

***In-situ* genomic prediction using low-coverage Nanopore sequencing**

Harrison J. Lamb¹, Ben J. Hayes¹, Imtiaz A. S. Randhawa², Loan T. Nguyen¹, Elizabeth M. Ross¹

¹ Centre for Animal Science, Queensland Alliance for Agriculture and Food Innovation, The University of Queensland, Brisbane, QLD 4072, Australia

²School of Veterinary Science, The University of Queensland, QLD, 4343, Australia

Abstract

Most traits in livestock, crops and humans are polygenic, that is, a large number of loci contribute to genetic variation. Effects at these loci lie along a continuum ranging from common low-effect to rare high-effect variants that cumulatively contribute to the overall phenotype. Statistical methods to calculate the effect of these loci have been developed and can be used to predict phenotypes in new individuals. In agriculture, these methods are used to select superior individuals using genomic breeding values; in humans these methods are used to quantitatively measure an individual's disease risk, termed polygenic risk scores. Both fields typically use SNP array genotypes for the analysis. Recently, genotyping-by-sequencing has become popular, due to lower cost and greater genome coverage (including structural variants). Oxford Nanopore Technologies' (ONT) portable sequencers have the potential to combine the benefits genotyping-by-sequencing with portability and decreased turn-around time. This introduces the potential for in-house clinical genetic disease risk screening in humans or calculating genomic breeding values on-farm in agriculture. Here we demonstrate the potential of the later by calculating genomic breeding values for four traits in cattle using low-coverage ONT sequence data and comparing these breeding values to breeding values calculated from SNP arrays. At sequencing coverages between 2X and 4X the correlation between ONT breeding values and SNP array-based breeding values was > 0.92 when imputation was used and > 0.88 when no imputation was used. With an average sequencing coverage of 0.5x the correlation between the two methods was between 0.85 and 0.92 using imputation, depending on the trait. This demonstrates that ONT sequencing has great potential for in clinic or on-farm genomic prediction.

Author Summary

Genomic prediction is a method that uses a large number of genetic markers to predict complex phenotypes in livestock, crops and humans. Currently the techniques we use to determine genotypes requires complex equipment which can only be used in laboratories. However, Oxford Nanopore Technologies' have released a portable DNA sequencer, which can genotype a range of organisms in the field. As a result of the device's higher error rate, it has largely only been considered for specific applications, such as characterising large mutations. Here we demonstrated that despite the devices error rate, accurate genomic prediction is also possible using this portable device. The ability to accurately predict complex phenotypes such as the predisposition to schizophrenia in humans or lifetime fertility in livestock *in-situ* would decrease the turnaround time and ultimately increase the utility of this method in the human clinical and on-farm settings.

Introduction

Complex traits in livestock, crop and human genetics are primarily polygenic, that is, a large number of variants contribute to genetic variation. These traits, which are often continuous (e.g., height [1], weight [2] or temperament [3]) are heritable to varying degrees [4]. The contribution to the overall phenotype of each variant typically ranges between very small effects for common variants to larger effects for some low frequency variants [5]. Each allele an individual carries at loci associated with the complex trait contributes to an increase or decrease in the phenotype [4]. By exploiting this relationship, it is possible to predict the complex phenotype of a genotyped individual without a phenotype if the variant effects are estimated using appropriate statistical methods such as genomic best linear unbiased prediction, Bayesian methods, or polygenic risk scores in humans [6, 7].

In livestock and crops, genomic predictions are termed genomic estimated breeding values (GEBVs) and are used to increase the accuracy and intensity of selection in a process referred to as genomic selection [6]. The benefits of genomic prediction include more accurate selection of complex traits and a decrease in generation interval, which ultimately leads to accelerated genetic gain, evidenced in the poultry [8] and dairy industries [9]. The method allows for the polygenic nature of complex traits by genotyping a number of markers spread evenly throughout the genome and assigning an effect to each marker.

Today, almost all livestock species have SNP arrays designed to enable genomic selection. In Australia, the turnaround time between sample collection and GEBV results is between 6-8 weeks. This turnaround time has prevented the adoption of genomic selection in Australia's northern beef industry, where cattle are generally only handled once a year for a handful of days. This means management decisions based on the GEBV results from SNP array genotyping cannot be implemented until the following year, when the cattle are handled again. There is also significant demand for faster GEBV turnaround time in the sheep and southern beef industries to allow point of management decisions, for example, allocation of feeding regimes based on genetic potential.

In human genetics, complex disease traits have been the focus of this type of quantitative genetics, although pharmacogenetics also holds some potential. For complex diseases such as bipolar disorder, Type II Diabetes, and Crohn's disease [10] polygenic risk scores are used to evaluate an individual's genetic predisposition to the disease. Poly-genic risk scores are based on the genome wide association study (GWAS) correlation between significant SNP markers and the disease [11]. To-date the large-scale utility of polygenic risk scores for individuals has been limited [5]. This has been attributed in part to the current turnaround time, and cost. Instead, family history is used to estimate an individual's predisposition to a particular disease.

However, family history information has limited utility for determining relative risk for complex diseases in some cases. For example, in schizophrenia a family history of the disease is reported in less than a third of cases [12]. A Swedish national study [13] reported family history of the disease in only 3.81% of cases, taking into account first-, second- and third-degree relatives. Moreover, despite half the genetic variance for schizophrenia occurring within family, siblings with a family history of the disease are given the same risk using family history alone. By genotyping the individual, polygenic risk scores far more accurately report an individual's genetic risk than family history alone. In the case of schizophrenia, this information could be used to differentiate between individuals vulnerable to environmental risk factors [14].

Studies exploring the polygenic nature of pharmacological response to chemotherapeutics [15] and asthma treatments [16] have also yielded promising results. Decreasing the turnaround time of genotyping could increase the utility of polygenic risk scores for individuals. For example, a clinician may like to know the genetic risk of a particular disease to assess the priority of diagnostic tests, predict the efficacy of a potential drug treatment, or predict the patient's likelihood of having a severe adverse reaction to a particular treatment, as quickly as possible before the treatment is due to be administered. Decreased turnaround time for polygenic risk scores would also help clinicians diagnose complex diseases in patients while in their prodromal phase, which is vital in early intervention [14].

Genomic prediction in both medicine and agriculture has previously relied heavily on SNP array technology. SNP arrays are a low-cost method to genotype thousands to hundreds-of-thousands of genetic markers. Once costing \$400 USD for 10,000 SNP genotypes, SNP arrays now cost less than \$50

USD for 50,000-800,000 SNP genotypes. However, the technology requires large, expensive laboratory equipment. Therefore, the turnaround time is, at a very optimistic minimum, the time for a sample to reach the laboratory.

Recently, there has been a rise in the popularity of genotyping-by-sequencing. A widely applied approach is to use restriction enzymes to reduce the complexity of genomes for low-coverage sequencing on short-read sequencing platforms [17]. This method is not only becoming cheaper than SNP array genotyping, but it also allows for simultaneous genome wide marker discovery. These advantages have allowed an increased number of traits and associated variants to be studied [4], which has resulted in the ability to accurately predict many new complex traits from genotypes. Still, genotyping-by-sequencing is laboratory based and requires expensive equipment. A potential solution to reduce the turnaround time of genomic prediction while incorporating the benefits of genotyping-by-sequencing, is to use Oxford Nanopore Technologies' (ONT) portable nucleotide sequencer, the MinION [18].

The MinION could be used to sequence samples *in-situ* for genomic prediction in livestock and human settings. Despite initial reports of poor sequencing accuracy and yield, steady developments in flow cell chemistry, library preparation and base calling algorithms has seen reported sequencing yields increase from less than 3 GB to greater than 40 GB and sequencing accuracy increase from 68.4% [19] to over 98% [20, 21]. A major advantage of ONT sequencing technology is the ability for ONT reads to map more accurately to complex genomic regions [22]. This is largely a result of the read length produced by the technology, which has no theoretical limit. This reduces read mapping bias [23] and eliminates the need for restriction enzyme digestion used for genotyping-by-sequencing with short read sequencing technology.

The aim of this study was to evaluate the accuracy of genomic predictions calculated from low-coverage ONT sequence data and compare the results to SNP array-based predictions. We calculated correlations and prediction bias for ONT genomic predictions against the SNP array predictions for various sequencing depths (4x, 2x, 1x, 0.5x) and tested three different methods of imputing missing genotypes. Our results suggest ONT's MinION sequencer could be a useful tool for *in-situ* genomic prediction, with human clinical or on-farm applications.

