sheep wool, skin wrinkle and breech cover traits. Genet Sel Evol.
925  2021;53:58.

926  73. van den Berg I, Xiang R, Jenko J, Pausch H, Boussaha M, Schrooten C, et al.
927  Meta-analysis for milk fat and protein percentage using imputed sequence variant
928  genotypes in 94,321 cattle from eight cattle breeds. Genet Sel Evol. 2020;52:37.

929  74. van den Berg I, Boichard D, Lund MS. Sequence variants selected from a multi-
930  breed GWAS can improve the reliability of genomic predictions in dairy cattle. Genet
931  Sel Evol. 2016;48:83.

932  75. Yoshida GM, Yáñez JM. Multi-trait GWAS using imputed high-density
933  genotypes from whole-genome sequencing identifies genes associated with body traits
934  in Nile tilapia. BMC Genomics. 2021;22:57.

935  76. Yang J, Fritsche LG, Zhou X, Abecasis G. A Scalable Bayesian Method for
936  Integrating Functional Information in Genome-wide Association Studies. Am J Hum
937  Genet. 2017;101:404–16.

938  77. Li J, Mukiibi R, Wang Y, Plastow GS, Li C. Identification of candidate genes and
939  enriched biological functions for feed efficiency traits by integrating plasma
940  metabolites and imputed whole genome sequence variants in beef cattle. BMC
941  Genomics. 2021;22:823.

942  78. Xiang R, Berg I van den, MacLeod IM, Hayes BJ, Prowse-Wilkins CP, Wang M,
943  et al. Quantifying the contribution of sequence variants with regulatory and
944  evolutionary significance to 34 bovine complex traits. Proc Natl Acad Sci.
945  2019;116:19398–408.

946  79. Xavier A, Xu S, Muir W, Rainey KM. Genomic prediction using subsampling.
947  BMC Bioinformatics. 2017;18:191.

948  80. Edwards SM, Sørensen IF, Sarup P, Mackay TFC, Sørensen P. Genomic
949  Prediction for Quantitative Traits Is Improved by Mapping Variants to Gene Ontology
950  Categories in *Drosophila melanogaster*. Genetics. 2016;203:1871–83.

951  81. Bian C, Prakapenka D, Tan C, Yang R, Zhu D, Guo X, et al. Haplotype genomic
952  prediction of phenotypic values based on chromosome distance and gene boundaries
953  using low-coverage sequencing in Duroc pigs. Genet Sel Evol. 2021;53:78.

954  82. Li H, Zhu B, Xu L, Wang Z, Xu L, Zhou P, et al. Genomic Prediction Using LD-
955  Based Haplotypes Inferred From High-Density Chip and Imputed Sequence Variants
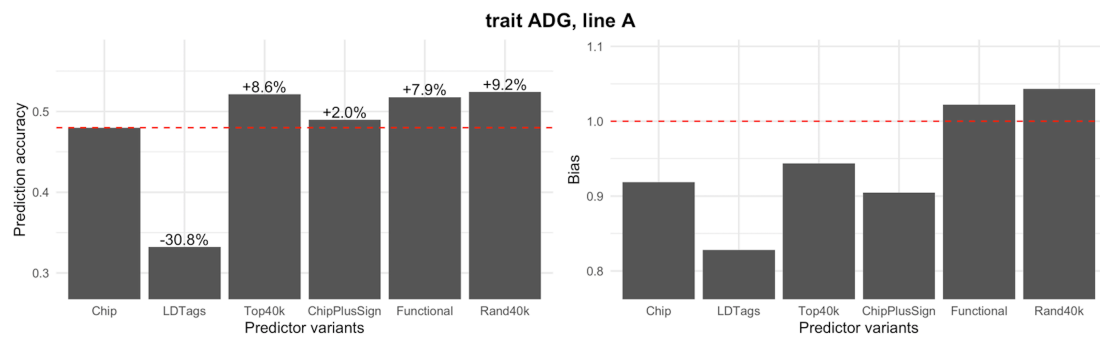956  in Chinese Simmental Beef Cattle. Front Genet. 2021;12:665382.

