



**HAL**  
open science

## Combining genomics and semen microbiome increases the accuracy of predicting bull prolificacy

Pâmela Alexandre, Silvia Rodríguez-Ramilo, Núria Mach, Antonio Reverter

► **To cite this version:**

Pâmela Alexandre, Silvia Rodríguez-Ramilo, Núria Mach, Antonio Reverter. Combining genomics and semen microbiome increases the accuracy of predicting bull prolificacy. *Journal of Animal Breeding and Genetics*, 2024, 10.1111/jbg.12899 . hal-04694588

**HAL Id: hal-04694588**

**<https://hal.inrae.fr/hal-04694588>**

Submitted on 11 Sep 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

# Combining genomics and semen microbiome increases the accuracy of predicting bull prolificacy

Pâmela A. Alexandre<sup>1,2</sup>  | Silvia T. Rodríguez-Ramilo<sup>3</sup>  | Núria Mach<sup>4</sup>  | Antonio Reverter<sup>1,2</sup> 

<sup>1</sup>CSIRO MOSH-Future Science Platform, St Lucia, Queensland, Australia

<sup>2</sup>CSIRO Agriculture & Food, St Lucia, Queensland, Australia

<sup>3</sup>GenPhySE, Université de Toulouse, INRAE, ENVT, Castanet Tolosan, France

<sup>4</sup>IHAP, Université de Toulouse, INRAE, ENVT, Toulouse, France

## Correspondence

Pâmela A. Alexandre, CSIRO MOSH-Future Science Platform, St Lucia, Qld, Australia.

Email: [pamela.alexandre@csiro.au](mailto:pamela.alexandre@csiro.au)

## Funding information

CSIRO-INRAE linkage Travel Grant

## Abstract

Commercial livestock producers need to prioritize genetic progress for health and efficiency traits to address productivity, welfare, and environmental concerns but face challenges due to limited pedigree information in extensive multi-sire breeding scenarios. Utilizing pooled DNA for genotyping and integrating seminal microbiome information into genomic models could enhance predictions of male fertility traits, thus addressing complexities in reproductive performance and inbreeding effects. Using the Angus Australia database comprising genotypes and pedigree data for 78,555 animals, we simulated percentage of normal sperm (PNS) and prolificacy of sires, resulting in 713 sires and 27,557 progeny in the final dataset. Publicly available microbiome data from 45 bulls was used to simulate data for the 713 sires. By incorporating both genomic and microbiome information our models were able to explain a larger proportion of phenotypic variation in both PNS (0.94) and prolificacy (0.56) compared to models using a single data source (e.g., 0.36 and 0.41, respectively, using only genomic information). Additionally, models containing both genomic and microbiome data revealed larger phenotypic differences between animals in the top and bottom quartile of predictions, indicating potential for improved productivity and sustainability in livestock farming systems. Inbreeding depression was observed to affect fertility traits, which makes the incorporation of microbiome information on the prediction of fertility traits even more actionable. Crucially, our inferences demonstrate the potential of the semen microbiome to contribute to the improvement of fertility traits in cattle and pave the way for the development of targeted microbiome interventions to improve reproductive performance in livestock.

## KEYWORDS

amplicon sequencing, hologenomics, inbreeding, microbial relationship matrix

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2024 The Author(s). *Journal of Animal Breeding and Genetics* published by John Wiley & Sons Ltd.

## 1 | INTRODUCTION

Growing concerns with animal health, welfare, and the environment encourage livestock producers to accelerate genetic progress and prioritize selection for efficiency-related traits. This proves challenging for commercial livestock enterprises in extensive operations, primarily due to the lack of comprehensive pedigree information. The alternative is to determine the relationship between individuals using genomic information, albeit incurring the high cost associated with genotyping individual animals. In this context, pooling DNA from a group of animals to genotype is a promising cost-effective and practical solution (Alexandre, Porto-Neto, et al., 2019; Alexandre, Reverter, et al., 2019; Baller et al., 2022; Bell et al., 2017). Further, it is possible to assess sire reproductive performance based on the contribution of sires to each pool (Baller et al., 2020; Bennett et al., 2021).

Herd bulls are routinely joined to between 20 and 50 females per year and can have a working life of over 5 years. Thus, males with low reproductive performance can significantly affect herd pregnancy rates and cause inbreeding-related problems, such as reduction of mean fitness of an individual (a condition known as inbreeding depression) and losses in genetic diversity (Doekes et al., 2021). To assess reproductive performance, farmers often buy bulls with bull breeding soundness evaluation (BBSE) data that indicates sperm quality and mating ability. The actual performance of bulls in a multi-sire setting is usually unknown, as the number of progeny per sire is difficult to determine. Undoubtedly, reproductive performance entails complex multifactorial processes whose mechanisms are still not fully understood. Several conditions can contribute to poor reproductive performance, including animal behaviour, environmental factors, scrotal circumference, sperm morphology abnormalities, and the fertilizing efficiency of sperm (Alkhawagah et al., 2022; Corte Pause et al., 2022; Rowe et al., 2020). Perhaps not surprisingly, inbreeding depression has also been associated with poor semen quality and other fertility-related traits (Antonios et al., 2021; Ben Braiek et al., 2021; Makanjuola et al., 2020).

New evidence shows variation in the seminal microbiome and its associated metabolites can impact fertilization dynamics and pregnancy outcomes, including sperm aberrant motility, deficient mitochondrial function, and loss of DNA integrity (Altmäe et al., 2019). In healthy bulls, differences in the seminal microbiota have been related to fertility rates (Cojkic et al., 2021). For instance, *Veillonellaceae*, *Campylobacter*, *Methanobacterium*, and *Lawsonella* microbial signatures relate to sperm quality impairment, while *Bacteroides*, *Trueperella*, *Methanosphaera*, and *Methanobrevibacter* improve seminal parameters. Similarly, the microbial composition of

seminal fluid from bulls with satisfactory or unsatisfactory semen quality, assessed as poor sperm motility and morphology, exerted synergistic or antagonistic effects on sperm quality, depending on the bacterial genus (Kozioł et al., 2022).