Results

Sequencing

ONT sequencing of 19 cattle from tail hair yielded an average read length of 1,797 bp and an average flow cell yield of 22.57 Gb over the 96-hr run, with the highest yield being 41.13 Gb (Table 1). The average Phred scaled base quality was 20.54 ± 0.16 with a maximum of 22.3 and minimum of 18.2. On average $86.4\% \pm 0.6$ of reads were effective (i.e., mapping quality Phred score > 0). For each animal, a random subset of sequence data representing 4x, 2x, 1x, and 0.5x sequencing coverage of the 2.7 Gb bovine genome was used for genomic predictions.

Table 1: Read length and flow cell yield for each sample sequenced on the MinION.

| Animal ID | Average Read Length (Bp) | Flow Cell Yield (Gb) | Average base quality ¹ | Percentage of effective reads ² |
|-----------------------|--------------------------|----------------------|-----------------------------------|--|
| DM1 | 879 | 25.87 | 20.5 | 83.2 |
| DM2 | 1,698 | 21.95 | 20.2 | 84.9 |
| DM3 | 1,710 | 19.74 | 20.7 | 86.3 |
| DM4 | 2,004 | 15.98 | 20.8 | 87.3 |
| DM5 | 1,990 | 21.15 | 20.9 | 87.9 |
| DM6 | 1,007 | 26 | 20.2 | 81.3 |
| DM7 | 1,678 | 20.7 | 20.5 | 89.1 |
| DM8 | 2,570 | 13.57 | 21.1 | 89.2 |
| DM9 | 2,456 | 19.22 | 20.1 | 87.4 |
| DM10 | 2,504 | 16.93 | 20.5 | 87.6 |
| DM11 | 1,637 | 23.69 | 21.2 | 85.2 |
| DM12 | 1,941 | 14.9 | 18.2 | 81.2 |
| DM13 | 1,999 | 18.59 | 20.9 | 86.6 |
| DM14 | 1,729 | 34.96 | 20.9 | 87.1 |
| DM15 | 1,847 | 24.83 | 20.7 | 86.9 |
| DM16 | 2,597 | 17.35 | 20.7 | 88.0 |
| DM17 | 1,374 | 21.54 | 20.3 | 85.6 |
| DM18 | 1,132 | 30.75 | 20.2 | 85.6 |
| DM19 | 1,351 | 41.13 | 21.6 | 92.1 |
| Average | 1,795 | 22.57 | 20.54 | 86.4 |
| standard error | 113.67 | 1.56 | 0.16 | 0.6 |

¹ Quality is the sum of all Phred base qualities divided by the number of bases

² $(1 - \frac{(MQ0\ reads + unmapped\ reads)}{total\ raw\ sequences}) * 100$

A minimum allele observation based genotyping approach was used to genotype 641,163 SNP markers. This method accounted for the random probability of sampling alleles in a diploid species by grouping loci with similar coverage together and calling genotypes using a minimum allele count specific to each coverage group. Three methods (non-imputed, imputed-AF & imputed-Beagle) were tested to genotype loci with less than one overlapping ONT read (where only one read was observed

all heterozygous loci were called as homozygous regardless). Although the first two methods are obviously not accurate, they are very fast which could be useful on farm or in a clinical setting. Before imputation, at 0.5x sequencing coverage 64% of ONT genotypes were correctly called, 34% were incorrect in one allele (i.e., called homozygous rather than heterozygous) and 0.02% were incorrect in both alleles (opposing homozygous; Figure 1). At 4x sequencing coverage 74% of calls were correctly called in both alleles, 24% were called incorrectly in one allele and 2% were called incorrectly in both alleles.

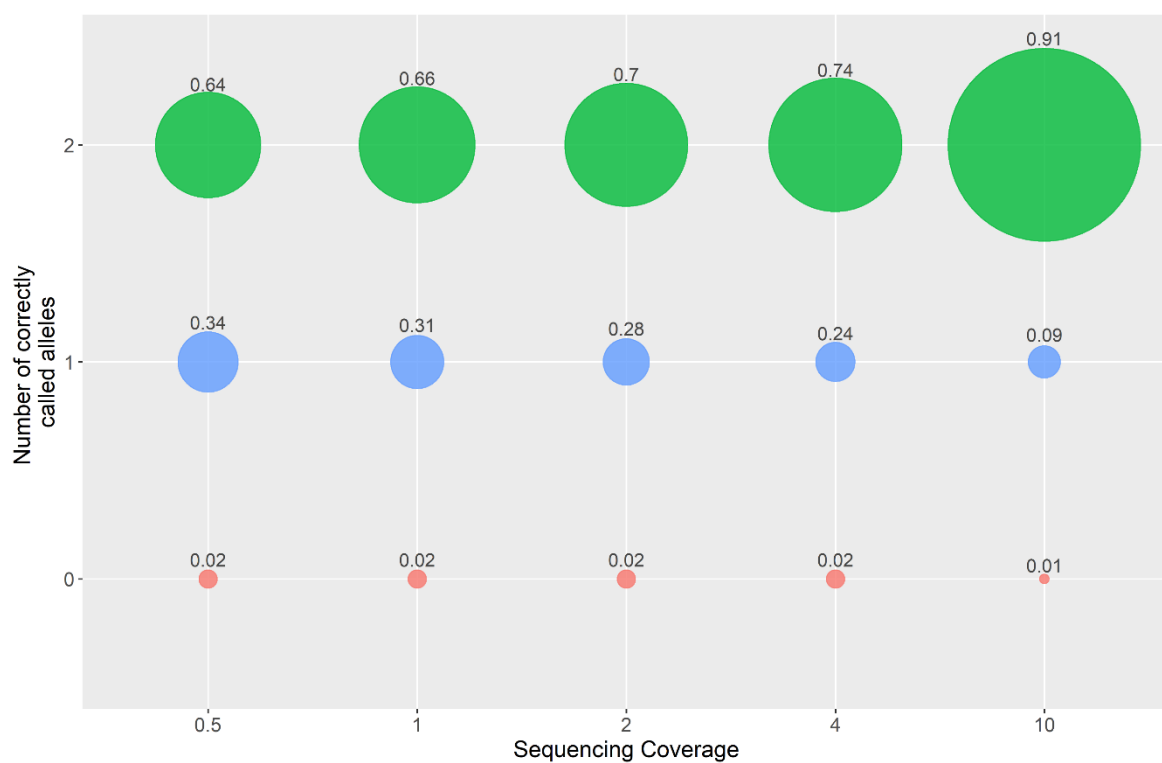


Figure 1: Proportion of genotype calls with both alleles correct (2), one allele incorrect (1) and two alleles incorrect (0) at different sequencing coverages. Two correctly called alleles (green) indicates no difference between the ONT genotype and SNP array genotype. One correctly called allele (blue) indicates one allele of the ONT genotype is correct and the other is incorrect. If both alleles are incorrect the (i.e., the alternate homozygous has been called) the number of correctly called alleles is zero (red).

Marker Effects

Marker effects (BLUP solutions) for predicting GEBV were derived from 26,145 female cattle genotyped for 641,163 SNP markers from Hayes, Fordyce & Landmark [24] for each of the four traits: body condition score, body weight, corpus luteum at 600 days (CL600) and hip height. The marker effects were highly polygenic (Figure 2) with the largest marker effect for each trait less than 0.007% of the total effect. The hip height and body weight marker effects had the largest standard deviations: 0.0011 and 0.00063. While CL600 and body condition score had the smaller standard deviations: 0.00048 and 0.00041, respectively. These SNP marker effects were used to calculate genomic predictions for each of the four traits using ONT genotypes. The effect of read length, base quality, effective mapping percentage and number of reads mapping with quality 0 on ONT genomic prediction accuracy was tested. None of these covariates had a significant effect (linear model; $P > 0.01$) on the prediction accuracy.