957

958

959

# Figures

960



961
962 **Figure 1.** Prediction accuracy for each set of predictor variants for trait ADG in line
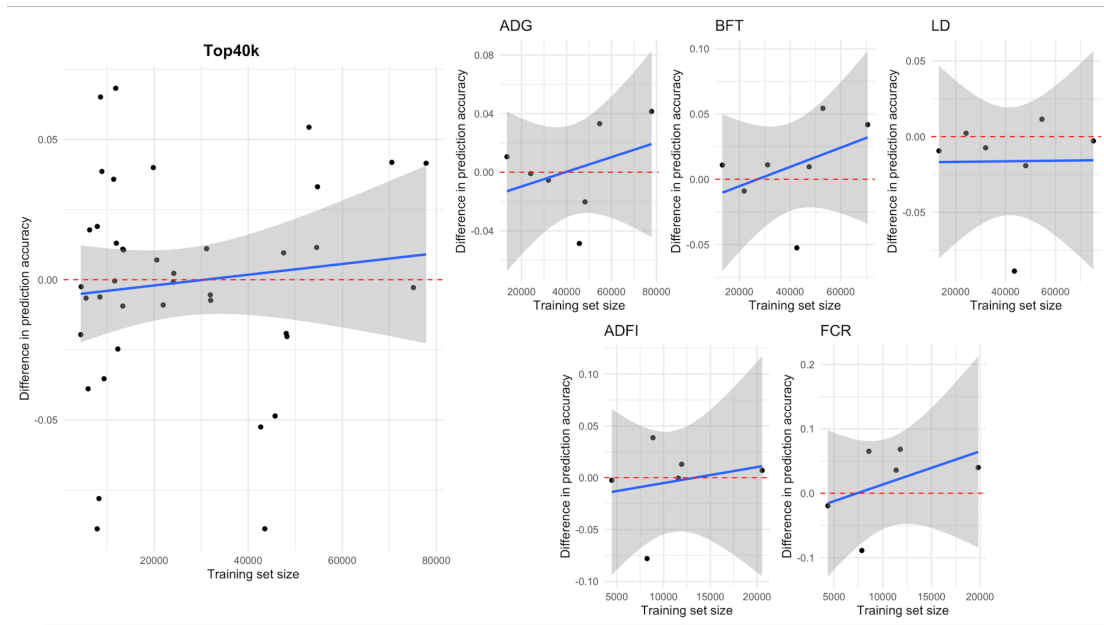
963 A. Left: Correlation (left). Dashed line at value of Chip as a reference. Values indicate

964 relative difference to Chip. Right: Bias. Dashed line at the ideal value.

