

## The Effect of Training Population Size and Chip Density on Accuracy and Bias of Genomic Predictions in Broiler Chickens

J. Ilska<sup>1</sup>, A. Kranis<sup>1,2</sup>, J. A. Woolliams<sup>1</sup>

<sup>1</sup>The Roslin Institute, R(D)SVS, University of Edinburgh, Scotland, <sup>2</sup>Aviagen Ltd., Edinburgh, United Kingdom

**Abstract:** A large dataset of genotyped broilers (over 23,500 individuals, genotyped for 600k SNPs) with phenotypic records on body weight (BWT), feed intake (FI) and hen house production (HHP) was analysed to test the effect of training population (TRN) size and chip density on the accuracy of breeding value prediction. Using 4 cross-validation scenarios and 7 chip densities, the size of TRN was determined as a limiting factor. With large enough TRN, increasing chip density over 20k SNPs did not bring further increases to accuracy. On the other hand, increasing TRN size improved accuracy across all chip densities, most clearly so for BWT and FI.

**Keywords:** accuracy, chip density, training population size

### Introduction

Accuracy of genomic selection (GS) depends on the heritability of the trait, number of genotyped animals with phenotypic records and number of independent chromosomal segments segregating in the population (Daetwyler et al. (2008)). While theoretical expectations regarding the size of the training (TRN) and validation (TST) sets have been compared to empirical accuracies in species like cattle (Luan et al. (2009)), until recently numbers of chickens genotyped at high density were too low to allow such analyses on real data. Most of the published studies on GS in chickens were limited to less than 5k genotyped individuals per line, genotyped at densities up to around 50k SNPs (Wolc et al. (2010; Chen et al. (2011; Wolc et al. (2011)). This study addresses the importance of size of the TRN population and marker density on accuracy of genomic predictions. The effect of these two parameters was tested in three key traits in broiler production.

### Materials and Methods

**Data** The data used in analysis and provided by Aviagen consisted of 23,583 records from broiler chickens, spread over up to 8 generations, genotyped using Affymetrix 600k SNP chip. The pedigree was limited to genotyped individuals only. The phenotypes were: juvenile body weight (**BWT**), recorded at 35 days of age; feed intake (**FI**) recorded between days 13-35; and hen house

production (**HHP**), the accumulated egg production during the whole laying period.

After the quality control performed on the genotypes in Plink (Purcell et al. (2007)) , 412k (412,692) SNPs remained. To evaluate the effect of chip density on accuracy of prediction, 6 additional chips were created (number in brackets is the actual number of SNPs on the chip): 2k (2,337), 7k (7,606), 19k (19,019), 40k (40,052), 70k (70,612), and 134k (134,924). Starting from chip 2k, each consecutive chip was created by adding a random sample of SNPs to the markers contained in the smaller chip.

**Design** In order to investigate the effect of the TRN size, traits such as BWT and FI were used, because for these traits larger numbers of genotyped animals that also have a phenotype were available. For the purposes of the current study, TST animals were assumed not phenotyped, contrary to the breeding practice, to emulate a trait where selection precedes phenotype recording. A sex limited trait, HHP, was also included in the study. To evaluate the effect of the size of the TRN population on the accuracy of predictions, 4 scenarios were tested. From the population available, TST populations of different sizes were created by progressively masking phenotypes of individuals without offspring in the data (i.e. descendants and siblings of TRN individuals only). The validation set was selected from individuals in such a way, that it accounted for 40%, 30%, 20% and 10% of the total population for scenarios SI, SII, SIII and SIV respectively, with the remaining individuals placed in TRN. No consideration was given to the numbers of generations behind TST individuals, with the only requirement being no progeny records.

**Analysis** Breeding values of TST individuals were predicted using the phenotypes and genotypes of their relatives in TRN population and the relationships they shared. This was achieved by fitting mixed linear models in ACTA software package (Gray et al. (2012)) for genomic predictions (GBLUP). The random effect fitted for all traits was the random effect of the animal. The fixed effect fitted to BWT and FI combined several environmental and husbandry factors, while for HHP fixed effect fitted was limited to hatch week. Relationship matrices were calculated using Van

Raden (2008) Method 2 (**G**, GBLUP). **G** varied across the chips due to the different sets of markers.

Accuracies of predictions were calculated in GenStat (Payne et al. (2009)) as a residual correlation between recorded phenotype and predicted BV, corrected for fixed effects, divided by the square root of heritability (BWT  $h^2=0.4$ , FI  $h^2=0.48$ , HHP  $h^2=0.29$ ). The bias of prediction was calculated in GenStat as the regression coefficient of phenotypes regressed on the predicted BVs, with fixed effects accounted for in the model.