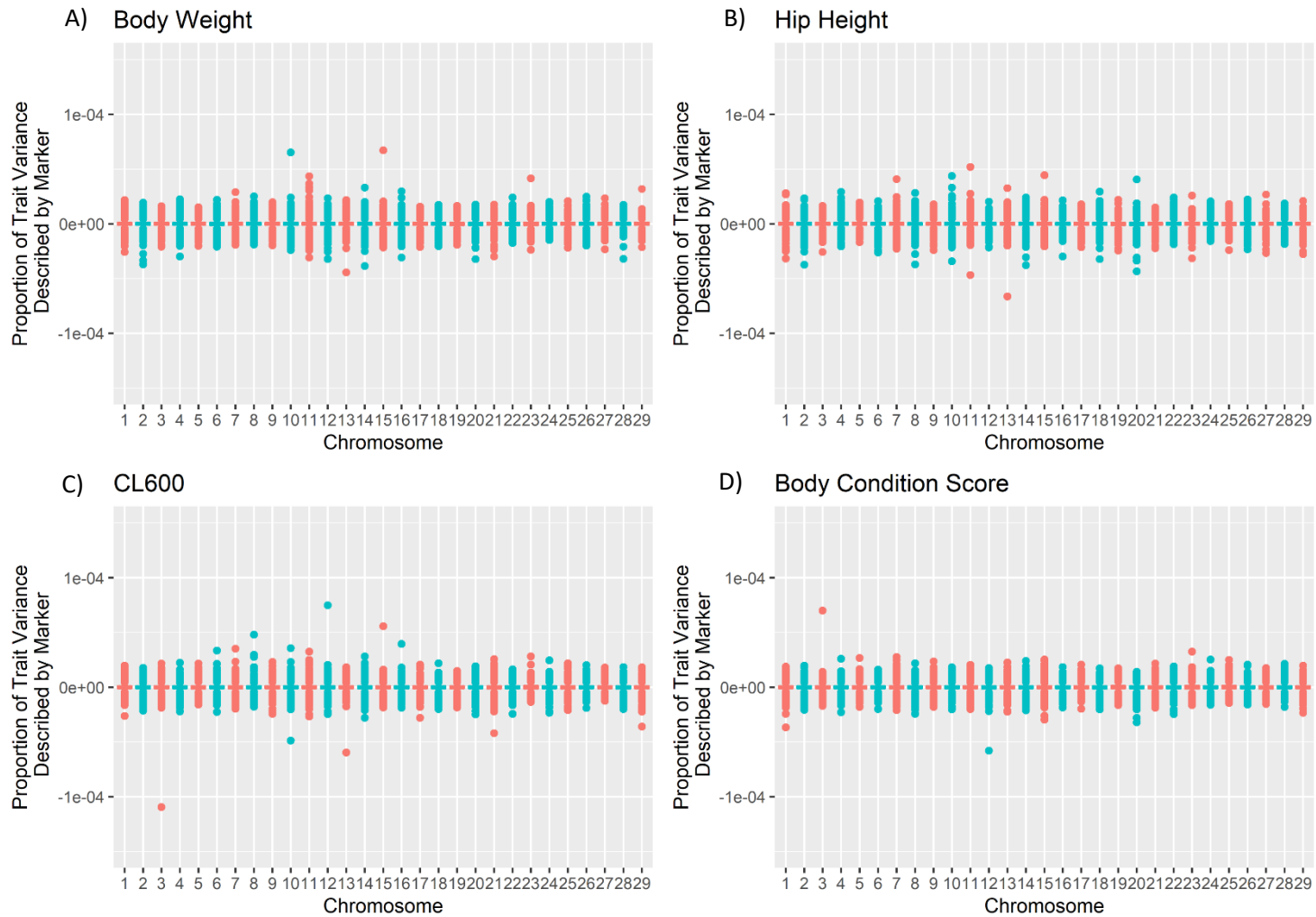


Figure 2: Proportion of trait variance described by the markers across the autosomes for each trait, where a positive value indicates a positive effect on the trait and a negative value a negative effect on the trait. The proportion of trait variance was calculated as $\frac{E_i}{\sum|E|}$ where E_i is the effect of the i_{th} marker and $\sum|E|$ is the sum of the magnitudes of all markers for a particular trait. A) Body weight B) Hip Height C) CL600 D) Body Condition Score.

2 *Genomic prediction accuracies based on non-imputed genotypes*

3 The non-imputed genotyping method, where uncalled loci were assigned a homozygous reference
4 genotype, was as expected the least accurate for genomic prediction across all traits (Figure 3). The
5 ONT non-imputed genomic predictions had reasonable correlations with the SNP array predictions
6 from 4x down to 1x sequencing coverage (0.98 – 0.95 at 4x and 0.95 – 0.85 at 1x). These correlations
7 dropped significantly at 0.5x sequencing coverage, 0.65 for body weight, 0.81 for CL600, 0.72 for body
8 condition score and 0.69 for hip height (Figure 3). For sequencing coverages below 4x the prediction
9 bias (regression coefficient of SNP array GEBV on ONT GEBV) for the non-imputed method was
10 consistently greater than 1, indicating the ONT GEBV are under-estimating compared to the SNP array
11 GEBV. At 2x sequencing coverage the bias ranged between 1.1 for body weight and 1.54 for body
12 condition score (Figure 4; Table 2). At 0.5x sequencing coverage the bias for these two traits increased
13 to 1.45 for body weight and 1.79 for body condition score. For the CL600 trait the bias at 0.5x was 2.7.

14

15 *Genomic prediction accuracies based on imputed-AF genotypes*

16 The correlation coefficients for the imputed-AF method, which used the population allele frequency
17 to genotype missing loci, were the highest of the three methods. At 4x sequencing coverage the
18 correlations were 0.99 for body condition score & CL600, 0.96 for body weight and 0.97 for hip height.
19 This dropped slightly at 2x sequencing coverage, down to 0.95 for body weight, 0.98 for CL600, 0.97
20 for body condition score and 0.95 for hip height. The correlations remained above 0.9 for all traits at
21 1x sequencing coverage except for hip height, which was 0.88. At 0.5x sequencing coverage the
22 correlations for this method remained above 0.8 with the lowest being 0.84 for body condition score.
23 The prediction bias for this method was consistently greater than 1, ranging between 1.65 for hip
24 height and 1.96 for body condition score at 1x sequencing coverage. This increased to between 2.2
25 and 3.1 for hip height and CL600 respectively at the lowest sequencing coverage.

26

27 *Genomic prediction accuracies based on imputed-Beagle genotypes*

28 The final genotyping method, which used the imputation package Beagle v5.1 [25], had correlations
29 greater than 0.85 for all traits at sequencing coverages as low as 0.5x. At 4x sequencing coverage the
30 correlations for body weight, CL600, body condition score and hip height were 0.97, 0.99, 0.98 and
31 0.98 respectively. The correlations decreased slightly with the decrease in sequencing coverage to
32 0.85 for body weight, 0.91 for CL600, 0.89 for body condition score and 0.85 for hip height at 0.5x
33 sequencing coverage. Although the correlations for this genotyping approach were, on average, not
34 as high as the imputed-AF approach, the prediction bias was significantly reduced. The prediction bias
35 at 4x sequencing coverage was around 1 for all traits and decreased slightly as the sequencing
36 coverage decreased. The decrease in prediction bias was most notable in the body weight and hip
37 height traits, which had a bias of 0.88 and 0.74 respectively at 0.5x sequencing coverage. For the same
38 sequencing coverage CL600 and body condition score has a bias of 1.02 and 1.06.