965

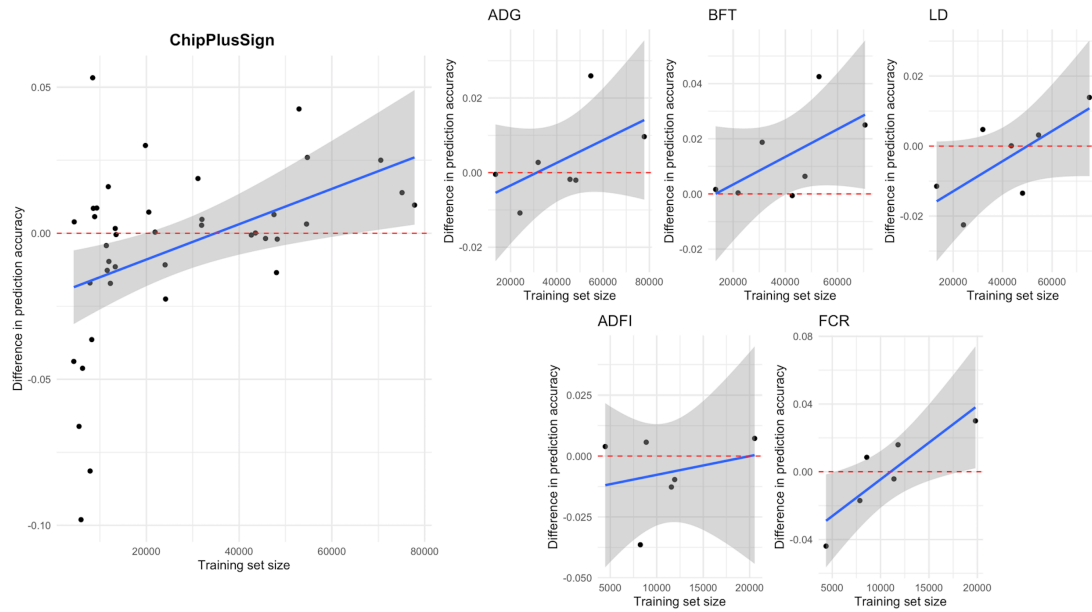**Figure 2.** Genomic prediction accuracy with the Top40k predictor variants for the real traits. The difference of prediction accuracy between Top40k and Chip is shown, for all traits and lines (left) or by trait (right). Red dashed line at 'no difference'.
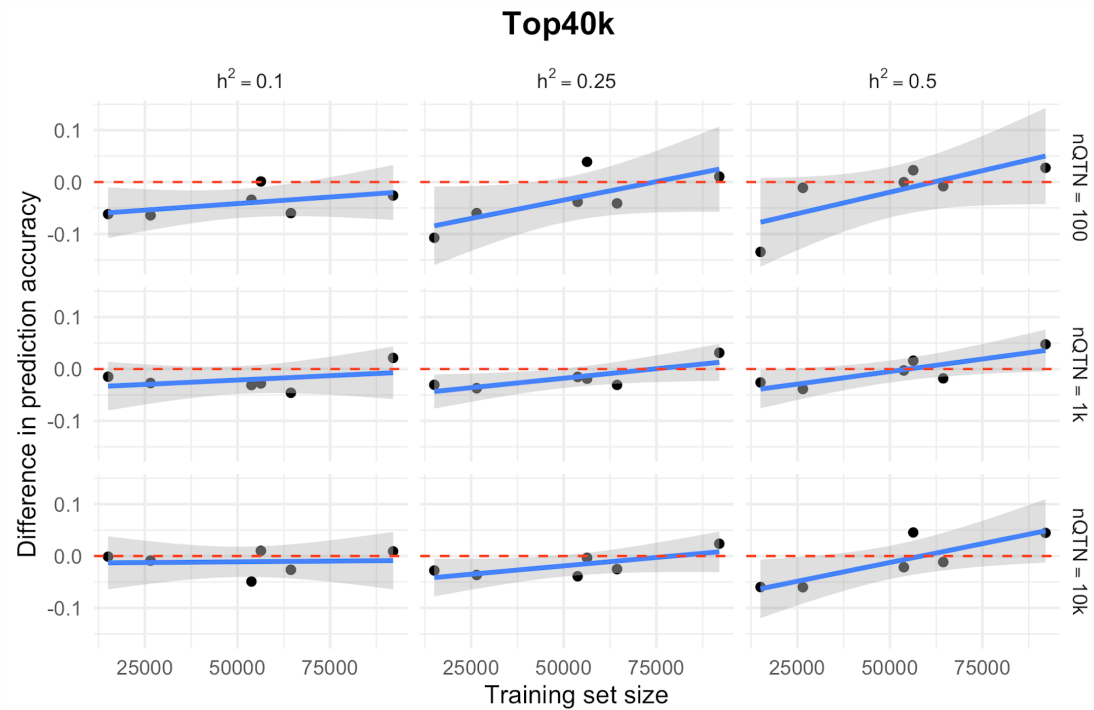
**Figure 3.** Genomic prediction accuracy with the ChipPlusSign predictor variants for the real traits. The difference of prediction accuracy between ChipPlusSign and Chip is shown, for all traits and lines (left) or by trait (right). Red dashed line at 'no difference'.