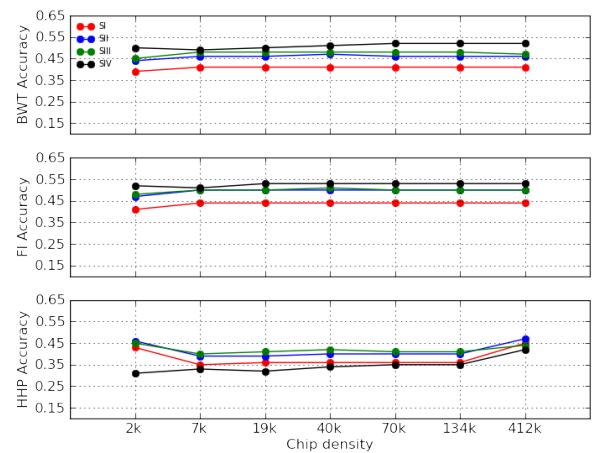
## Results and Discussion

Table 1 gives the numbers of individuals with records in TRN and TST for all scenarios. The size of the TRN population for prediction of GEBVs to achieve high accuracies in excess of 0.7 for unrelated individuals has been previously conjectured as  $2N_eL$ , where  $N_e$  is the effective population size and  $L$  is the size of genome in Morgans (Meuwissen (2009)).  $L$  depends on the recombination rates, which vary greatly across chicken genome (Hillier et al. (2004)). Assuming  $N_e=100$  and  $L=32M$ , (Groenen et al. (2009)),  $2N_eL$  for broiler chickens should exceed 6,400 individuals, which in this study was fulfilled for BWT and FI (see Table 1) but not met for HHP in SIV. However this empirical conjecture was based upon  $h^2=0.8$ . Since accuracy is a function of the product of the size of TRN ( $N$ ) and  $h^2$  (Daetwyler et al. (2008)) this would suggest a TRN size of  $1.6N_eL/h^2$ ; resulting in values of 12,800, 10,300 and 17,100 for BWT, FI and HHP. This is exceeded for all scenarios for BWT, SIII and SIV for FI, and not at all for HHP.

**Table 1. Numbers of individuals with records for the 3 traits and 4 scenarios with varying proportion of individuals assigned to TRN and TST populations.**

Scenario		BWT	FI	HHP
SI	TRN	14 150	7 984	4 167
	TST	9 433	5 744	1 252
SII	TRN	16 508	9 632	4 535
	TST	7 075	4 096	884
SIII	TRN	18 866	10 820	4 828
	TST	4 717	2 908	591
SIV	TRN	21 225	12 276	5 150
	TST	2 358	1 452	269
<b>TOTAL</b>		<b>23 583</b>	<b>13 729</b>	<b>5 420</b>

Figure 1 shows the achieved accuracies, none exceeded this threshold of 0.7. At the highest chip density and largest TRN size considered here, the magnitudes of the genomic accuracies observed ranged from 0.42 for HHP to 0.53 for FI. The rankings of accuracy among the traits increased as expected with  $Nh^2$ , i.e. HHP < BWT < FI. Although numbers were considered to be adequate, in some cases, for accuracies of 0.7, the lower values observed confirm that real data has more complexities than are envisaged in simulation models.



**Figure 1. Accuracy of GEBV prediction using different chip densities and different sizes of the TRN population.**

For BWT and FI the effect of increasing  $N$  was found to have a consistent effect in increasing the accuracy of prediction (Fig 1). For HHP, a similar trend was observed for SI, SII and SIII, however, the accuracy of SIV dropped below the estimates of SI (Fig 1). The number of data points in TST for HHP SIV was relatively small and this lack of precision in estimation is one explanation for this observation.

In contrast increasing chip density had only small effect in BWT and FI, with negligible increases above 19k density. Again results for HHP, with smaller numbers in TRN and TST, were less consistent than for BWT and FI, and the increased standard errors associated with the HHP accuracies make comparisons across scenarios difficult. However, in SIV increasing chip density appeared to lead to increased accuracy.

The regressions of phenotype on GEBV are presented in tables 2, 3 and 4. They demonstrate

that chip density has very small effect on the bias, as with accuracy. The standard errors for the regressions reflect the numbers in the TST set, and the heritability. When using the phenotype to assess accuracy the lower heritability implies a greater environmental noise and will also tend to lead to lower variance of the GEBVs. Therefore the regressions were well estimated for BWT and FI, but had very large standard errors for HHP. The regressions for BWT and FI indicate that the GBLUP model used in this analysis was likely to overestimating differences in true breeding values. For HHP, regressions were almost all consistent with a value of 1, but as noted the standard errors were large.