39

40

41

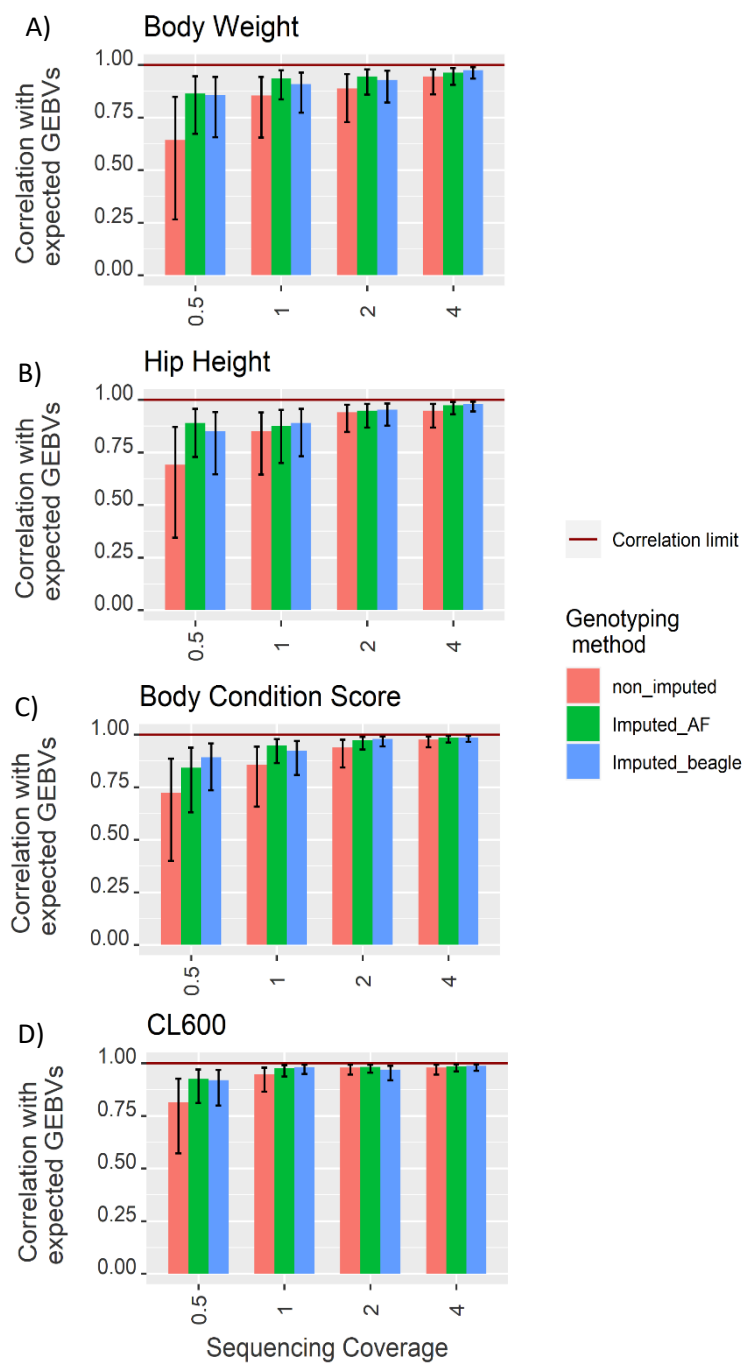


Figure 3: Correlations between the ONT genomic predictions and the expected GEBVs based on SNP array genomic predictions at four sequencing coverages for the four traits; A) body weight, B) hip height, C) body condition score, and D) CL600.

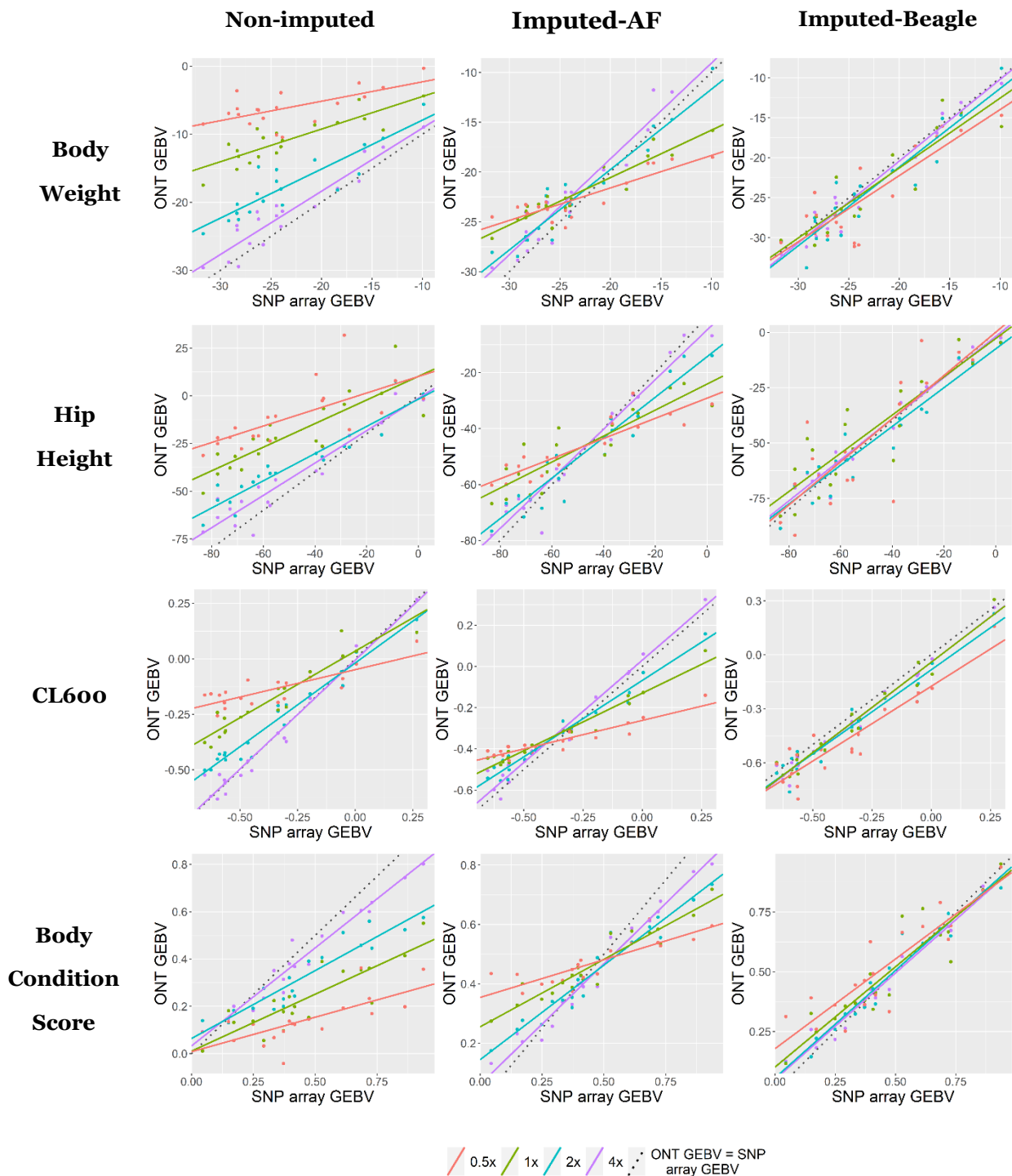


Figure 4: Correlations between the ONT genomic prediction and the SNP array genomic predictions for each of the four traits and three genotyping methods use.

42

43

44

45

46 Table 2: Prediction bias for the different genotyping methods and sequencing
 47 coverages across the four traits.

| | Sequencing Coverage (x of bovine genome) | Bias (Regression coefficient of the SNP array-GEBV on the ONT-GEBV) | | |
|-----------------------------|---|--|------------|----------------|
| | | Non-imputed | Imputed-AF | Imputed-Beagle |
| Body weight | 4 | 0.96 | 0.96 | 0.93 |
| | 2 | 1.1 | 1.11 | 0.88 |
| | 1 | 1.52 | 1.85 | 0.94 |
| | 0.5 | 1.45 | 2.27 | 0.88 |
| CL600 | 4 | 0.98 | 0.98 | 0.96 |
| | 2 | 1.27 | 1.3 | 1.01 |
| | 1 | 1.49 | 1.72 | 0.96 |
| | 0.5 | 2.7 | 3.13 | 1.02 |
| Body condition score | 4 | 1.15 | 1.2 | 1.09 |
| | 2 | 1.54 | 1.49 | 1.06 |
| | 1 | 1.52 | 1.96 | 1.02 |
| | 0.5 | 1.79 | 2.78 | 1.06 |
| Hip height | 4 | 1.06 | 1.08 | 1.03 |
| | 2 | 1.23 | 1.23 | 1.03 |
| | 1 | 1.16 | 1.64 | 0.92 |
| | 0.5 | 1.1 | 2.22 | 0.74 |

48 **Discussion**

49 In this study we sequenced 19 Droughtmaster heifers on ONT's portable MinION sequencer, in order
50 to compare genomic predictions from ONT data to SNP array genotyping, the current standard for
51 genomic prediction. We investigated the accuracy of the ONT genomic predictions at various
52 sequencing coverages (4x, 2x, 1x, and 0.5x) as well as using three different methods to genotype loci
53 with no sequencing coverage: non-imputed, imputed-AF and imputed-Beagle. Beagle [25] imputation
54 produced highly correlated GEBVs across all sequencing coverages. At 0.5x sequencing coverage this
55 method had GEBV correlations greater than 0.85 and there was no evidence of prediction bias. On the
56 other hand, the non-imputed and imputed-AF methods demonstrated significant prediction bias at
57 low sequencing coverages.

58

59 The average sequencing yield of 22.57 Gb represents a significant improvement on flow cell yield from
60 previous ONT studies in cattle [26] and other species [27, 28]. This is largely a result of advances in our
61 ability to fully utilise the sequencing capacity of each flow cell due to the release of a flow cell wash
62 kit (Oxford Nanopore Technologies, Oxford). The wash kit is able to increase the yield from a single
63 flow cell by unblocking 'unavailable' pores [29]. Pores on the flow cell can become blocked during
64 sequencing by contaminants or tertiary DNA structures. To the best of our knowledge, the sequencing
65 yield from one of our samples (41.13 Gb) is one of the largest MinION flow cell yields yet published.
66 The proposed 2022 release of ONT's shoebox sized PromethION P2 [30] could see portable sequencing
67 flow cell yields increase 2-3-fold [31], through the increased density of pores on PromethION flow
68 cells. This increase in the amount and speed of data generation would make *in-situ* genomic prediction
69 more achievable, as the requirement to wait for data acquisition would be drastically reduced.

