



Study on the concordance between different SNP-genotyping platforms in sheep

H. Marina , P. Chitneedi , R. Pelayo , A. Suárez-Vega , C. Esteban-Blanco ,
B. Gutiérrez-Gil and J. J. Arranz

Departamento de Producción Animal, Facultad de Veterinaria, Universidad de León, Campus de Vegazana s/n, León 24071, Spain.

Summary

Different SNP genotyping technologies are commonly used in multiple studies to perform QTL detection, genotype imputation, and genomic predictions. Therefore, genotyping errors cannot be ignored, as they can reduce the accuracy of different procedures applied in genomic selection, such as genomic imputation, genomic predictions, and false-positive results in genome-wide association studies. Currently, whole-genome resequencing (WGR) also offers the potential for variant calling analysis and high-throughput genotyping. WGR might overshadow array-based genotyping technologies due to the larger amount and precision of the genomic information provided; however, its comparatively higher price per individual still limits its use in larger populations. Thus, the objective of this work was to evaluate the accuracy of the two most popular SNP-chip technologies, namely, Affymetrix and Illumina, for high-throughput genotyping in sheep considering high-coverage WGR datasets as references. Analyses were performed using two reference sheep genome assemblies, the popular Oar_v3.1 reference genome and the latest available version Oar_rambouillet_v1.0. Our results demonstrate that the genotypes from both platforms are suggested to have high concordance rates with the genotypes determined from reference WGR datasets (96.59% and 99.51% for Affymetrix and Illumina technologies, respectively). The concordance results provided in the current study can pinpoint low reproducible markers across multiple platforms used for sheep genotyping data. Comparing results using two reference genome assemblies also informs how genome assembly quality can influence genotype concordance rates among different genotyping platforms. Moreover, we describe an efficient pipeline to test the reliability of markers included in sheep SNP-chip panels against WGR datasets available on public databases. This pipeline may be helpful for discarding low-reliability markers before exploiting genomic information for gene mapping analyses or genomic prediction.

Keywords Affymetrix, genotype concordance, high-throughput genotyping technologies, Illumina, sheep, whole-genome resequencing

Background

High-throughput SNP genotyping platforms have proven to be efficient tools to analyse large populations at an affordable cost to estimate the extension of linkage disequilibrium in livestock genomes (Mokry *et al.* 2014; Chitneedi *et al.* 2017), perform gene mapping studies based on different methodologies, from linkage analysis to genome-

wide association analyses (GWAS; Wu *et al.* 2014; Atlija *et al.* 2016), and allow the practical implementation of genomic selection in many commercial livestock populations, increasing the efficiency of classical breeding (Weller & Ron 2011; Martin *et al.* 2018). For most domestic animal species, the advancements of this technology have significantly increased the number of markers included in the analysed panel, from the available medium-density chips (~50–70 SNP-chips) to the later available high-density panels (~700K SNP-chips). These genomic tools provide a much higher gene mapping accuracy than previous genome scans based on microsatellite markers (Gutiérrez-Gil *et al.* 2009; McClure *et al.* 2013). Hence, since the use of SNP-chips for gene mapping has become routine, many

Address for correspondence

J. J. Arranz, Facultad de Veterinaria, Dpto Producción Animal, Campus de Vegazana, s/n. Universidad de León, 24071 León, Spain.
E-mail: jjarrs@unileon.es

Accepted for publication 28 August 2021

thousands of SNP associations with complex traits have been reported for traits of economic interest in livestock populations (Animal QTL database; Hu *et al.* 2019), some of which have paved the way for the subsequent report of the corresponding causal genetic variants, such as, for example, related to reproductive efficiency and Weaver syndrome in dairy cattle (Adams *et al.* 2016; Kunz *et al.* 2016), to dilated cardiomyopathy in the Doberman pinscher (Meurs *et al.* 2012) or to milk protein percentage in dairy sheep (García-Gómez *et al.* 2012). Regarding the use of SNP-chips in genomic selection programs, high-throughput SNP genotyping platforms also offer the interesting option of exploiting low-density chips (e.g., 3K SNP-chips) in combination with imputation strategies to extend the potential of exploiting genomic information in livestock populations, where the generation of genomic data is difficult to afford.

The most common commercial SNP-chip providers are Illumina (Illumina Inc) and Affymetrix (Affymetrix Inc). Both platforms offer commercial species-specific SNP-chips differing substantially in genotyping technology. Illumina microarray technology uses silica microbeads coated with specific oligos that fit into patterned microwells bearing highly multiplexed SNP genotyping array. Infinium assays are based on a two-colour single base extension from a single hybridisation probe (50-mer) per SNP marker, with allele calls ranging from 3K to over 5 million per sample (Steemers & Gunderson 2007). In contrast, Affymetrix Axiom technology is a two-colour, ligation-based assay using 30-mer oligonucleotide probes that allow simultaneous genotyping of 384 samples with 50K SNPs or 96 individuals with 650K SNPs (Hoffmann *et al.* 2011). The development of specific custom chips that are adapted to new versions of the genome or particular situations, such as the imputation of microsatellites used in determining paternity in specific animal populations, has become increasingly frequent (Nicolazzi *et al.* 2015; Marina *et al.* 2020). Based on these factors, one can easily identify the need to jointly analyse SNP-chip datasets generated with different platforms for different purposes (e.g., meta-analyses in different populations, merging of low- and high-density chips genotyped in the same population).

Additionally, the advancements and affordability of next-generation sequencing techniques, such as from whole-genome resequencing (WGR) datasets, offer a parallel approach to provide, after the corresponding variant calling analysis, SNP genotype datasets that can be used to conduct high-resolution gene mapping or genomic predictions with increased accuracy (de los Campos *et al.* 2013; Sanchez *et al.* 2019). Although there has been a marked decrease in sequencing cost in recent years, SNP genotyping remains the most cost-effective approach when conducting large population genomic studies (Marguerat *et al.* 2008; Bonetta 2010). In this scenario, imputation approaches may infer whole-genome sequence genotypes for a population genotyped with medium- or high-density SNP-chip technology to

increase the accuracy and detection power of different analyses (Al Kalaldehy *et al.* 2019; Sanchez *et al.* 2019; Van Den Berg *et al.* 2019). Therefore, it is crucial to ensure high concordance and reliable reproducibility among the resulting genotypes when merging genotypes generated through different platforms and technologies. This will facilitate the design of meta-analyses based on collaborative projects or publicly available datasets generated with different platforms.

In this context, we should consider that different studies have reported variable amounts of genotyping errors in SNP-chip datasets (Berry *et al.* 2016, 2021; Wu *et al.* 2019) that may impact the results of the subsequent analyses. For example, genotyping errors have been found to decrease the power to detect genuine associations between phenotype and genotype data (Gordon *et al.* 2002; Sun *et al.* 2004). Additionally, the presence of 3% genotyping errors can significantly influence the accuracy in the estimations of linkage disequilibrium extent (Akey *et al.* 2001b), which applies to a wide variety of topics, including disease-gene mapping (Collins *et al.* 1997; Akey *et al.* 2001a), delineation of the demographic history of populations (Laan & Pääbo 1997), and testing of hypotheses of human, cattle, and sheep evolution (Tishkoff & Williams 2002; Kijas *et al.* 2012; Pérez O'Brien *et al.* 2014). All of these factors justify the interest of estimating the average error rate for a given SNP-chip panel to improve the corresponding probe or to remove the problematic markers before proceeding into further analyses.

Therefore, the objective of this study was to compare the genotypes of SNP markers generated for 31 animals through two different medium-density chips (50K SNP-chips) based on the Illumina and Affymetrix array platforms and considering as a quality reference the genotypes determined after the corresponding variant calling bioinformatic analysis of WGR datasets generated for the same animals. The analyses have been made considering two sheep reference genomes: (i) the Oar_v3.1, a popular reference genome for sheep used in many studies since 2012 (https://www.ensembl.org/Ovis_aries_rambouillet/Info/Strains); and (ii) the newest available version of the sheep genome, Oar_rambouillet_v1.0 (https://www.ensembl.org/Ovis_aries_rambouillet/Info/Index). This new assembly has been produced using long-read sequencing technology and has better contiguity (contig N50, LG50) than Oar_v3.1 (Salavati *et al.* 2020). The results of this work will help us: (i) identify the best SNP-chip platform for our future analyses in sheep populations by investigating the differences in the frequency of genotyping errors in the analysed datasets; and (ii) quantify the impact of using different reference genome assemblies on the genotype reliability generated by two SNP-chip platforms. Furthermore, we present a reliable approach based on the minimum amount of WGR information needed to consider the sequence variants as reliable genomic information in this work.