Given the potential impact of reproductive microbiomes on host fertility and fitness, incorporating information about the seminal microbiome into genomic prediction models may accelerate genetic improvements in farm animals. This new biological scale has paved the way towards a new field of research referred to as hologenomics, which aims at integrating the genomic features of both the host and its microbiota (Saborío-Montero et al., 2021; Weishaar et al., 2020). In this regard, the heritable component of the microbiota (e.g., the proportion of total phenotypic variation in the host population that is due to variation in microbiome-encoded genetic factors that can be transmitted from one host to another) can be incorporated as a new source of information to explain phenotypic variation for reproductive performance traits (Venegas et al., 2023). Accordingly, simulations and real-data analyses performed by Pérez-Enciso et al. (2021) have shown an increased accuracy in predicting breeding values for livestock when including predicted microbial values. More recently, Hess et al. (2023) have shown that metagenome profiles from rumen samples in sheep improve the prediction accuracy of production efficiency and health traits.

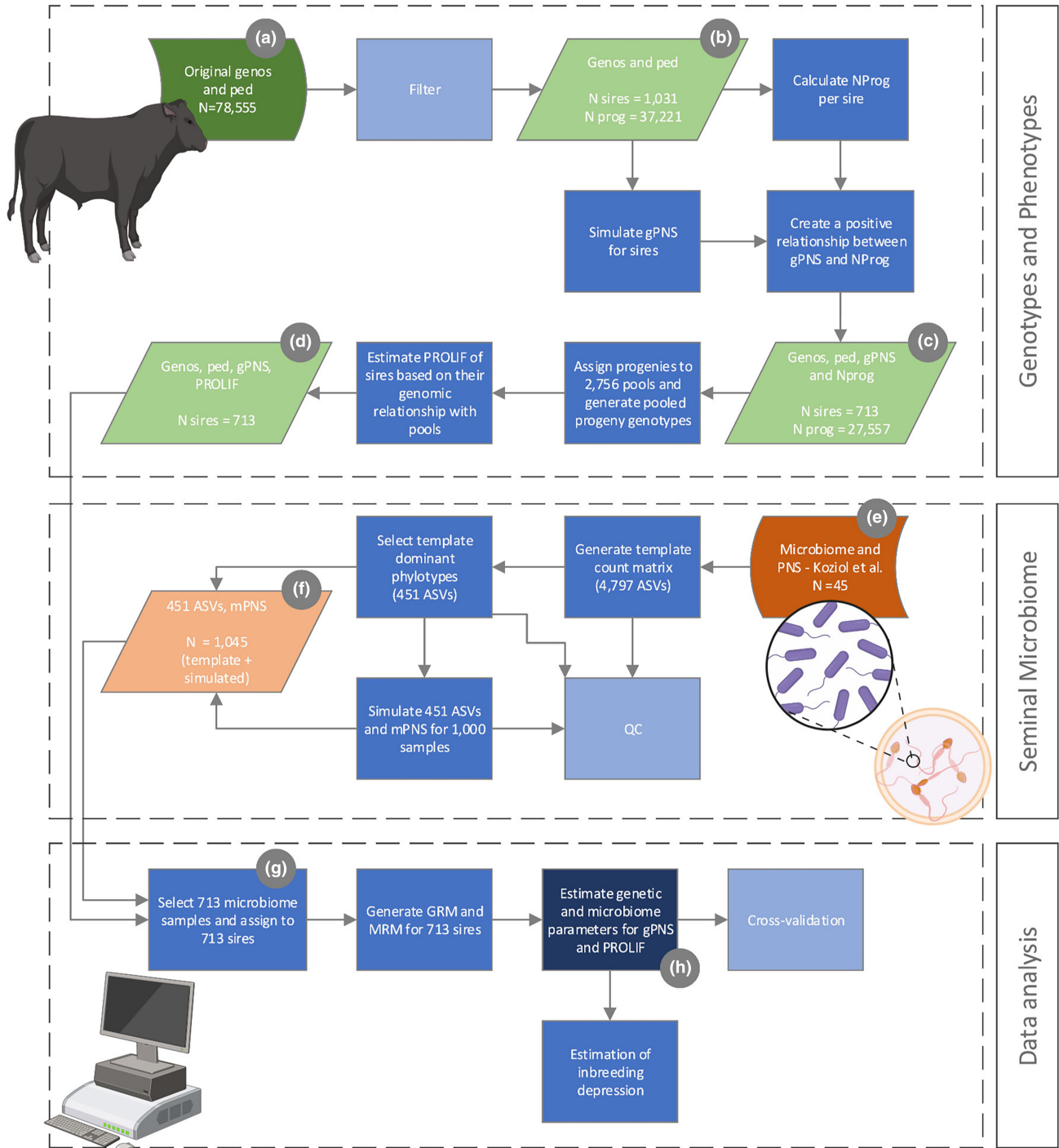
To our knowledge, there are no studies investigating the potential of combining seminal microbiome information in genomic predictions for male fertility traits. Here, we empirically show that inbreeding and reproductive microbiomes can have significant effects on the reproductive functions and performance of bulls in extensive conditions. We argue that knowledge of the reproductive microbiome is fundamental to our ability to predict reproductive performance.

## 2 | MATERIALS AND METHODS

### 2.1 | Genotypes and pedigree

Genotypes were sourced from the Angus Australia database and comprised 78,555 animals born from 2011 onwards, with imputed genotypes for 45,364 autosomal SNPs (Figure 1a). Pedigree records were examined to select genotypes from sires with at least five genotyped progeny. Genotyped progeny that were themselves sires were removed to avoid the inflated genomic relationship between the sire and the pool of progeny where the sire itself is represented.

These initial edits resulted in genotypes for 1031 sires and 37,221 progeny. The average number of progeny



**FIGURE 1** Data analysis pipeline. The original genotypes (genos) and pedigree (ped) were sourced from the Angus Australia database comprising 78,555 animals (a) and went through a preliminary filter to select 1031 sires and 37,221 progeny (b). This data was used to calculate the number of progeny per sire (Nprog) and simulate the percentage of normal sperm for sires based on their genotypes (gPNS). Then, a positive relationship between NProg and gPNS was simulated by removing some individuals, resulting in 713 sires and 27,557 progeny (c). Progeny genotypes were pooled, and sire prolificacy (PROLIF) was calculated based on the genomic relationship between sires and pools (d). The microbiome data was originally generated by Koziol et al. and comprised 16S amplicon sequencing for semen samples and PNS measures for 45 bulls (e). The resulting count matrix containing 4797 ASVs was reduced to the 451 dominant phylotypes and, together with the real PNS values, was used as a template to simulate data for 1000 individuals. The resulting dataset of 1045 microbiome samples (f) was then filtered to exclude samples with microbiome-simulated PNS (mPNS) outside of the ranges of the real data. The resulting microbiome samples (real + simulated) were assigned to sires based on gPNS and mPNS ranks, and only gPNS was kept for further analysis (g). This final dataset was used to generate a microbiome relationship matrix (MRM) and a genomic relationship matrix (GRM) for the 713 sires, which were then used to calculate inbreeding depression and genetic and microbiome parameters (h). Results were evaluated using a cross-validation scheme.