979
980 **Figure 4.** Genomic prediction accuracy with the Top40k predictor variants for the

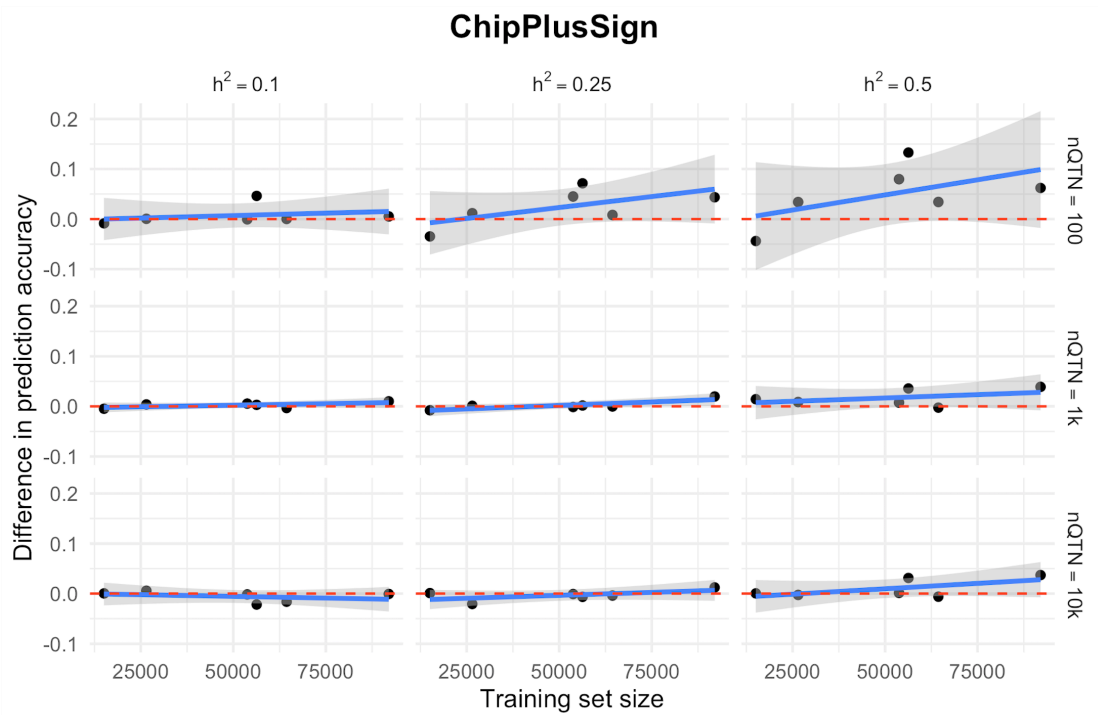981 simulated traits. The difference of prediction accuracy between Top40k and Chip is

982 shown by heritability ($h^2$) and number of quantitative trait nucleotides (nQTN) of the

983 simulated traits. Red dashed line at 'no difference'.

984

985

986
987 **Figure 5.** Genomic prediction accuracy with the ChipPlusSign predictor variants for

988 the simulated traits. The difference of prediction accuracy between ChipPlusSign and

989 Chip is shown by heritability ($h^2$) and number of quantitative trait nucleotides (nQTN)

990 of the simulated traits. Red dashed line at 'no difference'.