Methods

SNP-chip genotype datasets

All the analyses included in this work included 31 Spanish Churra rams. The DNA of these animals was extracted from semen samples following standard procedures (Sambrook *et al.* 1983). The 31 DNA samples were first genotyped with the commercially available Illumina Ovine SNP50 BeadChip (Illumina Inc) composed of 54 241 SNPs 53 747 SNPs have a known position according to the sheep reference genome Oar_v3.1 (<https://www.ensembl.org/index.html>). Furthermore, 45 444 SNPs of this array have a known position according to the latest sheep reference genome RAMBOUILLET version 1.0 (Oar_rambouillet_v1.0), following Brauning (2019). Raw signal intensities of the Illumina BeadChip array were transformed into genotype calls using GENOME STUDIO software (Illumina Inc).

All the considered DNA samples were also genotyped with a custom 50K Affymetrix Chip that includes 49 702 SNPs (Affymetrix Inc). This Affymetrix custom chip includes 43 406 SNPs and 33 806 SNPs shared with the commercial Illumina SNP-chip considering the Oar_v3.1 and Oar_rambouillet_v1.0 sheep reference genome annotations, respectively. The probe pairs that compose each SNP of the custom 50K SNP-chip were designed based on the sheep reference genome Oar_v3.1. In addition, the corresponding positions of these probes in the latest available sheep reference genome were inferred through their alignment against Oar_rambouillet_v1.0 using the program Burrows–Wheeler Aligner (Li & Durbin 2009). The positions of the probe pairs were accepted when both probes were mapped at the same genomic position, excluding multi-mapped probes. Finally, the SNP positions were verified using the ensemble variation file based on the Oar_rambouillet_v1.0 ovine reference genome release 102 (Yates *et al.* 2020). Following this procedure, a total of 44 456 pair probes were mapped according to the sheep reference genome Oar_rambouillet_v1.0. The raw signal intensities of the Affymetrix chip were transformed into genotype calls through Axiom Analysis Suite software (Applied Biosystems).

Variant calling analysis of WGR datasets

In addition to SNP array genotyping, the 31 DNA sheep samples considered in this work were subjected to WGR using paired-end Illumina sequencing technology (Illumina HiSeq 2000 and HiSeq 2500 sequencers). From the raw sequencing data, genotypes for SNP markers obtained by performing a variant calling analysis, following Marina *et al.* (2020), were considered for further comparison with the genotypes generated by SNP-chip genotyping. The complete variant calling analysis performed has been previously described by our

research group (see Marina *et al.* 2020, for details) and included the following steps: (i) the quality evaluation of the raw paired-end reads was performed with the FASTQC program (Andrews *et al.* 2015); (ii) the poor quality reads were filtered with Trimmomatic (Bolger *et al.* 2014), using filter parameters to paired-end samples (-phred33, LEADING:5, TRAILING:5 SLIDINGWINDOW:4:20, MINLEN:36 ILLUMINACLIP: Trimmomatic-0.33/adapters/TruSeq 3-PE.fa:2:30:10); (iii) sample alignments against the Oar_v.3.1 and Oar_rambouillet_v3.1 ovine reference genome assemblies (<https://www.ensembl.org/index.html>) were performed with the program Burrows–Wheeler Aligner (Li & Durbin 2009) using the algorithm of maximal exact matches (mem); (iv) data manipulation and statistics analysis were performed using SAMTOOLS (Li *et al.* 2009), Picard (Wysoker *et al.* 2019) and Genome Analysis Toolkit_v4.0 (GATK) (McKenna *et al.* 2010); (v) the variant identification was carried out with GATK version_v4.0 (McKenna *et al.* 2010), using the *HaplotypeCaller* tool following GATK Best Practices recommendations; (vi) the low-quality variants were removed with the program SNPSIFT (Cingolani *et al.* 2012) from the variants identified with GATK (DP > 10 & QUAL > 30 & MQ > 30 & QD > 5 & FS < 60); (vii) the BCFtools utilities (Li *et al.* 2009) were used to add the identifier code for each of the known variants through the Ensembl database as the reference; and (viii) the SnpEff program (Cingolani *et al.* 2012) was used to extract genotypes for variants also genotyped through the two considered SNP-chip platforms. Finally, the high-quality WGR variants were subjected to an additional quality control based on the genotype probabilities, which were calculated from allele read counts of the reference (nRef) and the alternative allele (nAlt) following Ros-Freixedes *et al.* (2018). Hence, the probabilities for the reference homozygote (0), heterozygote (1), and alternative homozygote (2) genotypes were calculated as follows:

$$p(0) = (1 - e)^{n_{\text{Ref}}} \cdot e^{n_{\text{Alt}}},$$

$$p(1) = 0.5^{n_{\text{Ref}}} \cdot 0.5^{n_{\text{Alt}}}, \text{ and}$$

$$p(2) = e^{n_{\text{Ref}}} \cdot (1 - e)^{n_{\text{Alt}}}$$

where e is defined as the sequencing error rate, which was assumed to be 0.01 (Ros-Freixedes *et al.* 2018). Last, only those genotypes with at least 99% certainty were considered reliable and were selected for comparison with the SNP-chip-derived datasets.

Genotype concordance analyses

Moreover, the genotypes generated with the two SNP-chip platforms and the WGR technology were converted to PLINK format and standardised to the same strand direction for the three platforms (Purcell *et al.* 2007). To maintain

congruency while comparing the allelic information of markers across SNP-chips and WGR data, all the multiallelic markers were excluded from our analyses. For both Affymetrix and Illumina SNP-chip raw datasets, we performed a quality control (QC) per individual first, keeping those samples with call rates >90%, and second, a QC per marker, considering a marker call rate >95%, a minor allele frequency (MAF) >5%, and genotype frequencies in HWE ($P > 0.05$), as described by Atlija *et al.* (2016). After this filtering, all the shared markers among the three technologies (Affymetrix SNP-chip, Illumina SNP-chip, and WGR) were considered for the subsequent genotype concordance rate calculation. The genotype concordance rate for each SNP (C_i) among the three datasets was estimated through pairwise comparison of the different platforms, considering the raw and QC-filtered data, for each SNP marker, as follows:

$$C_i = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(G_1=G_2)}$$

where n is the number of individuals considered in this study and G_1 and G_2 are the individual genotypes of each SNP marker (i) considered in the pairwise platform comparison. Finally, the global genotype concordance rate (C_g) among the three platforms was calculated as the average overall SNP genotype concordance (C_i).

Results

Raw dataset statistics

For the newest ovine genome reference assembly, Oar_rambouillet_v1.0, the Affymetrix and Illumina platforms included 44 456 SNPs and 45 444 SNPs with known positions, respectively (Fig. 1). The global genotyping call rate (number of non-missing genotypes across all individuals and all SNPs) was higher on the Affymetrix platform (99.35%) than on the Illumina platform (94.92%). However, considering only the 33 806 shared SNPs, the Illumina platform (99.96%) showed a slightly higher call rate than Affymetrix (99.46%).

The 31 considered WGR datasets showed read lengths ranging between 36 and 126 bp and 228 798 583 raw reads per sample on average. The numbers of variants detected with GATK software for the Oar_v3.1 and Oar_rambouillet_v1.0 reference genome assemblies were 44 866 545 and 38 916 540, respectively. Following the QC filtering steps, the numbers of variants retained were 37 674 392 and 36 514 565 for the Oar_v3.1 and Oar_rambouillet_v1.0 assemblies, respectively. The genotypes for these high-quality variants were used for comparison and concordance estimations with each of the two SNP-chip datasets analysed here, considering both reference ovine genome assemblies (Fig. 1).