(NProg) per sire was 36.10, ranging from 5 to 715. For subsequent data analyses and to approximate normality, the base-2 log (Log<sub>2</sub>) was applied to NProg. In the Log<sub>2</sub> scale, the average NProg was 4.18 and ranged from 2.32 to 9.48.

## 2.2 | Phenotype simulation

The genotype file for the 1031 sires (Figure 1b) was used to simulate observations for the percentage of normal sperm (PNS). For the simulation, 100 equally spaced QTLs were assumed to have an effect sampled from a standard normal distribution. Following recently reported values (Porto-Neto et al., 2023), PNS mean, phenotypic variance, and heritability were assumed to be 75%, 600%, and 0.25, respectively. An in-house source code in FORTRAN95 was written to undertake the simulation. Genotype-simulated PNS values were bounded to 0 and 100%. To assess the quality of the simulation, a genomic relationship matrix (GRM) comprising all SNP, including those assigned to be QTL, across the 1031 sires was built using Method 1 of VanRaden (2008) and variance components were estimated using a GBLUP model in BlupF90 (Misztal et al., 2018).

Although there is no single physical or reproductive trait that will accurately determine a bull's ability to sire calves in natural mating settings, PNS has been shown to be positively related to calf output (Fitzpatrick et al., 2002; Holroyd et al., 2002). In our study, the original Nprog could not be considered a measure of prolificacy (or high fertility) because bulls were of different ages and management groups. Moreover, progeny were often not a result of natural mating but artificial insemination, which also can unevenly influence the number of Nprog. Hence, there was a need to create an artificial correlation between the genotype-simulated PNS and the number of progeny per sire to be kept in the database for further analysis. To do that, we examined the scatter plot of PNS and NProg across the 1031 sires (Figure S1A). In the first instance, the correlation between PNS and NProg was estimated at  $-0.089 \pm 0.031$ , indicating independence. Based on the observed ranges for both variables, a line of reference expectation was created, passing through the {x,y} points of {0,2} and {100,10}. For 81 sires (or 7.86%), their observed NProg was higher than the reference expectation, so they were regressed to meet this line. For the remaining 950 sires (or 92.14%) with NProg below the reference expectation, the NProg averaged 4.07 with an SD of 1.93 (Log<sub>2</sub> scale). Of these, 318 sires (or 30.84% of the 1031 total) and their progeny were removed due to having a difference in NProg between the reference expectation and the observed below the mean plus half an SD. The remaining 713 sires (of 27,557 progeny) showed a positive correlation between PNS and NProg estimated at  $0.418 \pm 0.034$

(Figure S1B) and were kept for subsequent analyses. This approach is similar to acceptance-rejection methods of simulation, in which a random sampling is performed on a two-dimensional Cartesian graph and only samples under a pre-defined curve are kept. Reassuringly, this correlation falls between what was observed in Santa Gertrudis (0.37) and Brahman (0.64) cattle (Holroyd et al., 2002).

## 2.3 | Simulation of progeny pools

The 27,557 progeny (Figure 1c) were randomly grouped into 2756 pools, including 2755 pools of 10 progeny each plus one pool of 7 progeny. To combine individual genotypes into pools, the frequency of the B allele was computed for all SNPs in each pool. Then, genotypes were determined as proposed by Alexandre, Porto-Neto, et al. (2019) according to four rules: (1) if the B-allele frequency  $\leq 0.17$ , then SNP genotype = 0; (2) if the B-allele frequency  $> 0.25$  and  $\leq 0.75$ , then SNP genotype = 1; (3) if the B-allele frequency  $> 0.82$ , then SNP genotype = 2; (4) if the B-allele frequency  $> 0.17$  and  $\leq 0.25$  or  $> 0.75$  and  $\leq 0.82$ , then a "flipping coin" Markov function assigned SNP genotype to 0 or 1, and to 1 or 2, respectively.

To assess the quality of the pooled genotypes, as well as the quality of the pedigree information, we built a GRM with the genotypes of the 713 sires plus those from the 2756 pools, and a Principal Component Analysis was performed. The prolificacy (PROLIF) of each sire was defined by the number of pools with a genomic relationship  $> 2.5\%$ , representing half of the relationship between sire and offspring (0.5) divided by pool size (Rowe et al., 2020). The resulting dataset is represented in Figure 1d.

For the present study, and because genotyping every single progeny would be prohibitively expensive, we anticipate that every sire would have been measured for PNS; however, its real NProg is not available, and only a measure of its PROLIF can be obtained via its genomic relationship with pools of potential progeny. Previous literature supports this approach, showing not only accurate bull prolificacy when estimated using genotypes from DNA pools of calves but also a high repeatability across years (Bennett et al., 2021).

## 2.4 | Microbiome simulation: $\alpha$ - and $\beta$ -diversity

The simulation of the semen microbiome for the 713 sires was based on the template microbiome dataset published by Koziol et al. (2022) with associated PNS phenotype for 45 bulls. The aim here was to use the microbial population structure and its relationship to real PNS to expand

the number of individuals for which we could have microbial profile information associated with a specific PNS to match our 713 sires. The Koziol dataset comprised 16S rRNA amplicon sequencing data for 45 beef bulls with ages between 12 months and 6 years of age (Figure 1e). Samples were collected from 11 pure breeds and breed crosses, including Angus, Simmental, Simmental-Angus Crosses, Gelbvieh Cross, Gelbvieh, Beefmaster, Chianina Cross, Crossbred, Hereford, and Shorthorn. Semen samples were collected as part of routine breeding soundness exams, which include the evaluation of semen quality via PNS, among other parameters. According to PNS values, 31 bulls were classified as “satisfactory” and 14 as “unsatisfactory”. One of the bulls classified as “unsatisfactory” had a missing PNS value, and, in this case, we used the average PNS for the other 13 animals classified as “unsatisfactory”.