70

71 The average read length from these tail hair samples was significantly shorter than the average read
72 length from previously sequenced tissues [26, 28]. The average Phred base quality of 20.52 represents
73 a read error-rate of 0.9%, similar quality scores were reported by Runtuwene et al. [32] also using R9.4
74 flow cell chemistry. This shows the improvement in ONT sequencing from the initial reports using the
75 R7.3 flow cell chemistry [32, 33] of base quality scores between 6.88 and 9.4 which relate to an error
76 rate of 20.5% and 11.5% respectively. Further improvements are likely to arise with additional
77 advances in base calling software, which have already shown significant progress [34, 35]. Continued
78 improvements in base call accuracy will lead to more reliable genotypes, and therefore more accurate
79 genomic predictions.

80

81 Although the sequencing effects of read length, base quality, effective mapping percentage and
82 number of reads mapping with quality 0 had no significant effect on the ONT genomic prediction
83 accuracy, this should still be investigated further. It is possible that a lack of variability between
84 samples in this study made it difficult to model the true effect of these factors. Calculating genomic
85 predictions using sequence data from different tissues or library preparation kits could help to
86 investigate these different factors. It is particularly critical to test these effects under field conditions
87 which reflect the real-world variation that may be experienced by either clinicians or producers
88 applying this technology *in-situ*.

89

90 The non-imputed genotyping method performed significantly worse than either of the two imputation
91 methods. This is likely because the missing genotypes were assigned as homozygous reference and
92 therefore received no marker effect, because a homozygous reference is coded as a 0 in the genotype
93 matrix, thereby giving 0 × the alternative allele effect. As the average sequencing coverage decreased
94 the number of missing loci increased and therefore the breeding values were under-estimated at low
95 coverages. This is supported by the decrease in the regression coefficient of ONT GEBV on SNP array

96 GEBV for the non-imputed method as the sequencing coverage decreased. While not unexpected, this
97 phenomenon has important implications: individuals with deeper sequencing coverage will have
98 higher genomic prediction values, as more of their loci will be represented as non-zero in the genotype
99 matrix. Therefore this method is unsuitable for any circumstance where the sequencing coverage is not
100 tightly controlled.

101

102 The allele frequency method of imputation had the highest correlations with the SNP array breeding
103 values across the range of coverages. However, overestimation of breeding values in animals at the
104 lower end for each trait and under estimation of breeding values for animals at the high end was
105 observed at low coverages (Figure 4). This type of prediction bias was also reported by Pimentel, Edel
106 [36] who investigated the effect of imputation error on genomic predictions in 3,494 dairy cows
107 genotyped on a low-density SNP array. They concluded this type of bias was an artifact of imputation
108 algorithms suggesting the most frequent haplotype within a population whenever a haplotype cannot
109 be observed unambiguously. A similar bias was likely seen here because assigning the missing
110 genotypes using allele frequencies (i.e., the average genotype within the population) has regressed
111 the breeding value of each animal toward the population mean. At high coverages (4x and 2x) this is
112 less prominent because there are fewer missing genotypes. However, at the low coverages (1x and
113 0.5x) a significant number of markers are missing ONT genotypes and therefore the regression toward
114 the mean for animals is more prominent. Therefore, this method is likely only useful when within-
115 group rankings are of interest, and not absolute predictions: for example, if a producer wanted to
116 identify the top 30% of animals for a certain trait within their herd.

117

118 The Beagle [25] imputation method for genotyping missing loci was the overall most accurate method
119 when considering both correlation and prediction bias. Although the correlations for some traits were
120 not as high as using the allele frequencies, the prediction bias was consistently close to 1 across all

121 coverages and traits, indicating there was little over or under estimation of breeding values. This is
122 likely because imputation with Beagle [25] uses linkage disequilibrium between genotyped markers
123 to establish haplotypes. This means that the missing genotypes are imputed using the full extent of
124 available information from the genotyped data (i.e., nearby markers genotyped with ONT reads). This
125 approach is very different to imputing the missing genotypes using the allele frequency without
126 considering nearby ONT genotyped markers. Imputation methods such as this are likely the most
127 useful approach, despite the added time required to run the imputation program, which is currently
128 substantially less time than required to produce the sequence data.

129

130 Here we tested the effect of a single imputation program, Beagle [25]. Further increases in the
131 accuracy of ONT genomic predictions could be achieved using purpose built, low-coverage imputation
132 packages such as the recently published program QUILT [37]. This particular imputation package
133 produced similar genotype accuracies as a SNP array imputed using Beagle v.5.1 with 1x ONT
134 sequencing coverage [37]. Additionally, imputation and genomic prediction accuracy could further be
135 increased by genotyping genome wide SNPs rather than exclusively SNPs from the BovineHD BeadChip
136 (Illumina, San Diego, CA) SNP array. This would have the advantage of incorporating the additional
137 information available from whole genome sequencing (a very large number of polymorphic loci) rather
138 than SNP array genotyping.

139

140 Other methods to increase the accuracy of the genomic prediction also include using sequence
141 trimming quality control tools such as Prowler [38]. These tools increase the average Phred base
142 quality of reads which reduces alignment errors and increases the number of effective reads. Studies
143 based on short read information have found that optimisation of imputation strategies can have large
144 effects on the accuracy of imputed genotypes [39, 40]. It is likely that appropriately optimised
145 imputation strategies will increase the accuracy of imputation from ONT data further than what we

146 have reported here and may allow for the calculation of accurate genomic predictions from ultra-low
147 (less than 0.5X) sequence coverage, making the method more cost effective.

148

149 Targeted sequencing methods also hold enormous potential to increase the accuracy of genomic
150 prediction using ONT sequencing. A promising method for *in-situ* applications is ONT's adaptive
151 sequencing, which can target loci within a genome for sequencing without molecular intervention [41,
152 42]. This method only requires the location of genetic markers and no additional laboratory steps,
153 therefore would be ideal for on-farm or clinical applications of genomic prediction with the MinION
154 [41, 42].

155

156 Although each flow cell was run for 96-hrs, the data used here for genomic prediction represent a
157 small subset of the data acquired from each flow cell. For example, given the size of the bovine
158 genome, 1x and 0.5x sequencing coverage represent 3 Gb and 1.5 Gb of sequence data respectively.
159 This means that accurate genomic prediction for each animal would have been achievable in a fraction
160 of the 96-hr sequencing run. Xu, Ge [43] reported an average yield of 0.48 Gb of sequence data per
161 hour on a flow cell. For 0.5x sequencing coverage it would take a little over 3 hours to obtain sufficient
162 data for accurate genomic prediction. With the above-mentioned optimisation steps, it is possible that
163 the minimum required sequence data could fall further, decreasing the sequencing time. ONT also
164 provide rapid library preparation protocols which are capable of preparing sequencing libraries in 15
165 minutes. This could make ONT genomic prediction far more rapid than traditional methods by
166 eliminating the need to send samples to a laboratory. Studies have already demonstrated that ONT
167 sequencing can significantly decrease turnaround times for pathogen identification by providing *in-*
168 *situ* sequencing [44-46], genomic predictions could be the next step in *in-situ* genomic diagnosis.

169

170 Studies have demonstrated the ability for ONT data to characterise structural variants [47, 48], provide
171 rapid pathogen identification [44-46] and assemble both large and small genomes [33, 35, 49]. In
172 cattle, ONT data has successfully characterised the poll allele [26] and a novel structural variant in the
173 *ASIP* gene controlling coat colour in Nellore cattle [50]. More recently in cattle, ONT sequencing was
174 also used to annotate novel transcript isoforms by multiplex sequencing 32 bovine tissues on a single
175 ONT flow cell [51]. Using cattle as an example agricultural species, we have successfully demonstrated
176 yet another application of ONT data. Further optimisation of genotyping-by-sequencing methods with
177 ONT sequence data by combining adaptive sequencing and optimised imputation methods, could see
178 the required coverage for accurate genomic prediction decrease further. Given the average MinION
179 flow cell yield achieved here, more than 20 human/cattle genome sized samples could be multiplexed
180 on a single flow cell, making this approach cost-effective. This study is the first demonstration of
181 genomic predictions using ONT sequencing, which has applications not only in agriculture, but in the
182 clinical setting also.

183

184 **Materials & Methods**

185 *UQ Ethics*

186 Tail hair samples from 19 Droughtmaster animals were collected under the University of Queensland
187 ethics approval numbers SVS/301/18 and SVS/465/18.