Concordance rates among raw datasets

A total of 33 806 SNPs (Oar_rambouillet_v1.0) were shared between the two SNP-chips (Fig. 1). After removing the multiallelic variants and performing the certainty quality control based on the number of reads supporting the WGR-based genotypes, the WGR datasets provided a total of 32 493 variants shared with the SNPs genotyped on the Illumina and Affymetrix genotyping platforms (Fig. 1). For these markers and based on all 31 DNA samples, the pairwise comparison estimates for C_g were estimated between the Affymetrix SNP-chip and WGR (97.32%), between the Illumina SNP-chip and WGR (98.33%), and between both SNP-chip platforms (98.07%) (Table 1). The SNP classification based on the concordance among the raw and filtered chip data compared with WGR data is depicted in Table 2, considering the 31 animals and the two sheep reference genomes included in this study (Oar_v3.1 and Oar_rambouillet_v1.0). Globally, considering the 772 598 genotypes compared in the raw dataset of the Oar_rambouillet_v1.0 reference genome assembly, the C_g estimated across the three considered platforms was 96.88%, as represented in Table 2 and Fig. S1. The remaining non-concordant genotypes (3.12%) were classified into three concordance groups (CG): CG1, where all platform genotypes differ; CG2, where only genotypes of the SNP-chip platforms agree; and CG3, where only one of the SNP-chip platforms agrees with the WGR-based genotype, as shown in Table 2 and Fig. S1. For all 32 493 markers shared among the three considered platforms, the Affymetrix and Illumina arrays showed 4099 and 310 missing genotypes, respectively (genotypes with no information). As shown in Table 2, the majority of the genotyping discrepancies between the SNP-chips and WGR (CG2) resulted from single allotyping errors (heterozygotes vs. homozygotes), following by double allotyping errors. In addition, the third group of discordances (CG3), composed of genotypes only concordant between Affymetrix-WGR or Illumina-WGR, highlighted the difference for the C_g estimated from each of the SNP-chip platforms and the WGR dataset, showing a higher concordance for Illumina than for Affymetrix [3.30 (Illumina):1(Affymetrix); Fig. S1], henceforth referenced here as the CG3 ratio (discordance ratio).

When comparing the genotypes from both SNP-chip platforms with the WGR genotypes, the C_i distribution across the genome shows the number and distribution of the SNPs with C_i values lower than 95% (red line) for both SNP-chip platforms (Fig. 2). The C_i values lower than 95% for the pairwise comparison among the three platforms here are represented in Table S1. In addition, the C_i values for SNPs included in the Illumina Ovine SNP50 BeadChip are depicted in Table S2, considering both ovine genome assemblies included in this study. In particular, when considering the genotype concordance rate per chromosome, the value estimated between the Illumina and WGR

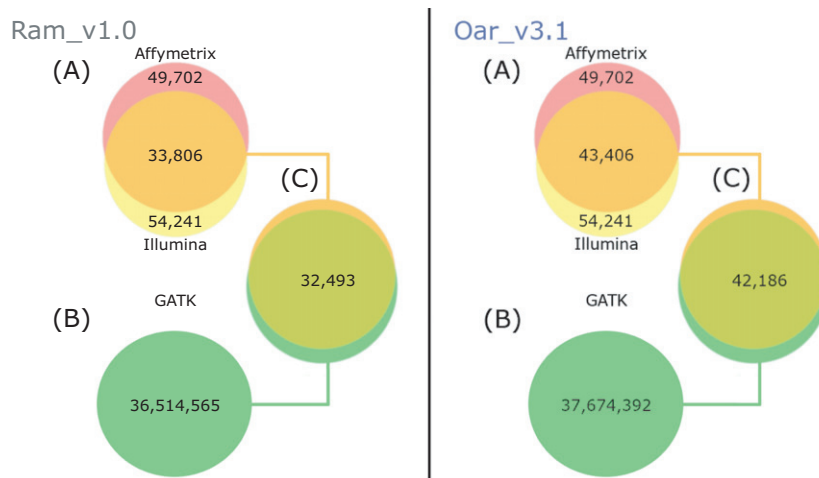


Figure 1 Venn diagrams. The Venn diagrams represent the shared and not shared SNPs among the three reference platforms, considering both ovine reference genomes included in this study (Oar_v3.1 and Oar_rambouillet_v1.0). (a) Venn diagram showing the shared and not shared SNPs in raw data of Affymetrix chip and Illumina chip. (b) Venn diagram showing high-quality variants obtained through the whole genome resequencing (WGR) pipeline. (c) The number of variants shared between SNP-chips and WGR data after removing the multiallelic variants from WGR data.

Table 1 Global genotype concordance rates.

	Before the QC filtering (%)	After the QC filtering (%)
Affymetrix-Illumina	98.07	98.87
Affymetrix-WGR	97.32	98.01
Illumina-WGR	98.33	98.39

This figure summarises the global genotype concordance rate (C_g) among the three technologies before and after the quality control (QC) filter was performed on the SNP-chip genotypes, considering the whole-genome resequencing data aligned against the Rambouillet ovine reference genome assembly (Oar_rambouillet_v1.0).

platforms was slightly higher than that obtained between Affymetrix and WGR, as can be seen in Fig. 2. The slightly lower genotype concordance rate on the X chromosome compared to the autosomes might be because the animals considered in this study belong to heterogametic sex. Further studies are required to compare the average genotype concordance rate between the autosomes and the X chromosome in the homogametic sex.

Concordance rates among QC-filtered datasets

The described QC applied on the SNP-chip datasets did not discard any individuals. In contrast, 2294 and 58 markers genotyped on Affymetrix and Illumina, respectively, were eliminated based on a call rate lower than 95%. Likewise, 1.83% and 2.38% of the markers with MAFs lower than 5% and 3.37% and 3.25% of the markers with a P -value < 0.0001 in the HWE test were removed for the Affymetrix and Illumina platforms, respectively. Fig. 3 shows the average and standard deviation of C_i values for Affymetrix-WGR (A) and Illumina-WGR (B) pairwise comparison for all the variants that passed the QC (HQ). To better understand which of the filtering criteria applied in the QC had a major impact on the control of markers with

a C_i value < 1 , we also provide the C_i averages and standard deviations of the SNP markers discarded based on the call rate (GENO), the HWE test, and the MAF. As shown in Fig. 3, most of the markers that passed the QC showed a C_i value equal to 1; nevertheless, the Affymetrix platform showed a higher number of markers with C_i values lower than 1 compared to the Illumina platform (Table S1). The parameter that filtered out most of the markers with low C_i values was a call rate $> 95\%$ (GENO). Increasing the SNP threshold value for call rate, the most useful parameter to discard SNPs with low C_i values, did not improve the number of SNPs discarded with a $C_i < 95\%$. However, this action eliminated markers with a high C_i value (95–100%, data not shown). The HWE and MAF parameters also discarded SNPs with low C_i values, but several remained after the QC procedure.

In summary, after applying a QC filtering to the SNP-chip datasets based on the Oar_rambouillet_v1.0 positions, a total of 29 849 of the remaining markers were shared among the three platforms. Based on these markers, the concordance accuracy estimated among the three platforms was 97.65%. The remaining non-concordant genotypes (2.35%) were categorised into the three CG categories previously defined, as illustrated in Table 2 and Fig. S2. The QC filtering procedure reduced the number of missing genotypes to 1629 (60.26% less than before QC) on the Affymetrix chip and to 42 genotypes (86.45% less than before QC) on the Illumina chip, as seen within the CG3 group (Table 2, Figs. S1 and S2). The QC procedure increased the C_g values between the Affymetrix and WGR platforms by 0.70%, between Illumina and WGR by 0.06%, and between the two SNP-chips (Affymetrix and Illumina) by 0.80% from the initial estimation based on the raw genotypes (Table 1). Furthermore, as we can appreciate in Table 2, the discordance ratio, which quantifies the number of genotypes only concordant between the Affymetrix-WGR or the Illumina-WGR datasets, strongly decreased by 38.43% after the QC procedure for both SNP-chip platforms

Table 2 Results of genotype concordance.