Raw 16S rRNA sequencing data were downloaded from the NCBI SRA repository under the Bioproject number PRJNA747921 and biosample numbers from SAMN20300345 to SAMN20300393. The Divisive Amplicon Denoising Algorithm (DADA) was implemented using the DADA2 plug-in for QIIME 2 (v. 2028.8) (Bolyen et al., 2019) to perform quality filtering and chimera removal and to construct a feature table consisting of read counts per amplicon sequence variants (ASVs) by sample. Differently from Koziol et al. (2022), taxonomic assignments were given to ASVs by importing the Greengenes2 16S rRNA Database (McDonald et al., 2023) and extracting the regions of interest based on the primers used to generate the amplicons (515R – GTGCCAGCMGCCGCGTAA/806R – GGACTACHVGGGTWTCTAAT) to QIIME 2. The representative ASVs were selected using the naive Bayes q2-feature-classifier plug-in. The phyloseq (v.1.36.0), vegan (v.2.5.7) and microbiome (v. 1.15.3; <http://microbiome.github.io>) packages were used in R (v.4.1.0) for the downstream steps of analysis. A total of 5,566,078 high-quality sequence reads were recovered for the 45 bulls of the template study (mean per sample:  $129,443.7 \pm 119,915.5$ , range: 9351–508,793). Reads were clustered into 4797 chimera- and singleton-filtered ASVs at 99% sequence similarity.

To circumvent the problem of false-positive species predictions due to misalignment and contamination, we selected the common and dominant phylotypes in at least 10% of the samples for downstream analysis. Using this occurrence threshold of microbial taxa across multiple samples of the same cohort, we likely selected the most ecologically and functionally important seminal microbial taxa. To assess the congruency of taxonomic and structure data between all microbiota communities and dominant phylotype subsets, we performed Procrustes analysis on Euclidean distances on raw data using the R package *vegan*. The *protest()* function was used to perform repeated symmetric analyses and estimate if the degree

of association of the two matrices is greater than that expected by chance alone. We also used Mantel tests as a complementary analysis to examine correlations between the whole and dominant phylotypes  $\beta$ -distance matrices at the individual level, where each value represents the beta distance between a pair of individuals, using *vegan::mantel()* with the Spearman correlation method.

Based on those selected dominant phylotypes, we simulated 1000 seminal microbiomes using the function *synth\_comm\_from\_counts()* from the R package *SpiecEasi* (Sparse InversE Covariance estimation for Ecological Association and Statistical Inference, v. 1.1.2), which used a Normal to Anything (NorTA) approach (Kurtz et al., 2015). This function accounts for the sparsity, overdispersion, and compositionality found in microbiome data. Starting from non-normalized and non-rarified ASVs' count data, the function fits parameters based on a zero-inflated negative binomial distribution using ASV margins and simulates a new community with those properties. Because we expected that simulated PNS values would extend beyond the limits of the real dataset (i.e., 3%–100%), we simulated more samples ( $n = 1000$ ) than we needed ( $n = 713$ ) so that we could select the ones presenting realistic simulated values.

To ascertain how well the simulated data resemble the template data in terms of the relationship between PNS and ASVs, we compared the Pearson correlation coefficient distribution between PNS and ASVs in the two datasets (template vs. simulated data). Then, we calculated the sparsity (% zeros) and the overdispersion (coefficient of variation), and we compared the  $\alpha$ - and  $\beta$ -diversity between both template and simulated datasets. The  $\alpha$ - and  $\beta$ -diversity were calculated using the microbiome R package, which allowed us to study global indicators of the seminal ecosystem state, including measures of evenness, dominance, divergences, and abundance.  $\beta$ -diversity in both datasets was estimated via Bray-Curtis dissimilarity using the phyloseq R package. The  $\beta$ -diversity was visualized using the non-metric dimensional scaling (NMDS) in the *vegan* R package through the *metaMDS()* function. Then, the ecological community structure between template and simulated data was compared using the PerMANOVA test (a non-parametric multivariate analysis of variance based on pairwise distances) implemented in the *adonis2()* function from the *vegan* R package. The significance of the effect of the data type was assessed in an *F*-test based on the sequential sum of squares estimated from a 10,000 permutations procedure. The significance threshold was chosen at an adjusted  $p < 0.05$ . We also used Mantel tests as a complementary analysis to assess the relationship between the phylogenetic distance of pairwise ASVs and the Euclidean distances using Mantel correlations with 999 randomizations via the *mantel.correlog()* function in the *vegan* R package. Pairwise comparisons

of mean Bray-Curtis distances to group centroids among datasets were assessed using the permutational analysis of multivariate dispersion, *permdisp()* function in the *vegan* package. The core microbial group in real and simulated datasets was defined as the ASVs present in all 30% of the individuals, using a detection threshold of 0.1% in the microbiome R package. The disparity in the abundance of microbial species in the template and simulated communities was further assessed using commonly used indices to evaluate microbial communities, such as the Simpson and the Chao1 (Kim et al., 2017).

## 2.5 | Assignment of microbiome data to sires and generation of a microbial relationship matrix

Semen microbiome data, comprising both the 45 template samples and the 1000 simulated samples ( $n=1045$ , Figure 1f), was assigned to the 713 bulls (Figure 1d) by first removing samples with microbiome-simulated PNS outside of expected values (3%–100%). Then, to select 713 microbiomes, we ranked the samples based on the microbiome-simulated PNS and selected the top 356 and the bottom 357. Finally, we sorted bulls based on the genotype-simulated PNS and then by NProg (i.e., within a PNS value) and merged the two datasets, so that the higher-ranked bull based on the genome-simulated PNS was assigned the microbiome data associated with the higher microbiome-simulated PNS (Figure 1g). The sorting by NProg within genotype-simulated PNS was done so that a stronger link is created between the microbiome and NProg. The correlation between microbiome-simulated and genotype-simulated PNS was 0.97, and we kept for further analysis only the genome-simulated PNS.