188

189 *DNA preparation*

190 DNA for this study was extracted from tail hairs collected from 19 Droughtmaster heifers. The hairs
191 were pulled and stored at room temperature for more than a year prior to DNA extraction. A subset
192 of the hairs was used for genotyping on the 777k BovineHD BeadChip (Illumina, San Diego, CA). The
193 remaining hairs were used for ONT sequencing. Genomic DNA for sequencing was extracted using the

194 Gentra Puregene Tissue Kit (Qiagen) according to the manufacturer's instructions with modifications.
195 Briefly, 20-30 hair samples were lysed in 300 µl of Cell lysis solution (Gentra® Puregene® Tissue Kit)
196 and 1.5 µl of Proteinase K solution (20mg/ml) for 5 hours at 55°C. RNA was then digested by addition
197 of 1.5 µl of RNase A Solution, following 1-hour incubation at 37°C. Samples were placed on ice for 5
198 minutes after adding 100 µl Protein Precipitation Solution (Gentra® Puregene® Tissue Kit) and spun at
199 14000 x g for 3 minutes. 300 µl of Isopropanol was used to precipitate DNA. Samples were centrifuged
200 at 14000 x g for 3 minutes. DNA pellets were washed in 300 µl of 70% ethanol, air-dried for 5 minutes
201 and resuspended in 55 µl of DNA Hydration Solution (Gentra® Puregene® Tissue Kit).

202

203 DNA concentrations were measured using the Qubit dsDNA Broad Range assay kit (Thermo Fisher
204 Scientific). The purity of the extracted DNA was determined with the NanoDropND 1000 (v.3.5.2,
205 Thermo Fisher Scientific), assessing the 260/280 nm and 260/230 nm ratios. The size of extracted DNA
206 was examined using pulsed-field gel electrophoresis (Sage science, USA) with a 0.75% Seakem Gold
207 agarose gel (Lonza, USA) in 0.5X Tris/Borate/EDTA (TBE) running buffer, run for 16 hours at 75 V. The
208 gel was stained after the electrophoresis with SYBR Safe dye (10000x) and visualized using Quantity
209 One analysis software (Bio-rad).

210

211 *Sequencing methodology*

212 Extracted DNA samples were prepared using a ligation kit (SQK-LSK 109, Oxford Nanopore
213 Technologies) based on the manufacturer's instruction with some modifications. Starting with 4 – 8
214 µg of DNA produced enough sequencing library to provide up to 4 flow cell loads from a single library.
215 Samples were diluted at the clean-up points with nuclease-free water to prevent bead clumping during
216 0.4x AmPureXP purifications. End-prep reaction and ligation incubation times were increased to 30
217 minutes and 1 hour respectively. During the 96-hour sequencing run the flow cell was washed three
218 times using the nuclease-flush kit (Oxford Nanopore Technologies) and then reloaded with the same

219 library. The use of DNase in flushing and flow-cell re-fuelling helps to remove blocking DNA and
220 increases the sequencing output.

221

222 *Genotyping methods*

223 Base calls were made from the raw current disruption data using GUPPY (version 4.2.2, Oxford
224 Nanopore Technologies) on the University of Queensland high performance computing infrastructure.
225 The sequence data (fastq files) from each animal were randomly subsampled down to 4x, 2x, 1x & 0.5x
226 sequencing coverage using a bash script. Minimap2 v2.17 [52] was used to align the reads to the *Bos*
227 *taurus* reference genome ARS-UCDv1.2 [53], using the default ONT settings with the flag *-x ont-aln*.
228 Samtools mpileup v1.3 [54] was used to create a pileup of the reads at 641,163 SNP loci that were also
229 genotyped on the 777k BovineHD BeadChip (Illumina, San Diego, CA) and are segregating in Australia's
230 northern beef population.

231

232 A variable allele count genotyping method was used to genotype the SNP loci. Briefly, this method
233 grouped loci with similar coverage together and called genotypes using a separate minimum allele
234 count for each group (Figure 5). The minimum allele counts for each coverage group were chosen such
235 that there was at least a 95% probability of observing both alleles given the random sampling nature
236 of sequencing a diploid species.

237

238

239

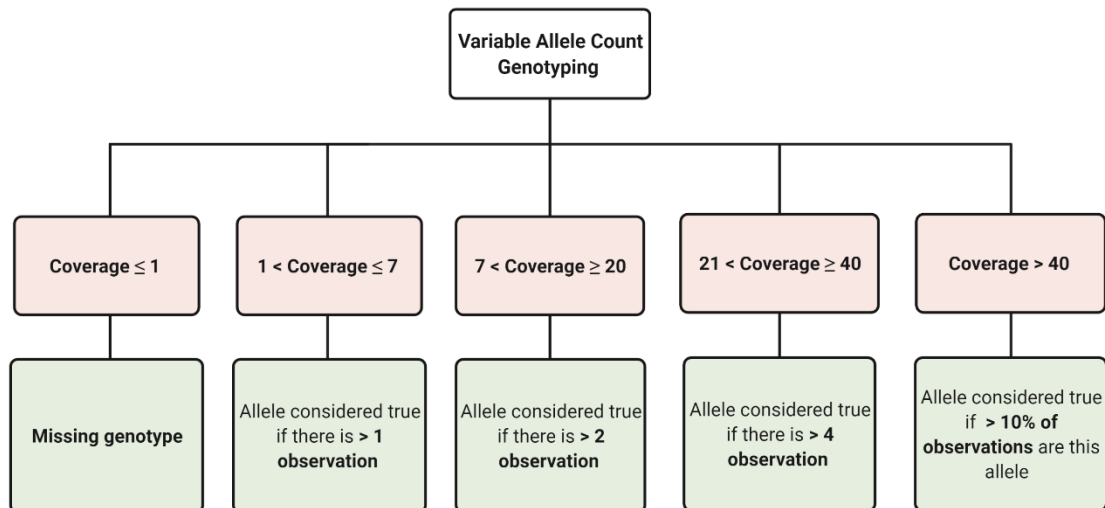


Figure 5: Variable allele count genotyping flow chart

240

241 Three methods were used to assign genotypes at loci with less than 2x coverage: non-imputed,
242 imputed-AF and imputed-Beagle. The non-imputed approach assigned missing loci as homozygous
243 reference, coded as 0 in the genotype matrix. The imputed-AF method used the allele frequency
244 within the population, calculated from SNP array data, to genotype the missing loci. The sum of the
245 probabilities of each genotype given Hardy-Weinberg equilibrium multiplied by the genotype codes,
246 0 for homozygous reference, 1 for heterozygous and 2 or homozygous alternate were used to fill
247 missing genotypes. Using continuous genotype values between 0 and 2 for this method allowed some
248 accounting for the uncertainty of genotyping from low coverage sequence data. Finally, the imputed-
249 Beagle method used the imputation package Beagle v5.1 [25] to impute the missing genotypes given
250 a reference panel of 1,200 animals which represented a subset of animals in the 1000 bull genomes
251 project [54] with high *Bos indicus* content. The effective population and window parameters in Beagle
252 were set to 100,000 and 100 respectively based on Pook, Mayer [56]. For imputation of the 0.5x
253 sequencing coverage genotypes the window parameter was increased to 158 cM to ensure genotyped
254 markers overlapped with the reference panel on each chromosome.

255

256 *SNP-BLUP method for EBVs*

257 Genomic BLUP solutions from [24] were back solved to obtain SNP effects using GCTA [57] at 641,163
258 loci. The SNP effects were for four traits: hip height, body condition score, body weight and CL600
259 (presence of corpus luteum at 600 days). The presence of a corpus luteum is used as an indicator of
260 heifer puberty in beef cattle [58, 59].

261

262 The estimated phenotypes ($\hat{\mathbf{b}}$) for the four traits were calculated using SNP array genotypes, contained
263 in an n by m design matrix \mathbf{M} , where n is number of animals and m is number of markers (641,163).
264 The matrix contained 0, 1 or 2, designating homozygous, heterozygous and homozygous alternative
265 genotype calls from the SNP array. The estimated phenotypes ($\hat{\mathbf{b}}$) were then calculated using:

$$266 \quad \hat{\mathbf{b}} = \mathbf{M}\hat{\mathbf{g}}$$

267

268 The ONT genotypes at the 641,163 loci were then used to calculate genomic breeding values for the
269 four traits. The SNP marker information was contained in an n by m design matrix \mathbf{N} , where n is number
270 of animals genotyped using ONT data, and m is number of markers that were used (641k,
271 corresponding to the same loci as in matrix \mathbf{M}). The matrix contained genotype values between 0 and
272 2, according to the genotyping method used (non-imputed, imputed-AF or imputed-Beagle), and
273 consistent with the nomenclature in matrix \mathbf{M} . The matrix contained only 0, 1 or 2 for the non-imputed
274 and imputed-Beagle methods, however contained continuous genotypes between 0 and 2 for loci
275 imputed using the allele frequency. The estimated phenotypes ($\hat{\mathbf{a}}$) were then calculated using the
276 same estimated marker effects as were used for the SNP array data:

$$277 \quad \hat{\mathbf{a}} = \mathbf{N}\hat{\mathbf{g}}$$

278

279 Linear models were used to evaluate the correlation between \hat{a} and \hat{b} , being the estimated genetic
280 value derived from ONT and SNP array, respectively. The correlation between \hat{a} and \hat{b} was reported
281 as well as the prediction bias (regression coefficient of the SNP array-GEBV on the ONT-GEBV).