	Oar_rambouillet_v1.0_raw	Oar_rambouillet_v1.0_QC	OAR_v3.1_raw	OAR_v3.1_QC
Total genotypes	772 598	722 336	1 000 378	903 283
Concordant genotypes (%)	748 464 (96.88%)	705 347 (97.65%)	955 535 (95.52%)	878 711 (97.28%)
Non-concordant genotypes (%)	24 134 (3.12%)	16 989 (2.35%)	44 843 (4.48%)	24 572 (2.72%)
(CG1) A ≠ I ≠ WGR	282	107	650	170
A ≠ I ≠ WGR	117	57	303	101
A or I = NA	165	50	347	69
I ≠ WGR and A = NA	116	38	280	54
A ≠ WGR and I = NA	49	12	67	15
(CG2) A = I ≠ WGR	9228	8845	12 788	12 037
(A = I) = NA	6	0	12	2
(A = I) ≠ NA	9222	8845	12 776	12 035
Chip: Hom and WGR: Hom	1070	873	1408	1087
Chip: Hom and WGR: Het	48	34	62	37
Chip: Het and WGR: Hom	8104	7938	11 248	10 854
Chip: Het and WGR: Het	0	0	58	57
(CG3) A ≠ I and (A or I) = WGR	14 624	8037	31 405	12 365
I = WGR	11 226	5388	26 917	9115
A ≠ WGR	7249	3797	17 419	6780
A = NA	3977	1591	9498	2335
A = WGR	3398	2649	4488	3250
I ≠ WGR	3143	2619	4186	3215
I = NA	255	30	302	35

A, Affymetrix platform; Het, heterozygous genotype; Hom, homozygous genotype; I, Illumina platform; NA, missing genotypes were represented as 'NA' (not available); WGR, whole-genome resequencing.

Genotype concordance between the SNP-chip data (raw and QC data) of Affymetrix (A) and Illumina (I) platforms compared with genotypes determined through WGR analysis. The genotype concordance of shared SNPs between genotyping platforms was estimated considering the 31 animals and the two sheep reference genomes (Oar_v3.1 and Oar_rambouillet_v1.0). The non-concordant genotypes were classified as follows: CG1: all platforms yield different genotypes (A ≠ I ≠ WGR); CG2: genotypes of Affymetrix and Illumina were identical but different from WGR (A = I ≠ WGR); and CG3: genotypes were different between both SNPs genotyping platforms (Affymetrix ≠ Illumina), but one of them was coincident with WGR.

[2.03(Illumina):1(Affymetrix)] (Figs. S1 and S2). However, the numbers of genotyping discrepancies among all technologies (CG1) and between both SNP-chip platforms and WGR datasets (CG2) were barely reduced after the QC. In this last group (CG2), most of the identified errors were single allotyping errors (heterozygotes vs. homozygotes; Fig. S2).

How the ovine reference genome assembly affects concordance rates

After updating the positions at the ovine genome reference assembly, from Oar_v3.1 to Oar_rambouillet_v1.0, a total of 5246 SNPs for the Affymetrix platform and 8303 for the Illumina platform were unmapped, which directly decreased the number of markers compared in this study from 42 186 (Oar_v3.1) to 32 493 (Oar_rambouillet_v1.0). Comparing the general concordance results between these two ovine reference genomes, we appreciated a slight increase when the genome reference version was updated (from Oar_v3.1 to Oar_rambouillet_v1.0), especially before SNP QC filtering (1.35%), but also after the SNP QC filtering (0.36%; see Table 2). Additionally, the proportion of markers classified within the CG1 and CG2 non-concordance groups, which

included the total of markers compared between the Oar_v3.1 and Oar_rambouillet_v1.0 reference genomes, was slightly reduced by 0.03% and 0.08%, respectively (Table 2). Particularly, the discordance ratio drastically dropped off when the genome reference version was updated, especially before the SNP QC filtering (from 6.00 to 3.30) and after SNP QC filtering (from 2.80 to 2.03).

Discussion

The importance of having shared SNPs among several genotyping platforms lies in the fact that genotype data from different studies through different platforms are being generated and, in many cases, are publicly available. Furthermore, the interest in performing meta-analyses to increase the statistical power of gene-mapping analyses or the accuracy of genomic predictions makes it necessary in many cases that the genetic information from these different array platforms can be analysed together (Lopes *et al.* 2018). Some studies that compare the genotype discrepancies between the two most relevant platforms used in high-throughput SNP genotyping, Affymetrix and Illumina, have been previously reported in humans (Suarez *et al.* 2005; Mägi *et al.* 2007; Kim *et al.* 2009), cattle (Wu *et al.* 2019),

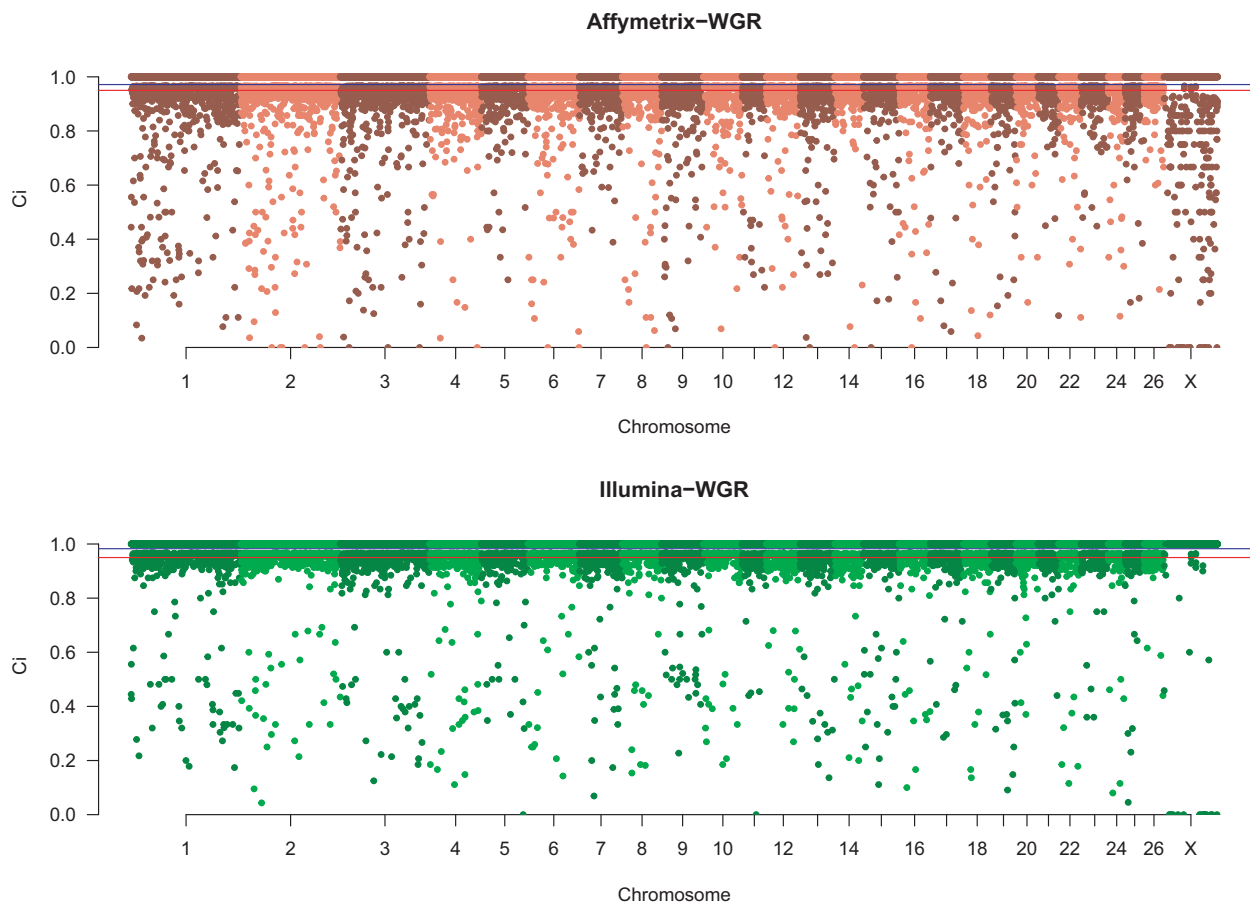


Figure 2 Genotype concordance rate for each SNP (C_i) comparing SNP platforms with whole genome resequencing (WGR). This figure represents the genotype concordance rate for each SNPs (C_i) in both chips platforms (Affymetrix and Illumina) with WGR data aligned against the rambouillet ovine reference genome assembly (Oar_rambouillet_v1.0). The X-axis represents the chromosome to which the SNP belongs, and the Y-axis shows the C_i analysed against WGR. The red lines represent 95% concordance, and the blue lines represent C_g per each SNP-chip platform.

and sheep (Berry *et al.* 2016). As a novel approach focusing on ovine DNA samples, we considered WGR datasets to identify genotyping errors in SNP-chip-derived genotypes. In this way, we think that estimating the genetic information reliability of these two popular SNP-chip technologies is more reliable than when the genotypes from two different SNP-chip platforms are contrasted. Generally, it is difficult to identify genotyping errors on a SNP-chip because two heterozygous parents would be compatible with any observed genotypes due to the diallelic nature of SNP markers (Hinrichs & Suarez 2005). Genotyping errors can reduce the accuracy of imputation and genomic predictions and can also determine false-positive associations in gene-mapping studies by masking the true segregation of alleles (Berry & Kearney 2011; Hong *et al.* 2012a).