Next, we generated a microbial relationship matrix for sires. Mainali et al. (2017) showed the inadequacy of Pearson's correlation coefficient as a measure of similarity in microbiome datasets and favoured the use of Jaccard's index of similarity, which uses the presence-absence data of the microbiome taxa. Consequently, in the present study, the microbiome relationship matrix (MRM) among sires was computed based on the Jaccard similarity index as follows: When comparing two sires,  $i$  and  $j$ , let  $T_i$ ,  $T_j$ , and  $T_{ij}$  denote the number of taxa present in sire  $i$  only, in sire  $j$  only, and co-present in both  $i$  and  $j$  sires, respectively. Then, the Jaccard's similarity index between sires  $i$  and  $j$  is defined as  $Jac(i,j) = T_{ij}/(T_i + T_j + T_{ij})$ . Note that for the diagonal elements of the MRM, when  $i=j$ , then  $T_i = T_j = 0$  and  $Jac(i,j) = 1$ . Similarly, off-diagonal elements of the MRM range from 0 when no taxa are co-present in both sires to 1 when no taxa are present in a sire-specific manner.

## 2.6 | Estimation of genetic and microbiome parameters

Three models were explored for the bi-variate analysis of PNS and PROLIF (Figure 1h):

Model 1 (GRM-ONLY) was:

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{g}_1 \\ \mathbf{g}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix}$$

Where  $\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}$  is the phenotype vector of length 713 + 713 for trait 1 (PNS) and 2 (PROLIF);  $\mathbf{1}$  is a 713 × 1 vector with all entries equal to 1;  $\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$  is the vector of population means for PNS and PROLIF;  $\mathbf{I}$  is an identity matrix of dimension 713 × 713;  $\begin{bmatrix} \mathbf{g}_1 \\ \mathbf{g}_2 \end{bmatrix}$  if the vector of random genomic values assumed to have a bi-variate normal distribution with mean zero and variance  $\mathbf{V}_g = \mathbf{H} \otimes \mathbf{G}$ , where  $\mathbf{H} = \begin{bmatrix} \sigma_{g1}^2 & \sigma_{g12} \\ \sigma_{g12} & \sigma_{g2}^2 \end{bmatrix}$ ,

$\mathbf{G}$  is the genomic relationship matrix (GRM) across the 713 sires built using the Method 1 of VanRaden (2008), and  $\otimes$

is the Kronecker product;  $\begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix}$  is the vector of random errors assumed to have a bi-variate normal distribution with mean zero and variance  $\mathbf{R} = \mathbf{R}_0 \otimes \mathbf{I}$ , where  $\mathbf{R}_0 = \begin{bmatrix} \sigma_{e1}^2 & \sigma_{e12} \\ \sigma_{e12} & \sigma_{e2}^2 \end{bmatrix}$ ;  $\sigma_{gi}^2$  and  $\sigma_{ei}^2$  represent the genetic and residual variance of trait  $i=1, 2$ ; and  $\sigma_{g12}$  and  $\sigma_{e12}$  are the genetic and residual covariance between trait 1 and trait 2.

Model 2 (MRM-ONLY) was:

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix}$$

Where elements are as defined before when applicable;  $\begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix}$  if the vector of random metagenomic values assumed to have a bi-variate normal distribution with mean zero and variance  $\mathbf{V}_m = \mathbf{J} \otimes \mathbf{M}$ , where  $\mathbf{J} = \begin{bmatrix} \sigma_{m1}^2 & \sigma_{m12} \\ \sigma_{m12} & \sigma_{m2}^2 \end{bmatrix}$ ,  $\mathbf{M}$  is the microbiome relationship matrix (MRM) across the 713 sires built using the Jaccard similarity index;  $\sigma_{mi}^2$  represents the metagenomic variance

of trait  $i=1, 2$ ; and  $\sigma_{m12}$  is the metagenomic covariance between trait 1 and trait 2.

Model 3 (GRM + MRM) was:

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \\ \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mathbf{m}_1 \\ \mathbf{m}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \\ \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix} \begin{bmatrix} \mathbf{g}_1 \\ \mathbf{g}_2 \end{bmatrix}$$

Where elements are as defined before.

In all cases, variance components and variance ratios (e.g., heritability, microbiability (Difford et al., 2018), and correlations) were estimated using both Qxpk5 (Pérez-Enciso & Misztal, 2011) and BLUPF90 (Misztal I, Tsuruta S, Lourenco D, Aguilar I, Legarra A, Vitezica Z. Manual for BLUPF90 family of programmes. University of Georgia; 2014).

## 2.7 | Estimation of inbreeding depression

Using the data from the 713 sires, a linear regression of the phenotypic values (for Prolificacy or PNS) or GEBV (for Prolificacy or PNS) on the inbreeding coefficients obtained from the diagonal of the GRM was performed to assess the magnitude of inbreeding depression (Doekes et al., 2021).

## 2.8 | Cross-validation accuracy of genomic and metagenomic predictions

For the cross-validation of genomic and metagenomic predictions, we created five validation datasets at random, each with the phenotypes from a random 20% set as missing (three with 143 records plus two with 142 records). For model M3 (GRM + MRM), predictions were based on either the GRM component only, the MRM component only, or the summation of the two. In each cross-validation schema, traditional (Bolormaa et al., 2013) and LR method (Legarra & Reverter, 2018) approaches were used to estimate accuracy, bias, and dispersion of predictions similarly to previous works (Alexandre et al., 2021). For bias and dispersion, we constructed 95% confidence intervals based on  $\pm 1.96$  SE around the observed means across the 10 scenarios, that is, 2 traits  $\times$  5 validation datasets. Finally, to translate prediction accuracies into real phenotype differences, animals in the validation population were ranked based on their predicted value, and the phenotypic differences between animals in the top and bottom quartiles (Q1Q4) were reported.

## 3 | RESULTS

### 3.1 | Relationship between percentage of normal sperm and sires' prolificacy

Starting from a dataset of 1031 sires and 37,221 progeny, the number of progeny per sire (Nprog) varied from 5 to 715 with an average of 36.10. For genotype-simulated percent normal sperm (PNS) using this same dataset, estimates of genetic variance and heritability were  $156.61 \pm 34.82\%$  and  $0.25 \pm 0.05$ , respectively, within one SE of reference values. After filtering sires and progeny to create a positive correlation ( $0.418 \pm 0.034$ ) between Nprog and PNS, a total of 27,557 progeny from 713 sires were kept for subsequent analysis. The average PNS for the remaining sires was 64.25% and ranged from 3.36% to 100%, while the NProg averaged 38.65 and ranged from 5 to 715.