282

283 Linear models were also used to evaluate the effect of read length, base quality, effective mapping
284 percentage and number of reads mapping with Phred mapping quality 0 on ONT genomic prediction
285 accuracy. To evaluate each covariate, \mathbf{y} was calculated as the residuals from the model $\hat{b} \sim \hat{a} + e$,
286 where \hat{a} was an array of SNP-array based genomic predictions (from above), and \hat{b} was an array of
287 ONT-based genomics predictions (from above). The significance of each covariate (read length, base
288 quality, effective mapping percentage or number of reads mapping with Phred mapping quality 0) on
289 the genomic prediction accuracy from ONT sequence data was modelled using:

$$290 \quad \mathbf{y} \sim \mathbf{x} + e$$

291 Where \mathbf{x} was an array of the covariate values (read length, base quality, effective mapping percentage
292 or number of reads mapping with Phred mapping quality 0) as a random effect; and \mathbf{y} was the
293 difference between the SNP-based genomic prediction and the ONT-based genomic prediction.

294

295

296 **Acknowledgements**

297 The authors would like to acknowledge the many contributions from friends and colleagues. Firstly,
298 we would like to thank the 54 collaborators in the Northern genomics project, as well as Mr James
299 Copley, Ms Shannon Speight and Dr Geoffrey Fordyce, who were instrumental in its coordination. We
300 would also like to thank Professor Michael McGowan and Dr Russell Lyons for their help in collecting
301 the tail hair samples for sequencing. We are grateful for funding from the MLA Donor Company and
302 University of Queensland through projects P.PSH.P0833, B.STU.2001, L.GEN.1713 and L.GEN.1808.

303 **References**

- 304 1. Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, et al. Defining the role of
305 common variation in the genomic and biological architecture of adult human height. *Nat Genet.*
306 2014;46(11): 1173-86.
- 307 2. Khera AV, Chaffin M, Wade KH, Zahid S, Brancale J, Xia R, et al. Polygenic Prediction of
308 Weight and Obesity Trajectories from Birth to Adulthood. *Cell.* 2019;177(3): 587-96 e9.
- 309 3. Haskell MJ, Simm G, Turner SP. Genetic selection for temperament traits in dairy and beef
310 cattle. *Front Genet.* 2014;5: 368.
- 311 4. Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, et al. 10 Years of GWAS
312 Discovery: Biology, Function, and Translation. *Am J Hum Genet.* 2017;101(1): 5-22.
- 313 5. Torkamani A, Wineinger NE, Topol EJ. The personal and clinical utility of polygenic risk
314 scores. *Nat Rev Genet.* 2018;19(9): 581-90.
- 315 6. Meuwissen THE, Hayes BJ, Goddard ME. Prediction of total genetic value using genome-wide
316 dense marker maps. *Genetics.* 2001;157(4): 1819-29.
- 317 7. Wray NR, Kemper KE, Hayes BJ, Goddard ME, Visscher PM. Complex Trait Prediction from
318 Genome Data: Contrasting EBV in Livestock to PRS in Humans: Genomic Prediction. *Genetics.*
319 2019;211(4): 1131-41.
- 320 8. Wolc A, Kranis A, Arango J, Settar P, Fulton JE, O'Sullivan NP, et al. Implementation of
321 genomic selection in the poultry industry. *Anim Front.* 2016;6(1): 23-31.
- 322 9. Hayes BJ, Bowman PJ, Chamberlain AJ, Goddard ME. Invited review: Genomic selection in
323 dairy cattle: progress and challenges (vol 92, pg 433, 2009). *J Dairy Sci.* 2009;92(3): 1313-.
- 324 10. Evans DM, Visscher PM, Wray NR. Harnessing the information contained within genome-
325 wide association studies to improve individual prediction of complex disease risk. *Hum Mol Genet.*
326 2009;18(18): 3525-31.
- 327 11. Lambert SA, Abraham G, Inouye M. Towards clinical utility of polygenic risk scores. *Hum Mol*
328 *Genet.* 2019;28(R2): R133-R42.
- 329 12. Tandon R, Keshavan MS, Nasrallah HA. Schizophrenia, "just the facts" what we know in
330 2008. 2. Epidemiology and etiology. *Schizophr Res.* 2008;102(1-3): 1-18.
- 331 13. Lichtenstein P, Yip BH, Bjork C, Pawitan Y, Cannon TD, Sullivan PF, et al. Common genetic
332 determinants of schizophrenia and bipolar disorder in Swedish families: a population-based study.
333 *Lancet.* 2009;373(9659): 234-9.
- 334 14. Wray NR, Visscher PM. Narrowing the boundaries of the genetic architecture of
335 schizophrenia. *Schizophr Bull.* 2010;36(1): 14-23.

- 336 15. Chhibber A, Mefford J, Stahl EA, Pendergrass SA, Baldwin RM, Owzar K, et al. Polygenic
337 inheritance of paclitaxel-induced sensory peripheral neuropathy driven by axon outgrowth gene sets
338 in *CALGB 40101* (Alliance). *Pharmacogenomics J.* 2014;14(4): 336-42.
- 339 16. McGeachie MJ, Stahl EA, Himes BE, Pendergrass SA, Lima JJ, Irvin CG, et al. Polygenic
340 heritability estimates in pharmacogenetics: focus on asthma and related phenotypes.
341 *Pharmacogenet Genomics.* 2013;23(6): 324-8.
- 342 17. Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, et al. A Robust, Simple
343 Genotyping-by-Sequencing (GBS) Approach for High Diversity Species. *Plos One.* 2011;6(5).
- 344 18. Jain M, Olsen HE, Paten B, Akeson M. The Oxford Nanopore MinION: delivery of nanopore
345 sequencing to the genomics community. *Genome Biol.* 2016;17(1): 239.
- 346 19. Ashton PM, Nair S, Dallman T, Rubino S, Rabsch W, Mwaigwisya S, et al. MinION nanopore
347 sequencing identifies the position and structure of a bacterial antibiotic resistance island. *Nat*
348 *Biotechnol.* 2015;33(3): 296-300.
- 349 20. Oxford Nanopore announces multiple releases, for high-accuracy, content-rich, high-
350 throughput whole-genome sequencing, and dynamic targeted sequencing [press release]. Oxford
351 Nanopore Technologies - News 28th October 2021 202.
- 352 21. Karst SM, Ziels RM, Kirkegaard RH, Sørensen EA, McDonald D, Zhu Q, et al. Enabling high-
353 accuracy long-read amplicon sequences using unique molecular identifiers with Nanopore or PacBio
354 sequencing. *bioRxiv.* 2020: 645903.
- 355 22. Kono N, Arakawa K. Nanopore sequencing: Review of potential applications in functional
356 genomics. *Dev Growth Differ.* 2019;61(5): 316-26.
- 357 23. Krishnakumar R, Sinha A, Bird SW, Jayamohan H, Edwards HS, Schoeniger JS, et al.
358 Systematic and stochastic influences on the performance of the MinION nanopore sequencer across
359 a range of nucleotide bias. *Sci Rep.* 2018;8(1): 3159.
- 360 24. Hayes BJ, Fordyce G, Landmark S. Genomic predictions for fertility traits in tropical beef
361 cattle from a multi-breed, crossbred and composite reference population. Australian Association for
362 Animal Breeding and Genetics; Armidale 2019.
- 363 25. Browning BL, Browning SR. Genotype Imputation with Millions of Reference Samples. *Am J*
364 *Hum Genet.* 2016;98(1): 116-26.
- 365 26. Lamb HJ, Ross EM, Nguyen LT, Lyons RE, Moore SS, Hayes BJ. Characterization of the poll
366 allele in Brahman cattle using long-read Oxford Nanopore sequencing. *J Anim Sci.* 2020;98(5).
- 367 27. Shin SC, Kim H, Lee JH, Kim HW, Park J, Choi BS, et al. Nanopore sequencing reads improve
368 assembly and gene annotation of the *Parochlus steinenii* genome. *Sci Rep-Uk.* 2019;9.