Reliability of WGR genotypes

Although we have taken WGR data here as a reference to determine genotyping errors, we acknowledge that WGR

technology is certainly not devoid of sequencing errors. The sequencing depth of WGR datasets is one factor that substantially affects the total coverage of the genome (Sims *et al.* 2014) and can also influence the reliability of the genotypes identified through a variant calling analysis. In this study, all the considered WGR-derived genotypes were subjected to an additional filtering step by considering only those genotypes with at least 99% of quality assurance, as suggested by Ros-Freixedes *et al.* (2018). Higher average sequencing depth requirements help detect more variants and reach higher genome coverage and, henceforth, many reliable genotypes (Taylor *et al.* 2016). The factors responsible for the removal of WGR-derived genotypes during this additional filtering can be attributed to the low quality of the reads during the sequencing procedure, to an insufficient sequencing depth supporting the genotype, or to the location of the variants in genome regulatory regions, which are associated with lower coverage (Wang *et al.* 2011). Thus, comparing genotyping data with a reliable reference (WGR) will help us make more reliable inferences

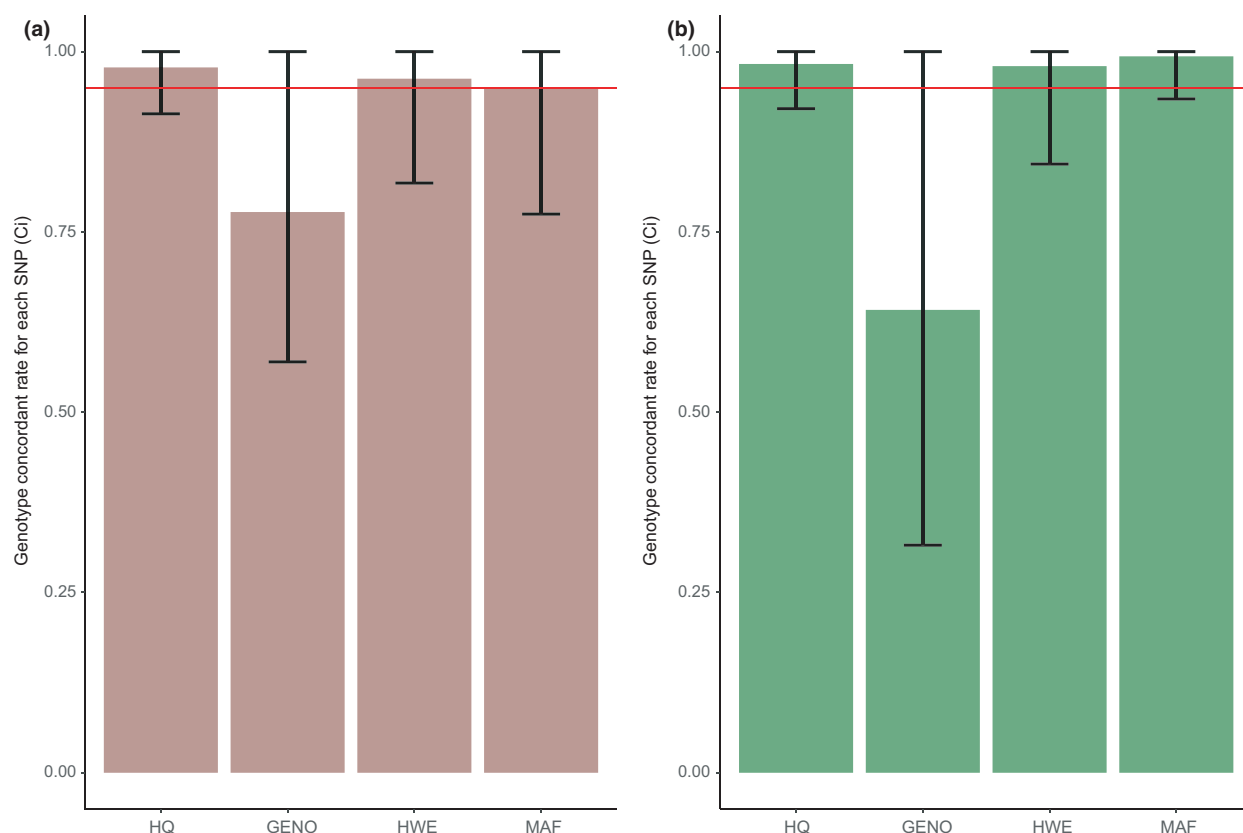


Figure 3 Bar plot of genotype concordance rate for each SNPs (C_i) categorised as concerned by before quality control (QC) filter procedure. The image represents the C_g of all SNPs shared in both SNP-chip platforms [Affymetrix (A) and Illumina (I)] compared to whole genome resequencing (WGR) data, aligned against the rambouillet ovine reference genome assembly (Oar_rambouillet_v1.0). The X-axis represents the four groups into which SNPs have been classified during QC procedure: HQ (high-quality SNPs), GENO (markers filtered by a call rate lower than 95%), HWE [markers filtered because their frequencies differ from the Hardy–Weinberg equilibrium (HWE P -value < 0.05)] and MAF (markers filtered by a minor allele frequency <5%). The Y-axis represents the C_g value and the standard deviation for each group, together with a red line that represents 95% concordance.

and draw robust conclusions regarding the SNP-chip platform that shows the lowest error rate.

Benefits of the SNP-chip QC procedure

The raw call rate estimates of the shared SNPs to the Affymetrix (99.46%) and Illumina (99.96%) SNP-chip platforms were in accordance with the call rate estimates reported in the work of Berry *et al.* (2016), where shared SNPs were also compared between two platforms for ovine DNA samples. In the present study, the C_g estimate for the Oar_rambouillet_v1.0 reference genome was calculated before and after QC filtering of the shared SNPs genotyped on the two SNP-chip platforms. Comparing the C_g before and after the QC helped us understand how removing low-quality SNPs affects the concordance rate and further improves the genotype reproducibility of SNPs across different platforms. The C_g among the three platforms substantially increased (nearly 1%) after implementing QC measures in the SNP-chip dataset (i.e., from 96.88% to

97.65%), as represented in Table 2. The increase in the concordance rate after SNP-chip QC was previously reported in humans (Hong *et al.* 2012b; Jiang *et al.* 2013).

Even though the QC filtering steps carried out over the SNP-chip data do not directly consider the individual genotype quality, we observed that a significant number of erroneous genotypes were removed. Previous studies have recommended SNP QC filtering before applying subsequent association analyses because it reduces the rates of genotyping errors and false-positive results (Zhao *et al.* 2018; Wu *et al.* 2019). As shown in Fig. 3, the parameters included in the QC procedure filtered SNPs with low C_i values. In particular, the genotype call rate filter was the parameter that had a larger impact on increasing the C_g value through the removal of SNPs with genotyping errors. The average C_i of the removed genotypes was 77.08%, which was significantly lower than the concordance rate reported for the two SNP-chip platforms. In the QC applied in this study to both SNP-chip datasets, those SNPs with a call rate <95% were removed, which is a very conservative

value compared to the threshold (85%) recommended by the Affymetrix and Illumina companies. However, the distribution of the SNP C_i values filtered by this threshold (SNP call rate >95%) was considerably lower than the distribution of SNPs that remained after the QC filtering steps (Fig. 3), which supports the need to apply such a conservative threshold to remove SNPs with genotyping errors. SNP genotyping errors are caused by multiple factors, such as errors in the DNA sequence, mutations in the probe complementary regions, errors due to the quality and quantity of the DNA sample, and experimental errors (Hinrichs & Suarez 2005).