In a real multi-sire setting, Nprog would be unknown, and only an estimation of prolificacy could be cost-effectively estimated based on the genomic relationship between each sire and the pools of progeny. Based on this approach, values of PROLIF averaged 34.06 and ranged from 1 to 210 (or 4.59, 0 and 7.71 for the same set of values in the log2 scale). The Pearson's correlation between the PROLIF of a sire based on the GRM and its real Nprog was  $0.908 \pm 0.016$ . Similarly, the correlation between PROLIF and PNS was  $0.383 \pm 0.035$ .

### 3.2 | The dominant basal seminal phylotypes represent the whole microbiota structure

After initial processing of the Koziol et al. (2022) seminal microbial data, we generated a count matrix containing 4797 ASVs and 45 samples (File S1). Then, we retained the most dominant phylotypes ( $n=451$  ASVs; File S2). These phylotypes harboured 170 unique genera and accounted for 69.9% of the annotated sequences. They were represented mainly by *Corynebacterium* (10%), *Bacteroides* (8.73%), *S5-A14a* (5.81%), and *Eremococcus* (4.63%). Importantly, Procrustes analysis showed a perfect alignment between the ordinations of the whole community (4791 ASVs) and the most dominant phylotypes (451 ASVs; Correlation in a symmetric Procrustes rotation: 0.9186). Likewise, Mantel tests indicated a high and positive relationship at individual-level microbiome distances (Mantel  $r \sim 0.8959$ ,  $p < 0.001$ ) between the whole and dominant community types for Bray-Curtis dissimilarity. Therefore, the set of 451 ASVs was selected for the downstream steps of analysis.



### 3.3 | Diversity and richness analysis in the template and simulated microbial datasets

Based on the dominant community ( $n=451$  ASVs), we then determined if our simulated microbial data (File S3) echoed that of Koziol's, using several measures. The Koziol data contained 78.06% zeros and was as sparse as the simulated data, which comprised 78.01% zeros (Figure 2a). However, the coefficient of variation (CV) of real data was 9.49%, while that of simulated data was 3.95%, indicating that the real data were more overdispersed. Then, the simulated communities were compared to the template regarding  $\alpha$  and  $\beta$  diversity. The SpiecEasi method successfully reproduced the overall  $\beta$  diversity, with no significant differences in both types of datasets (PerMANOVA,  $R^2=0.00029$ ,  $p=1$ ; Figure 2b,c). The intra- and interindividual microbiome variability was substantial in both types of data. However, the successional  $\beta$ -diversity dispersions were slightly higher in the simulated than in the real dataset (distance to the centroid of 0.49 and 0.69 for real and simulated data, respectively;  $p=2.5e-12$ ; Figure 2d). As a result, there was a larger disparity in the abundance of species in the simulated communities (Simpson index;  $p=2.07e^{-20}$ , Figure 2e), demonstrating the difficulty of simulating the seminal microbiome, which comprises many different species whose abundance profiles differ widely among samples. Nonetheless, regarding species occurrences, both scenarios capture the same rich structural complexity, with many rare species and only a few dominant common species (Chao1 index;  $p=0.98$ , Figure 2f). Although representing a skewed pattern in species abundance, the dominant taxa (Figure 2g) and the individual core microbiota, defined as any set of microbial taxa characteristic of the seminal fluid, remained similar between datasets. The core encompassed nine taxa, *Corynebacterium* dominated the assemblage, followed by *Escherichia*, *Bacteroides*, *Gemella*, *Eremococcus*, *S5-A14a*, *Methanobrevibacter*, *Streptococcus*, *Devosia\_A\_502124*, and *Aquicella\_A* (Figure 2e).

Finally, to assess the significant effects that seminal microbiomes can have on the reproductive function and fitness of males, we calculated the relationship between host PNS and the dominant phylotypes. Through their influence on host phenotypes, correlation values spanned a continuum in both datasets, from detrimental to beneficial. For the template data, the average correlation was 0, varying from  $-0.46$  to  $0.26$ , while for the simulated data, the average correlation was 0, varying from  $-0.35$  to  $0.20$ . Additionally, the correlation between the two correlation distributions was remarkably high at  $0.97$ .

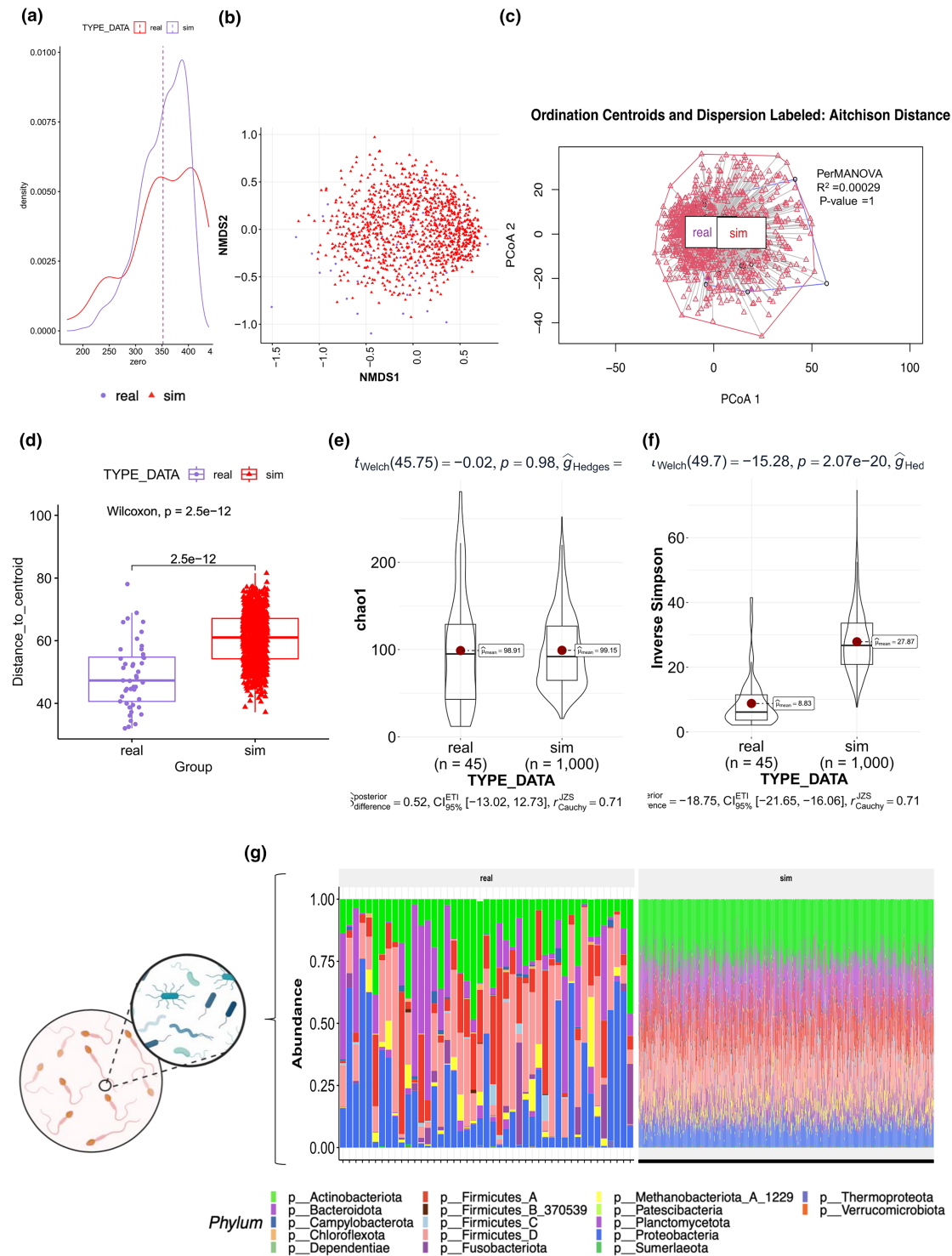
### 3.4 | Considering microbiome information to improve performance predictions