**Table 2. Estimates of bias (regression coefficient of phenotypes regressed on the predicted BVs) for BWT GEBV using different chip densities and different sizes of TRN population. Standard errors (SE) are given in brackets.**

	SI	SII	SIII	SIV
<b>2k</b>	0.76 (0.03)	0.85 (0.04)	0.83 (0.04)	0.89 (0.06)
<b>7k</b>	0.77 (0.03)	0.84 (0.03)	0.84 (0.04)	0.84 (0.05)
<b>19k</b>	0.77 (0.03)	0.84 (0.03)	0.84 (0.04)	0.85 (0.05)
<b>40k</b>	0.76 (0.03)	0.85 (0.03)	0.84 (0.04)	0.86 (0.05)
<b>70k</b>	0.77 (0.03)	0.85 (0.03)	0.84 (0.04)	0.88 (0.05)
<b>134k</b>	0.76 (0.03)	0.85 (0.03)	0.84 (0.04)	0.88 (0.05)
<b>412k</b>	0.77 (0.03)	0.85 (0.03)	0.83 (0.04)	0.88 (0.05)

**Table 3. Estimates of bias (see Table 2 for details) for FI GEBV using different chip densities and different sizes of TRN population. SE are given in brackets.**

	SI	SII	SIII	SIV
<b>2k</b>	0.72 (0.03)	0.85 (0.04)	0.85 (0.05)	0.89 (0.06)
<b>7k</b>	0.74 (0.03)	0.84 (0.04)	0.81 (0.04)	0.81 (0.06)
<b>19k</b>	0.74 (0.03)	0.83 (0.04)	0.82 (0.04)	0.83 (0.06)
<b>40k</b>	0.73 (0.03)	0.84 (0.04)	0.82 (0.04)	0.84 (0.06)
<b>70k</b>	0.73 (0.03)	0.85 (0.04)	0.82 (0.04)	0.85 (0.06)
<b>134k</b>	0.74 (0.03)	0.86 (0.04)	0.83 (0.04)	0.84 (0.06)
<b>412k</b>	0.74 (0.03)	0.86 (0.04)	0.83 (0.06)	0.85 (0.06)

**Table 4. Estimates of bias (see Table 2 for details) in HHP GEBV using different chip densities and different sizes of TRN population SE are given in brackets.**

	SI	SII	SIII	SIV
<b>2k</b>	1.07 (0.13)	1.11 (0.15)	1.01 (0.17)	0.55 (0.21)
<b>7k</b>	0.88 (0.13)	0.99 (0.16)	0.94 (0.18)	0.61 (0.21)
<b>19k</b>	0.92 (0.13)	0.98 (0.16)	0.98 (0.18)	0.60 (0.21)
<b>40k</b>	0.89 (0.13)	0.98 (0.15)	0.99 (0.18)	0.64 (0.21)
<b>70k</b>	0.94 (0.14)	1.04 (0.16)	1.02 (0.19)	0.68 (0.22)
<b>134k</b>	0.96 (0.14)	1.05 (0.16)	1.03 (0.19)	0.71 (0.23)
<b>412k</b>	1.07 (0.12)	1.05 (0.14)	0.97 (0.17)	0.72 (0.20)

## Conclusions

Analysis of large number of chicken phenotypes and genotypes revealed that increasing the TRN population size is more consistent in increasing the accuracy of genomic selection than increasing chip density. Once the latter exceeds 20k markers (in chickens) the changes observed in accuracy were small although for traits with low numbers of records available, using increased marker density appeared to recover some of the lost information. In contrast, the increases in TRN size continued to increase accuracy for BWT and FI; this effect was less clear in HHP but this trait had a much smaller dataset. Whilst all traits require more records to establish accuracies of 0.7 or more, the accuracies have already reached levels that deliver benefits to the breeding scheme.

## Literature Cited

- Chen, C. Y., Misztal, I., Aguilar, I., et al. (2011). *J Anim Sci*, 89: 23-28.
- Daetwyler, H. D., Villanueva, B. and Woolliams, J. A. (2008). *PLoS One*, 3: e3395.
- Gray, A., Stewart, I. and Tenesa, A. (2012). *Bioinformatics*, 28: 3134-3136.
- Groenen, M. A., Wahlberg, P., Foglio, M., et al. (2009). *Genome Res*, 19: 510-519.
- Hillier, L. W., Miller, W., Birney, E., et al. (2004). *Nature*, 432: 695-716.
- Luan, T., Woolliams, J. A., Lien, S., et al. (2009). *Genetics*, 183: 1119-1126.
- Meuwissen, T. (2009). *Genetics Selection Evolution*, 41: 35.

Payne, R. W., Murray, D. A., Harding, S. A., et al. (2009). <http://www.vsni.co.uk/downloads/genstat/release12/doc/IntroGuide.pdf>, VSN International  
Purcell, S., Neale, B., Todd-Brown, K., et al. (2007). The American Journal of Human Genetics, 81: 559-575.

Wolc, A., Arango, J., Settar, P., et al. (2011). Genetics Selection Evolution, 43: 23.  
Wolc, A., Stricker, C., Arango, J., et al. (2010). Breeding Value Prediction For Production Traits In Layers Using Pedigree and Marker Based Methods. World Congress on Genetics Applied to Livestock Production