The non-concordance rate before and after QC filtering dropped from 3.12% to 2.35% (Table 2, Figs. S1 and S2). A possible explanation for this decrease might be that SNPs with a call rate <95% were eliminated and that the C_i values for the discarded genotypes were lower than the average, as shown in Fig. 3. Within the classification of non-concordant genotypes, three subgroups were categorised. The first group (CG1), where all platforms provided different genotype information, was mainly composed of missing genotypes for one of the SNP-chip platforms, whereas the other SNP-chip platform did not match the WGR-based genotype. This group consisted of 0.48% of the total non-concordant genotypes before QC and 0.34% after QC. A possible explanation for the lack of concordance among these genotypes could be that the region containing these markers would be difficult to genotype or the reference genome is poorly annotated at these specific positions. The second group (CG2) refers to the genotypes generated with the Affymetrix and Illumina arrays that were concordant themselves but different from WGR variants. Following Wu *et al.* (2019), most non-concordant genotypes in this group may be due to a single allotyping error when comparing the SNP-chip genotypes versus the WGR genotypes (i.e., from heterozygous to homozygous). This explanation agrees with our observations as, in our case, multiple allotyping errors in the CG2 group of genotypes were also relatively rare (i.e., between opposing homozygotes). The high reliability of WGR genotypes obtained in this work suggests that the cause of this lack of concordance could be due to non-specific hybridisation of both marker probes, which will reduce the C_i values of the markers (Table S1). Moreover, the differently genotyped SNPs shared between the SNP-chip platforms and WGR could be due to breed-specific mutations in the genome compared to the reference genome or incorrect hybridisation of the probes. Particularly, this second group of non-concordant genotypes (GC2), where the genotyping information agrees between SNP-chip platforms, was not filtered by the QC procedure applied to the SNP-chip datasets, which highlights the importance that comparison of SNP-chip derived genotypes with trustworthy WGR genotypes may have to ensure the reliability of the SNP-chip genotyping information (Berry & Kearney 2011).

Benefits of updating the reference genome version

Comparing the ovine reference genome versions included in this study (Oar_v3.1 and Oar_rambouillet_v1.0), we appreciated that the new assembly increases the genotype concordance among the three platforms (Table 2) and in the pairwise SNP-chip comparison (Table 1). Using a more contiguous reference genome assembly has also reduced the number of non-concordant markers, especially by reducing the discordance ratio between the two SNP-chip platforms (Table 2), which quantifies the number of genotypes only concordant between the Affymetrix-WGR or the Illumina-WGR datasets. The Oar_rambouillet_v1.0 reference genome version, in comparison with the Oar_v3.1 genome, has strongly improved the contiguity (contig N50) [from 40 376 (Oar_v3.1) to 2 572 683 (Oar_rambouillet_v1.0) bp], and significantly reduced: (i) the smallest number of contigs whose length sum makes up half of genome size (contig L50) [from 18 404 to 313 bp]; (ii) the total ungapged length [from 2 534 327 564 to 2 869 531 333 bp]; and (iii) the number of contigs that compose the genome [from 130 764 to 7486] (Agarwala *et al.* 2018). This new assembly also improves the existing genome of the Texel sheep through a higher genomic representation (about 2% more genes represented in the RefSeq annotation) (Liu *et al.* 2016). Therefore, the higher quality of the new assembly has allowed us to discard incorrectly mapped markers in the Oar_v3.1 reference genome version, as described in this study and as can be appreciated in Tables S1 and S2. Accordingly, the upgrade to the Rambouillet reference genome version increases the reliability and concordance of genotypes and corrects the differences between the Affymetrix and Illumina platforms.

Comparison of the reliability of the Affymetrix and Illumina platforms

The C_g estimated between Affymetrix and Illumina before QC filtering was 98.07%, which was in agreement with the genotype concordance of 97.38% previously reported in a multibreed sheep study presented by Berry *et al.* (2016) when also contrasting Affymetrix- and Illumina-generated SNP-chip genotypes. Comparing the C_g between the SNP-chips and the WGR information before and after QC filtering, we found that the Illumina array had a slightly higher number of concordant genotypes per SNP marker than the Affymetrix-based platform (Table 1). As we commented above, the rate of non-concordant genotypes decreased from 3.12% to 2.35% after SNP-chip QC filtering (Table 2, Figs. S1 and S2). Regarding the three subgroups of the classification of non-concordant genotypes, the first and second subgroups (CG1 and CG2) are not helpful for comparing the SNP-chip platforms accurately. However, the third group (CG3), composed of genotypes only concordant between Affymetrix-WGR or Illumina-WGR, may be useful to estimate

the most accurate reliabilities of the SNP-chip platforms. The ratio of genotypes only concordant between Affymetrix-WGR or Illumina-WGR strongly decreased by 38.43% after QC filtering, highlighting the need for performing QC before switching between platforms. Additionally, we found that despite the higher number of missing genotypes removed from Affymetrix after QC filtering, the Illumina platform maintained a slightly higher concordance rate with the WGR genotypes than the Affymetrix platform when the two SNP-chip platforms disagreed, thus providing evidence of the slight advantage of the Illumina SNP-chip to the Affymetrix genotypes (0.4% in filtered data).

Furthermore, the C_i values estimated for both SNP-chip platforms, depicted in Fig. 2, showed that both platforms achieved C_g values >95% (represented by a red line). In addition, the Affymetrix-WGR comparison had more SNPs with C_i values <95% across the genome compared to the Illumina-WGR comparison (Fig. 2 and Table S1), which reduced the C_g values. The genotypes of those markers with a C_i value <95% need to be taken with caution for genomic imputation, genomic prediction, and especially for genome-wide association studies (Tables S1 and S2), as previously recommended by Wu *et al.* (2019).

The C_g values before and after the QC filtering steps among the three platforms indicate that genotyping a population on the Affymetrix or Illumina platforms could influence the outcomes of subsequent analyses, such as genome-wide association studies (Hong *et al.* 2012a) and genomic predictions (Berry & Kearney 2011), as SNP genotyping on both genotyping arrays can contain hundreds to thousands of SNPs with potential errors (Zhao *et al.* 2018). The results from this study will help minimise the differences arising between SNP-chip genotypes and will aid in assessing the reliability of their genotypes. As suggested by Wu *et al.* (2019), we recommend considering the genotyping concordance rate reported here for each SNP (C_i) before carrying out analyses based on SNP-chip datasets. This information can also be interesting when considering the design of future custom SNP-chip arrays, especially to maximise the number of quality markers to be included in a low-density chip that will be used later for imputation to a higher density array (Tables S1 & S2). Moreover, we consider that the approach presented here for genotyping platform comparisons should be suggested as a preliminary step to meta-analyses where genotype datasets from different technologies are merged to gain statistical power. Ensuring that markers with low C_i values are excluded from the meta-analysis will reduce the introduction of errors in subsequent analyses, such as genotype imputation, gene mapping associations and genomic prediction.

Conclusions

When the two different platforms most commonly used for high-throughput SNP genotyping were compared, a slightly

lower concordance rate with WGR data was observed for the Affymetrix platform than for the Illumina platform. The difference in the genotype concordance rate between SNP-chip platforms was reduced after SNP-chip QC filtering because the QC removed low-quality SNP makers, with almost two times more markers filtered in the Affymetrix array than in the Illumina array. The workflow presented here allowed us to identify makers with systematic discordances between SNP-chip platforms and WGR data currently being used to analyse commercial populations in north-west Spain. This list of markers can help avoid their use in subsequent studies to minimise the influence of genotyping errors on the corresponding results. Therefore, we suggest that before performing genomic analyses based on SNP-chip datasets or the manufacturing of a custom SNP-chip, concordance testing of SNP array-derived genotypes with WGR may help to select and relocate markers with low genotype concordance rates to provide an efficient and reliable genomic tool to accomplish guaranteed unbiased, accurate analyses such as GWAS, imputation and genomic prediction. Finally, the comparison of results presented here for the two considered sheep reference genome assemblies offers an opportunity to identify how the use of a new, more complete reference genome can influence the concordance rate of SNP genotypes generated by both SNP-chips or through analysis of WGR datasets.