The microbiome matrix based on the presence-absence of dominant phylotypes (Jaccard similarity approach) provided a new approach for determining the reproductive bull traits. Given that semen microbial ecosystems have highly skewed distributions, that is, there are many rare species and only a few common species, the Jaccard index was as informative as methods that take their abundance into account (the Jaccard similarity matrix showed 0.99 correlation with a matrix based on Bray-Curtis dissimilarities). Based on this Jaccard distance matrix, across all pairwise sires of the MRM built for the 713 sires, the 253,828 off-diagonal elements averaged 0.186 and ranged from 0.013 to 0.646. Similarly, for the GRM, diagonal elements averaged  $1.00 \pm 0.05$  and off-diagonal elements averaged  $0.00 \pm 0.05$ . Similar standard deviations for diagonals and off-diagonals indicate a single population.

Variance components for PNS and PROLIF can be found in Table 1. Higher heritability and genetic variance were identified for PROLIF (using the GRM) compared to the microbiability and microbiome variance (using the MRM). Conversely, for PNS, higher microbiability and microbiome variance were identified when compared to heritability and genetic variance, suggesting that seminal microbiome influences host PNS phenotype. However, for both phenotypes, simultaneously including the GRM and the MRM in the model resulted in higher heritability+microbiability and overall variance explained. Table S1 shows the parameter estimates for PROLIF and PNS considering a 5-fold cross-validation. The results are comparable to those in Table 1 using the entire dataset.

The quality of the predictions, evaluated through bias, dispersion, and accuracy, can be seen in Table 2. In the absence of bias and dispersion, the values approximate zero, which is the case for most of our results. The exception is an overdispersion for predictions of PNS from model 3 based on the GRM only ( $0.2986 \pm 0.0557$ ), while there is an underdispersion also for PNS based on the MRM only ( $-0.1700 \pm 0.0672$ ). However, when results from model 3 are combined (GRM + MRM), the dispersion gets closer to zero. The accuracy of predictions reflects the heritability found for the different models and traits, following the same pattern discussed before.

Table S2 presents the results from the inbreeding depression for PROLIF and PNS when the response variable was the actual phenotype or the GEBV obtained from the GBLUP in Model 1 (GRM only). When using the actual PNS phenotype, no inbreeding depression was observed (and none simulated,  $p=0.54$ ). In all other cases, inbreeding depression was estimated as significant ( $p$ -value  $\leq 0.002$ ).



**FIGURE 2** Comparison of simulated and real seminal 16S rRNA data. (a) Distribution of zero counts per OTU; (b) NMDS ordination analysis (Bray-Curtis distance) of the 541 ASV composition. Points denote individual samples that are coloured according to the real (red) and simulated datasets (violet); (c) NMDS ordination plot showing centroids and ellipses for real and simulated datasets. The ellipses represent a calculated region of error around each group's centroid. The confidence level is set with  $conf=0.95$ ; (d) Box plots showing the Bray-Curtis distance to the centroid of the seminal microbial ASVs between the real and simulated data. Boxplots show the median, 25th, and 75th percentiles, the whiskers indicate the minima and maxima, and the points lying outside the whiskers of box plots represent the outliers. Adjusted  $p$  values from two-sided Wilcoxon rank-sum test; (e, f) Violin plots of Chao1 and Inverse Simpson indices between the microbial semen of real and simulated individuals. Boxes show median and interquartile range, and whiskers indicate the 5th to 95th percentile; (g) Taxonomic bar plots of the phyla in the semen according to the real and simulated data. Colours denote microbial phyla. Taxonomic inference relied on the QIIME closed-reference approach against the Greengenes2 16S rRNA Database at a sequence similarity level of 99% for the dominant 451 ASVs.

Parameter	Trait	Models			
		M1: GRM	M2: MRM	M3: GRM + MRM <sup>a</sup>	
$V_{g/m}$	PROLIF	0.550 ± 0.097	0.475 ± 0.039	0.477	0.280
	PNS	144.72 ± 23.51	286.58 ± 37.28	21.84	285.17
$h^2/m^2$	PROLIF	0.415 ± 0.066	0.308 ± 0.027	0.351	0.206
	PNS	0.365 ± 0.062	0.866 ± 0.040	0.067	0.871
$r_{g/m}$		0.363 ± 0.261	0.540 ± 0.024	0.283	0.619
$r_e$		0.370 ± 0.058	0.220 ± 0.292	0.169	

Abbreviations: GRM, genomic relationship matrix;  $h^2$ , heritability;  $m^2$ , microbiability; MRM, microbiome relationship matrix;  $r_e$ , residual correlation between PROLIF and PNS;  $r_{g/m}$ , genetic or microbiome correlation between PROLIF and PNS;  $V_{g/m}$ , genetic or microbiome variance.