991

992

**Figure 6.** Genomic prediction accuracy with the ML-ChipPlusSign predictor variants for the real traits. The difference of prediction accuracy between ML-ChipPlusSign and ML-Chip is shown, for all traits and lines (left) or by trait (right). Red dashed line at 'no difference'.
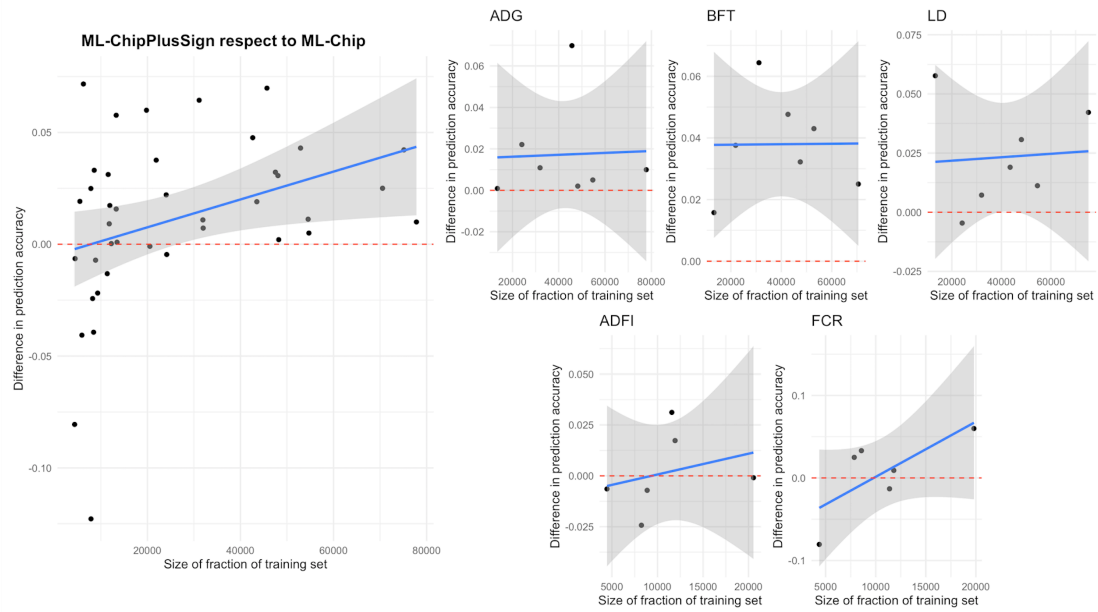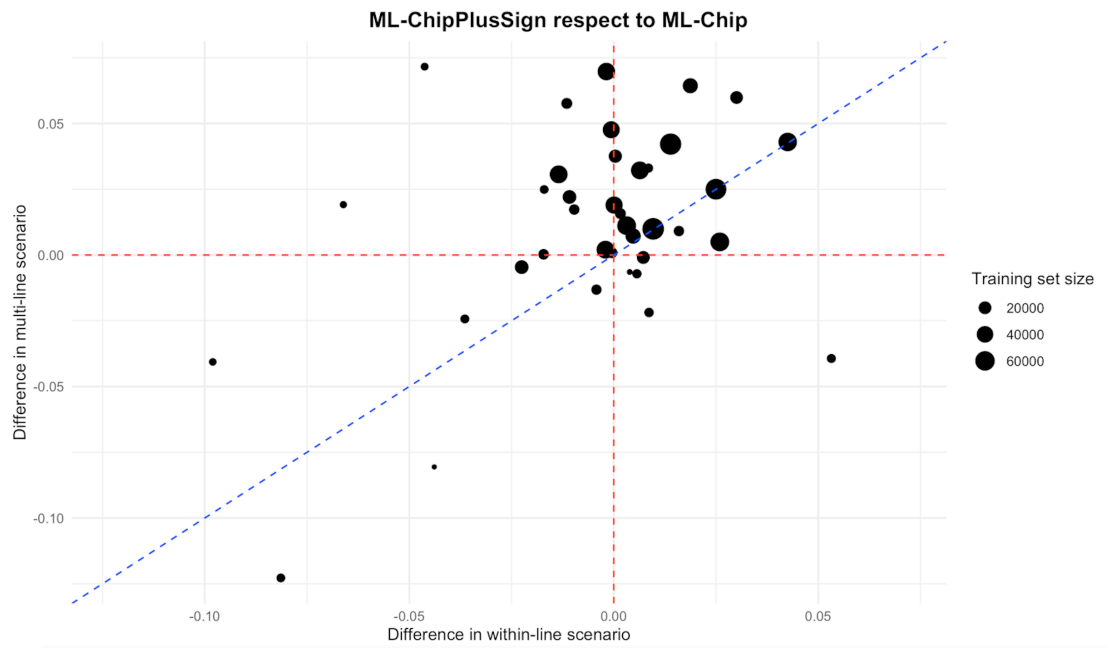
**Figure 7.** Comparison of the difference in genomic prediction accuracy in the multi-line scenarios (between ML-ChipPlusSign and ML-Chip) and in the within-line scenarios (between ChipPlusSign and Chip). Red dashed line at 'no difference'. Blue dashed line is the bisector.