Acknowledgements

The support and availability of the computing facilities of the Foundation of Supercomputing Center of Castile and León (SCAYLE; <https://www.scayle.es/>) are greatly acknowledged.

Funding

This research work was supported by the RTI2018-093535-B-I00 project funded by the 'Agencia Estatal de Investigación' of the Spanish Ministry of Science and Innovation (Madrid, Spain), co-funded by the European Regional Development Fund. H. Marina is funded by an FPU from the Ministry of Science, Innovation, and Universities (MICIU, Ref. FPU16/01161).

Conflict of interest

The authors declare no conflict of interest.

Authors' contributions

Conceived and designed the experiments: J.J.A. Performed the experiments: H.M. and J.J.A. Designed the analysis: J.J.A. Analysed the data: H.M. Wrote the paper: A.S.V., B.G.G., C.E.B., H.M., J.J.A., P.C. and R.P. All authors read and approved the final manuscript.

Ethical approval

According to the Research Ethics Committee of the University of León, formal ethical approval is not necessary because the animals have been sampled as part of the routine procedures performed on commercial farms.

Data availability statement

Supplementary material related to this article can be found in the online version. The datasets that support the results of this study are available from the authors on a reasonable request.

References

- Adams H.A., Sonstegard T.S., VanRaden P.M., Null D.J., Van Tassell C.P., Larkin D.M. & Lewin H.A. (2016) Identification of a nonsense mutation in APAF1 that is likely causal for a decrease in reproductive efficiency in Holstein dairy cattle. *Journal of Dairy Science* **99**, 6693–701.
- Agarwala R., Barrett T., Beck J. *et al.* (2018) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Research* **46**, D8–13.
- Akey J., Li J. & Xiong M. (2001) Haplotypes vs single marker linkage disequilibrium tests: what do we gain? *European Journal of Human Genetics* **9**, 291–300.
- Akey J.M., Zhang K., Xiong M., Doris P. & Jin L. (2001) The effect that genotyping errors have on the robustness of common linkage-disequilibrium measures. *American Journal of Human Genetics* **68**, 1447–56.
- Al Kalaldehy M., Gibson J., Duijvesteijn N., Daetwyler H.D., Macleod I., Moghaddar N., Lee S.H. & Van Der Werf J.H.J. (2019) Using imputed whole-genome sequence data to improve the accuracy of genomic prediction for parasite resistance in Australian sheep. *Genetics Selection Evolution* **51**, 32.
- Andrews S., Krueger F., Seconds-Pichon A., Biggins F. & Wingett S. (2015) FastQC. A quality control tool for high throughput sequence data. [WWW Document]. Babraham Inst. URL <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/> (accessed 7.3.20).
- Atlíja M., Arranz J.-J., Martínez-Valladares M. & Gutiérrez-Gil B. (2016) Detection and replication of QTL underlying resistance to gastrointestinal nematodes in adult sheep using the ovine 50K SNP array. *Genetics Selection Evolution* **48**, 4.
- Berry D.P., Dunne F.L., Evans R.D., McDermott K. & O'Brien A.C. (2021) Concordance rate in cattle and sheep between genotypes differing in Illumina GenCall quality score. *Animal Genetics* **52**, 208–13.
- Berry D.P. & Kearney J.F. (2011) Imputation of genotypes from low- to high-density genotyping platforms and implications for genomic selection. *Animal* **5**, 1162–9.
- Berry D.P., O'Brien A., Wall E., McDermott K., Randles S., Flynn P., Park S., Grose J., Weld R. & McHugh N. (2016) Inter- and intra-reproducibility of genotypes from sheep technical replicates on Illumina and Affymetrix platforms. *Genetics Selection Evolution* **48**, 86.
- Bolger A.M., Lohse M. & Usadel B. (2014) Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–20.
- Bonetta L. (2010) Whole-genome sequencing breaks the cost barrier. *Cell* **141**, 917–9.
- Brauning R. (2019) Mapping of ISGC SNP chip probes. figshare. Dataset. [WWW Document]. URL <https://doi.org/10.6084/m9.figshare.8424935.v2> (accessed 6.17.21).
- Chitneedi P.K., Arranz J.J., Suarez-Vega A., García-Gámez E. & Gutiérrez-Gil B. (2017) Estimations of linkage disequilibrium, effective population size and ROH-based inbreeding coefficients in Spanish Churra sheep using imputed high-density SNP genotypes. *Animal Genetics* **48**, 436–46.
- Cingolani P., Platts A., Wang L.L., Coon M., Nguyen T., Wang L., Land S.J., Lu X. & Ruden D.M. (2012) A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly* **6**, 80–92.
- Collins F.S., Guyer M.S. & Chakravarti A. (1997) Variations on a theme: cataloging human DNA sequence variation. *Science* **278**, 1580–1.
- de los Campos G., Vazquez A.I., Fernando R., Klimentidis Y.C. & Sorensen D. (2013) Prediction of complex human traits using the genomic best linear unbiased predictor. *PLoS Genetics* **9**, e1003608.
- García-Gámez E., Gutiérrez-Gil B., Sahana G., Sánchez J.-P., Bayón Y. & Arranz J.-J. (2012) GWA analysis for milk production traits in dairy sheep and genetic support for a QTN influencing milk protein percentage in the LALBA gene. *PLoS One* **7**, e47782.
- Gordon D., Finch S.J., Nothnagel M. & Ott J. (2002) Power and sample size calculations for case-control genetic association tests when errors are present: application to single nucleotide polymorphisms. *Human Heredity* **54**, 22–33.
- Gutiérrez-Gil B., Pérez J., Álvarez L., Martínez-Valladares M., de la Fuente L.-F., Bayón Y., Meana A., Primitivo F.S., Rojo-Vázquez F.-A. & Arranz J.-J. (2009) Quantitative trait loci for resistance to trichostrongylid infection in Spanish Churra sheep. *Genetics Selection Evolution* **41**, 46.
- Hinrichs A.L. & Suarez B.K. (2005) Genotyping errors, pedigree errors, and missing data. *Genetic Epidemiology* **29**, S120–4.
- Hoffmann T.J., Kvale M.N., Hesselton S.E. *et al.* (2011) Next generation genome-wide association tool: design and coverage of a high-throughput European-optimized SNP array. *Genomics* **98**, 79–89.
- Hong H., Xu L., Liu J. *et al.* (2012a) Technical reproducibility of genotyping SNP arrays used in genome-wide association studies. *PLoS One* **7**, e44483.
- Hong H., Xu L., Liu J. *et al.* (2012b) Technical reproducibility of genotyping SNP arrays used in genome-wide association studies. *PLoS One* **7**, e44483.
- Hu Z.L., Park C.A. & Reecy J.M. (2019) Building a livestock genetic and genomic information knowledgebase through integrative developments of animal QTLdb and CorrDB. *Nucleic Acids Research* **47**, D701–10.
- Jiang L., Willner D., Danoy P., Xu H. & Brown M.A. (2013) Comparison of the performance of two commercial genome-wide association study genotyping platforms in Han Chinese samples. *G3 Genes/Genomes/Genetics* **3**, 23–9.
- Kijas J.W., Lenstra J.A., Hayes B. *et al.* (2012) Genome-wide analysis of the world's sheep breeds reveals high levels of historic mixture and strong recent selection. *PLoS Biology* **10**, e1001258.
- Kim K.K., Won H.H., Cho S.S., Park J.H., Kim M.J., Kim S. & Kim J.W. (2009) Comparison of identical single nucleotide polymorphisms