- 369 28. Ge H, Lin KB, Shen M, Wu SQ, Wang YL, Zhang ZP, et al. De novo assembly of a chromosome-
370 level reference genome of red-spotted grouper (*Epinephelus akaara*) using nanopore sequencing
371 and Hi-C. Mol Ecol Resour. 2019;19(6): 1461-9.
- 372 29. Oxford Nanopore Technologies. New kit extends yields of flow cells 2019 August 18 [cited 09
373 July 2021]. In: Oxford Nanopore Technologies News [Internet]. Available from:
374 <https://nanoporetech.com/about-us/news/new-kit-extends-yields-flow-cells>
- 375 30. Oxford Nanopore Technologies. PromethION P2 2021 May 21 [cited 09 July 2021]. In: Oxford
376 Nanopore Technologies Products [Internet]. Available from: <https://nanoporetech.com/products/p2>
- 377 31. Shafin K, Pesout T, Lorig-Roach R, Haukness M, Olsen HE, Bosworth C, et al. Nanopore
378 sequencing and the Shasta toolkit enable efficient de novo assembly of eleven human genomes. Nat
379 Biotechnol. 2020;38(9): 1044-53.
- 380 32. Runtuwene LR, Tuda JSB, Mongan AE, Makalowski W, Frith MC, Imwong M, et al. Nanopore
381 sequencing of drug-resistance-associated genes in malaria parasites, *Plasmodium falciparum*. Sci
382 Rep. 2018;8(1): 8286.
- 383 33. Schmidt MH-W, Vogel A, Denton AK, Istace B, Wormit A, van de Geest H, et al. De Novo
384 Assembly of a New *Solanum pennellii* Accession Using Nanopore Sequencing. The Plant Cell.
385 2017;29(10): 2336-48.
- 386 34. Silvestre-Ryan J, Holmes I. Pair consensus decoding improves accuracy of neural network
387 basecallers for nanopore sequencing. Genome Biol. 2021;22(1).
- 388 35. Vereecke N, Bokma J, Haesebrouck F, Nauwynck H, Boyen F, Pardon B, et al. High quality
389 genome assemblies of *Mycoplasma bovis* using a taxon-specific Bonito basecaller for MinION and
390 Flongle long-read nanopore sequencing. BMC Bioinformatics. 2020;21(1).
- 391 36. Pimentel ECG, Edel C, Emmerling R, Gotz KU. How imputation errors bias genomic
392 predictions. J Dairy Sci. 2015;98(6): 4131-8.
- 393 37. Davies RW, Kucka M, Su D, Shi S, Flanagan M, Cunniff CM, et al. Rapid genotype imputation
394 from sequence with reference panels. Nat Genet. 2021.
- 395 38. Lee S, Nguyen LT, Hayes BJ, Ross E. Prowler: A novel trimming algorithm for Oxford
396 Nanopore sequence data. bioRxiv. 2021: 2021.05.09.443332.
- 397 39. Wang X, Su GS, Hao D, Lund MS, Kadarmideen HN. Comparisons of improved genomic
398 predictions generated by different imputation methods for genotyping by sequencing data in
399 livestock populations. J Anim Sci Biotechnol. 2020;11(1).
- 400 40. Brouard JS, Boyle B, Ibeagha-Awemu EM, Bissonnette N. Low-depth genotyping-by-
401 sequencing (GBS) in a bovine population: strategies to maximize the selection of high quality
402 genotypes and the accuracy of imputation. BMC Genet. 2017;18.

- 403 41. Kovaka S, Fan YF, Ni BH, Timp W, Schatz MC. Targeted nanopore sequencing by real-time
404 mapping of raw electrical signal with UNCALLED. *Nat Biotechnol.* 2021;39(4): 431-441.
- 405 42. Payne A, Holmes N, Clarke T, Munro R, Debebe BJ, Loose M. Readfish enables targeted
406 nanopore sequencing of gigabase-sized genomes. *Nat Biotechnol.* 2021;39(4): 442-450.
- 407 43. Xu F, Ge C, Luo H, Li S, Wiedmann M, Deng X, et al. Evaluation of real-time nanopore
408 sequencing for *Salmonella* serotype prediction. *Food Microbiol.* 2020;89: 103452.
- 409 44. Taxt AM, Avershina E, Frye SA, Naseer U, Ahmad R. Rapid identification of pathogens,
410 antibiotic resistance genes and plasmids in blood cultures by nanopore sequencing. *Sci Rep.*
411 2020;10(1): 7622.
- 412 45. Li P, Wang K, Qiu S, Lin Y, Xie J, Li J, et al. Rapid identification and metagenomics analysis of
413 the adenovirus type 55 outbreak in Hubei using real-time and high-throughput sequencing
414 platforms. *Infect Genet Evol.* 2021;93: 104939.
- 415 46. O'Donnell VK, Grau FR, Mayr GA, Samayoa TLS, Dodd KA, Barrette RW. Rapid Sequence-
416 Based Characterization of African Swine Fever Virus by Use of the Oxford Nanopore MinION
417 Sequence Sensing Device and a Companion Analysis Software Tool. *J Clin Microbiol.* 2020;58(1).
- 418 47. Norris AL, Workman RE, Fan Y, Eshleman JR, Timp W. Nanopore sequencing detects
419 structural variants in cancer. *Cancer Biol Ther.* 2016;17(3): 246-53.
- 420 48. De Coster W, De Rijk P, De Roeck A, De Pooter T, D'Hert S, Strazisar M, et al. Structural
421 variants identified by Oxford Nanopore PromethION sequencing of the human genome. *Genome*
422 *Res.* 2019;29(7): 1178-87.
- 423 49. Scott AD, Zimin AV, Puiu D, Workman R, Britton M, Zaman S, et al. The giant *sequoia* genome
424 and proliferation of disease resistance genes. *bioRxiv.* 2020: 2020.03.17.995944.
- 425 50. Beatriz Batista Trigo ATHU, Alvaro Fortunato; Marco Milanese, Rafaela Beatriz Pintor
426 Torrecilha, Harrison Lamb, Loan Nguyen, Elizabeth M. Ross, Ben Hayes, Rômulo Cláudio Morozini
427 Padula, Thayla Souza Sussai, Ludmilla Balbo Zavarez, Rafael Silva Cipriano, Maria Margareth
428 Theodoro Caminhas, Flavia Lombardi Lopes, Laiza Helena de Souza lung, Cassiano Pelle, Tosso Leeb,
429 Danika Bannasch, Derek Bickhart, Timothy P L Smith, José Fernando Garcia, Yuri Tani Utsunomiya,
430 Variants spanning *ASIP* contribute to darkness of hair coat in Nellore cattle. *Genet Sel Evol.*
431 2021;53(40).
- 432 51. Halstead MM, Islas-Trejo A, Goszczynski DE, Medrano JF, Zhou H, Ross PJ. Large-Scale
433 Multiplexing Permits Full-Length Transcriptome Annotation of 32 Bovine Tissues From a Single
434 Nanopore Flow Cell. *Front Genet.* 2021;12: 664260.
- 435 52. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34(18):
436 3094-3100.

- 437 53. Rosen BD, Bickhart DM, Schnabel RD, Koren S, Elvik CG, Tseng E, et al. De novo assembly of
438 the cattle reference genome with single-molecule sequencing. *Gigascience*. 2020;9(3).
- 439 54. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
440 Alignment/Map format and SAMtools. *Bioinformatics*. 2009;25(16): 2078-9.
- 441 55. Hayes BJ, Daetwyler HD. 1000 Bull Genomes Project to Map Simple and Complex Genetic
442 Traits in Cattle: Applications and Outcomes. *Annu Rev Anim Biosci*. 2019;7: 89-102.
- 443 56. Pook T, Mayer M, Geibel J, Weigend S, Cavero D, Schoen CC, et al. Improving Imputation
444 Quality in BEAGLE for Crop and Livestock Data. *G3 (Bethesda)*. 2020;10(1): 177-88.
- 445 57. Yang J, Lee SH, Goddard ME, Visscher PM. Genome-wide complex trait analysis (GCTA):
446 methods, data analyses, and interpretations. *Methods Mol Biol*. 2013;1019: 215-36.
- 447 58. Brooks AL, Morrow RE, Youngquist RS. Body-Composition of Beef Heifers at Puberty.
448 *Theriogenology*. 1985;24(2): 235-50.
- 449 59. Johnston DJ, Barwick SA, Corbet NJ, Fordyce G, Holroyd RG, Williams PJ, et al. Genetics of
450 heifer puberty in two tropical beef genotypes in northern Australia and associations with heifer- and
451 steer-production traits. *Anim Prod Sci*. 2009;49(5-6): 399-412.
- 452
- 453