**Figure 8.** Genome-wide association study results for trait BFT in line A. Only chromosome 1 is displayed as an example. In red, results for the variants in the marker array (Chip); in black, results for the whole-genome sequence data (WGS). The blue dashed line indicates significance threshold with Bonferroni's multiple test correction assuming that the markers from the marker arrays were independent (p-value $\leq 10^{-6}$).

## Tables

1015  **Table 1.** Number of sequenced pigs and pigs with imputed data.

| Line | Individuals sequenced | Individuals sequenced by coverage | | | | Individuals used in analyses | |
|---|---|---|---|---|---|---|---|
| | | 1x | 2x | 5x | 15–30x | Pedigree | Imputed |
| A | 1,856 | 1,044 | 649 | 73 | 90 | 122,753 | 104,661 |
| B | 1,366 | 685 | 545 | 44 | 92 | 88,964 | 76,230 |
| C | 1,491 | 628 | 728 | 54 | 81 | 84,420 | 66,608 |
| D | 731 | 362 | 311 | 16 | 42 | 79,981 | 60,474 |
| E | 760 | 394 | 274 | 27 | 65 | 50,797 | 41,573 |
| F | 381 | 193 | 137 | 16 | 35 | 35,309 | 29,330 |
| G | 445 | 217 | 176 | 15 | 37 | 21,129 | 17,224 |

1016

1017 **Table 2.** Number of phenotypic records per trait and line.

| Trait | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| ADG | 88,342 | 64,285 | 56,173 | 51,061 | 35,423 | 26,335 | 15,452 |
| BFT | 80,146 | 62,027 | 55,233 | 47,509 | 34,527 | 23,872 | 15,268 |
| LD | 85,233 | 64,141 | 56,026 | 48,509 | 35,495 | 26,453 | 15,274 |
| ADFI | 21,960 | 9,525 | 9,062 | 12,256 | 12,444 | 4,105* | 4,851 |
| FCR | 21,200 | 9,217 | 8,654 | 12,044 | 12,316 | 4,016* | 4,754 |
| TNB | 13,581 | 10,721 | 9,626 | 7,729* | 6,506* | - | 3,230* |
| LWW | - | 9,112 | 7,251 | - | - | - | 2,813* |
| RET | - | 6,978 | 6,327 | - | - | - | 1,669* |
| Simulated | 104,661 | 76,230 | 66,608 | 60,474 | 41,573 | 29,330 | 17,224 |

1018   *ADG* average daily gain, *BFT* backfat thickness, *LD* loin depth, *ADFI* average daily

1019   feed intake, *FCR* feed conversion ratio, *TNB* total number of piglets born, *LWW* litter

1020   weight at weaning, *RET* return to oestrus 7 days after weaning.