- genotyped by the GeneChip targeted genotyping 25K, Affymetrix 500K and Illumina 550K platforms. *Genomics* **94**, 89–93.
- Kunz E., Rothhammer S., Pausch H., Schwarzenbacher H., Seefried F.R., Matiasek K., Seichter D., Russ I., Fries R. & Medugorac I. (2016) Confirmation of a non-synonymous SNP in PNPLA8 as a candidate causal mutation for Weaver syndrome in Brown Swiss cattle. *Genetics Selection Evolution* **48**, 1–14.
- Laan M. & Pääbo S. (1997) Demographic history and linkage disequilibrium in human populations. *Nature Genetics* **17**, 435–8.
- Li H. & Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–60.
- Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G. & Durbin R. (2009) The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–9.
- Liu Y., Murali S.C., Harris R.A. *et al.* (2016) P1009 Sheep reference genome sequence updates: texel improvements and Rambouillet progress. *Journal of Animal Science* **94**, 18–9.
- Lopes F.B., Wu X.-L., Li H., Xu J., Perkins T., Genho J., Ferretti R., Tait R.G., Bauck S. & Rosa G.J.M. (2018) Improving accuracy of genomic prediction in Brangus cattle by adding animals with imputed low-density SNP genotypes. *Journal of Animal Breeding and Genetics* **135**, 14–27.
- Mägi R., Pfeufer A., Nelis M., Montpetit A., Metspalu A. & Remm M. (2007) Evaluating the performance of commercial whole-genome marker sets for capturing common genetic variation. *BMC Genomics* **8**, 159.
- Marguerat S., Wilhelm B.T. & Bähler J. (2008) Next-generation sequencing: applications beyond genomes. *Biochemical Society Transactions* **36**, 1091–6.
- Marina H., Gutiérrez-Gil B., Esteban-Blanco C., Suárez-Vega A., Pelayo R. & Arranz J.J. (2020) Analysis of whole genome resequencing datasets from a worldwide sample of sheep breeds to identify potential causal mutations influencing milk composition traits. *Animals* **10**, 1542.
- Martin P., Palière L., Maroteau C., Clément V., David I., Klopp G.T. & Rupp R. (2018) Genome-wide association mapping for type and mammary health traits in French dairy goats identifies a pleiotropic region on chromosome 19 in the Saanen breed. *Journal of Dairy Science* **101**, 5214–26.
- McClure M., Kim E., Bickhart D. *et al.* (2013) Fine mapping for Weaver syndrome in Brown Swiss cattle and the identification of 41 concordant mutations across NRCAM, PNPLA8 and CTTNBP2. *PLoS One* **8**, e59251.
- McKenna A., Hanna M., Banks E. *et al.* (2010) The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research* **20**, 1297–303.
- Meurs K.M., Lahmers S., Keene B.W., White S.N., Oyama M.A., Mauceli E. & Lindblad-Toh K. (2012) A splice site mutation in a gene encoding for PDK4, a mitochondrial protein, is associated with the development of dilated cardiomyopathy in the Doberman pinscher. *Human Genetics* **131**, 1319–25.
- Mokry F., Buzanskas M., de Alvarenga Mudadu M. *et al.* (2014) Linkage disequilibrium and haplotype block structure in a composite beef cattle breed. *BMC Genomics* **15**, 1–9.
- Nicolazzi E.L., Biffani S., Biscarini F., Orozco Ter Wengel P., Caprera A., Nazzicari N. & Stella A. (2015) Software solutions for the livestock genomics SNP array revolution. *Animal Genetics* **46**, 343–53.
- Pérez O'Brien A.M., Utsunomiya Y.T., Mészáros G., Bickhart D.M., Liu G.E., Van Tassell C.P., Sonstegard T.S., Da Silva M.V., Garcia J.F. & Sölkner J. (2014) Assessing signatures of selection through variation in linkage disequilibrium between taurine and indicine cattle. *Genetics Selection Evolution* **46**, 1–14.
- Purcell S., Neale B., Todd-Brown K. *et al.* (2007) PLINK: a tool set for whole-genome association and population-based linkage analyses. *American Journal of Human Genetics* **81**, 559–75.
- Ros-Freixedes R., Battagin M., Johnsson M., Gorjanc G., Mileham A.J., Rounsley S.D. & Hickey J.M. (2018) Impact of index hopping and bias towards the reference allele on accuracy of genotype calls from low-coverage sequencing. *Genetics Selection Evolution* **50**, 64.
- Salavati M., Caulton A., Clark R. *et al.* (2020) Global analysis of transcription start sites in the new ovine reference genome (Oar rambouillet v1.0). *Frontiers in Genetics* **11**, 1184. <https://doi.org/10.3389/fgene.2020.580580>
- Sambrook J., Fritsch E. & Maniatis T. (1983) Molecular cloning: a laboratory manual. *Immunology* **49**, 411.
- Sanchez M.P., Ramayo-Caldas Y., Wolf V. *et al.* (2019) Sequence-based GWAS, network and pathway analyses reveal genes co-associated with milk cheese-making properties and milk composition in Montbéliarde cows. *Genetics Selection Evolution* **51**. <https://doi.org/10.1186/s12711-019-0473-7>.
- Sims D., Sudbery I., Ilott N.E., Heger A. & Ponting C.P. (2014) Sequencing depth and coverage: key considerations in genomic analyses. *Nature Reviews Genetics* **15**, 121–32.
- Steemers F.J. & Gunderson K.L. (2007) Whole genome genotyping technologies on the BeadArray™ platform. *Biotechnology Journal* **2**, 41–9.
- Suarez B.K., Taylor C., Bertelsen S., Bierut L.J., Dunn G., Jin C.H., Kauwe JSK, Paterson A.D. & Hinrichs A.L. (2005) An analysis of identical single-nucleotide polymorphisms genotyped by two different platforms. *BMC Genomics* **1**(Suppl 1), S152.
- Sun J.K., Finch S.J., Haynes C. & Gordon D. (2004) Quantifying the percent increase in minimum sample size for SNP genotyping errors in genetic model-based association studies. *Human Heredity* **58**, 139–44.
- Taylor J.F., Whitacre L.K., Hoff J.L., Tizioto P.C., Kim J., Decker J.E. & Schnabel R.D. (2016) Lessons for livestock genomics from genome and transcriptome sequencing in cattle and other mammals. *Genetics Selection Evolution* **48**, 59.
- Tishkoff S.A. & Williams S.M. (2002) Genetic analysis of African populations: human evolution and complex disease. *Nature Reviews Genetics* **3**, 611–21.
- Van Den Berg S., Vandenplas J., Van Eeuwijk F.A., Bouwman A.C., Lopes M.S. & Veerkamp R.F. (2019) Imputation to whole-genome sequence using multiple pig populations and its use in genome-wide association studies. *Genetics Selection Evolution* **51**, 2.
- Wang W., Wei Z., Lam T.-W. & Wang J. (2011) Next generation sequencing has lower sequence coverage and poorer SNP-detection capability in the regulatory regions. *Scientific Reports* **1**, 55.
- Weller J.I. & Ron M. (2011) Invited review: quantitative trait nucleotide determination in the era of genomic selection. *Journal of Dairy Science* **94**, 1082–90.

- Wu X.-L., Xu J., Li H. *et al.* (2019) Evaluation of genotyping concordance for commercial bovine SNP arrays using quality-assurance samples. *Animal Genetics* **50**, 367–71.
- Wu Y., Fan H., Wang Y., Zhang L., Gao X., Chen Y., Li J., Ren H. & Gao H. (2014) Genome-wide association studies using haplotypes and individual SNPs in simmental cattle. *PLoS One* **9**, e109330.
- Wysoker A., Tibbetts K., McCowan M., Homer N. & Fennell T. (2019). Picard Toolkit [WWW Document]. Broad Institute, GitHub Repos. URL <http://broadinstitute.github.io/picard/> (accessed 7.3.18).
- Yates A.D., Achuthan P., Akanni W. *et al.* (2020) Ensembl 2020. *Nucleic Acids Research* **48**, D682–8.
- Zhao S., Jing W., Samuels D.C., Sheng Q., Shyr Y. & Guo Y. (2018) Strategies for processing and quality control of Illumina genotyping arrays. *Briefings in Bioinformatics* **19**, 765–75.

Supporting information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

Figure S1. Genotype concordance classification before the SNP quality control filtering.

Figure S2. Genotype concordance classification after the SNP quality control filtering.

Table S1. Pairwise of genotype concordance rate for each SNPs (C_i) between the three platforms before the SNP quality control filtering.

Table S2. Pairwise of genotype concordance rate for each SNPs (C_i) between the Illumina Ovine SNP50 BeadChip and the WGR before the SNP quality control filtering.