1021   *Included in multi-line scenarios, but excluded in within-line scenarios because of the

1022   limited size of the testing set.

1023

1024    **Table 3.** Number of significant variants from the whole-genome sequence data that

1025    were added to the marker array in ChipPlusSign.

| Trait | A | B | C | D | E | F | G | Multi-line |
|-------|-----|-----|-----|-----|------|-----|-----|------------|
| ADG | 646 | 581 | 424 | 498 | 279 | 219 | 143 | 4731 |
| BFT | 1083 | 758 | 664 | 518 | 1030 | 218 | 237 | 6149 |
| LD | 633 | 579 | 458 | 518 | 222 | 215 | 43 | 7247 |
| ADFI | 145 | 224 | 169 | 23 | 183 | - | 119 | 767 |
| FCR | 198 | 224 | 162 | 95 | 56 | - | 134 | 1369 |
| TNB | 71 | 117 | 161 | - | - | - | - | 248 |
| LWW | - | 32 | 73 | - | - | - | - | 480 |
| RET | - | 184 | 31 | - | - | - | - | 60 |

1026    *ADG* average daily gain, *BFT* backfat thickness, *LD* loin depth, *ADFI* average daily

1027    feed intake, *FCR* feed conversion ratio, *TNB* total number of piglets born, *LWW* litter

1028    weight at weaning, *RET* return to oestrus 7 days after weaning.

1029

1030 **Table 4.** Number of significantly associated genomic regions in the genome-wide
1031 association study for the simulated phenotypes that contained 0, 1 or 2 or more
1032 quantitative trait nucleotides (QTN).

| $h^2$ | nQTN | Line size | Chip | | Whole-genome sequence | | |
|---|---|---|---|---|---|---|---|
| | | | 0 QTN | 1 QTN | 0 QTN | 1 QTN | ≥2 QTN |
| 0.10 | 100 | 27k | 4 | 1 | 8 | 6 | 0 |
| | | 56k | 11 | 3 | 19 | 19 | 0 |
| | | 92k | 10 | 7 | 44 | 19 | 0 |
| | 1k | 27k | 1 | 0 | 4 | 0 | 1 |
| | | 56k | 1 | 0 | 16 | 3 | 1 |
| | | 92k | 1 | 0 | 283 | 9 | 0 |
| | 10k | 27k | 1 | 0 | 1 | 0 | 0 |
| | | 56k | 0 | 0 | 16 | 2 | 1 |
| | | 92k | 2 | 0 | 186 | 17 | 12 |
| 0.25 | 100 | 27k | 11 | 6 | 26 | 15 | 1 |
| | | 56k | 22 | 8 | 44 | 28 | 3 |
| | | 92k | 20 | 7 | 90 | 34 | 1 |
| | 1k | 27k | 0 | 0 | 8 | 1 | 3 |
| | | 56k | 3 | 0 | 34 | 15 | 6 |
| | | 92k | 6 | 0 | 692 | 49 | 16 |
| | 10k | 27k | 0 | 0 | 2 | 0 | 0 |
| | | 56k | 0 | 0 | 90 | 9 | 22 |
| | | 92k | 4 | 0 | 564 | 56 | 164 |
| 0.50 | 100 | 27k | 18 | 9 | 24 | 24 | 1 |
| | | 56k | 30 | 13 | 116 | 41 | 3 |
| | | 92k | 17 | 9 | 425 | 44 | 1 |
| | 1k | 27k | 6 | 0 | 22 | 9 | 6 |
| | | 56k | 5 | 1 | 238 | 59 | 32 |
| | | 92k | 11 | 1 | 903 | 169 | 120 |
| | 10k | 27k | 0 | 0 | 4 | 0 | 0 |
| | | 56k | 0 | 0 | 360 | 77 | 172 |
| | | 92k | 10 | 0 | 379 | 116 | 508 |

1033