<sup>a</sup>Estimates based on Qxpk software because of lack of convergence using Blupf90. So, no SE available.

Model	Component	Parameter	Trait	
			PNS	PROLIF
M1	GRM	Bias <sup>a</sup>	0.0306 ± 0.5381	0.0029 ± 0.0333
		Dispersion <sup>a</sup>	-0.0112 ± 0.1186	-0.0192 ± 0.07264
		ACC <sub>LR</sub>	0.4382	0.6289
		ACC <sub>T</sub>	0.3480	0.6063
		Q1Q4 <sup>b</sup>	9.43%	2.16 Progeny
M2	MRM	Bias	-0.1039 ± 0.8022	-0.0042 ± 0.0244
		Dispersion	-0.0568 ± 0.0591	0.0241 ± 0.0734
		ACC <sub>LR</sub>	0.8115	0.4786
		ACC <sub>T</sub>	0.8013	0.3964
		Q1Q4	36.96%	1.60 Progeny
M3	GRM	Bias	0.0474 ± 0.1692	0.0057 ± 0.0299
		Dispersion	0.2986 ± 0.0557	-0.0071 ± 0.0692
		ACC <sub>LR</sub>	0.3237	0.6131
		ACC <sub>T</sub>	0.5171	0.6577
		Q1Q4	9.67%	2.14 Progeny
M3	MRM	Bias	-0.1359 ± 0.8239	-0.0036 ± 0.0200
		Dispersion	-0.1700 ± 0.0672	0.0571 ± 0.0682
		ACC <sub>LR</sub>	0.8813	0.4487
		ACC <sub>T</sub>	0.9081	0.3853
		Q1Q4	37.28%	1.56 Progeny
M3	GRM + MRM	Bias	-0.0442 ± 0.4685	0.0010 ± 0.0219
		Dispersion	-0.1516 ± 0.0730	0.0258 ± 0.0825
		ACC <sub>LR</sub>	0.5818	0.3892
		ACC <sub>T</sub>	0.8239	0.5349
		Q1Q4	37.37%	2.34 Progeny

Abbreviations: GRM, genomic relationship matrix; MRM, microbiome relationship matrix.

<sup>a</sup>Bias and Dispersion are given with ± SE.

<sup>b</sup>Q1Q4: Average phenotype difference between the top (Q1) and the bottom (Q4) quartiles. For Prolificacy the value is given in real progeny units (i.e., Not log2 scale).

**TABLE 1** Estimates (±SE) of variance components and ratios for prolificacy (PROLIF) and percentage of normal sperm (PNS) based on 713 sires and three bivariate models.

**TABLE 2** Bias, dispersion, accuracy (ACC), and Q1Q4 of genomic (GRM) and microbiome (MRM) predictions (averages across the 5 validation datasets) for prolificacy (PROLIF) and percentage of normal sperm (PNS).

## 4 | DISCUSSION

In this study, we endeavour to understand the potential of the semen microbiome to aid male fertility

predictions and control the negative impacts of suboptimal fertility on inbreeding depression. We designed a simulation study recognizing the importance of accurately mimicking the real-world scenario envisioned for



the prospective commercial application of the approach in which it will contribute to more sustainable livestock farming systems. Our results demonstrate the potential contribution of the hologenome on fertility indicator traits (PNS and Prolificacy) considering that linear mixed models containing both genomic and microbiome information were able to explain a larger portion of the phenotypic variation compared to models containing a single source of information. That was then reflected in higher phenotypic differences expected for animals ranked at the top or bottom of hologenome-based prediction estimates. These results suggest a benefit in considering the hologenome for genomic prediction of fertility traits. Microbiome-encoded genetic factors affecting host traits also open the possibility of influencing seminal quality and reproductive performance via the production of tailor-made microbiomes in the semen through artificial selection, shaping the genetic composition of microbiomes independently of host genome selection, which could certainly improve animal performance and sustainability of the production system in the hologenomic era (Mueller & Linksvayer, 2022). However, in the future we also need to make sure these microbiomes are transmitted between hosts with sufficient fidelity for microbiome breeding to work.

The relationship structure of a population is affected by the accumulation of inbreeding through selection. However, not all pedigree-based genetic evaluation methods consider inbreeding explicitly. Consequently, failing to contemplate inbreeding when determining a relationship will affect, for example, how **A** versus **G** scales to generate **H** (Garcia-Baccino et al., 2017). Fortunately, this is not the case for **G** (Gowane et al., 2019). Because **G** considers inbreeding implicitly, estimates of inbreeding depression derived from GEBV are bound to be more significant than estimates derived from phenotypic data. In our case, this indicates that the GEBV for PROLIF and PNS will be lower, on average, in inbred sires. There are already some examples of significant inbreeding depression for traits related to bull prolificacy (Dorado et al., 2017; Ghoreishifar et al., 2023).

Indeed, significant inbreeding depression was observed for PROLIF when using phenotypes. An increase in inbreeding of 1% decreases PROLIF by 0.83% on average ( $p$ -value  $<0.0001$ ), which agrees with reported values of inbreeding depression on fertility. However, no inbreeding effect (depression or boosting) was modelled when simulating PNS data based on true genotypes. Consequently, no significant inbreeding depression was observed for this trait ( $p$ -value  $>0.53$ ). More accentuated significant inbreeding depression for both traits was estimated when using GEBV because genomic relationships automatically account for inbreeding. This means that an increase of 1%

in inbreeding in GEBV for Prolificacy will reduce this trait by 0.88% on average ( $p$ -value  $<2.2 \times 10^{-16}$ ). Accordingly, a 1% increase in inbreeding decreases PNS obtained from GEBV by 0.36% ( $p$ -value  $<0.0002$ ). Considering that the microbiome was related to genotype information through PNS, these results anticipate the added importance of microbiome data in the presence of inbreeding. Indeed, the importance of the microbiome information increases when inbreeding is taken into consideration. In extreme circumstances, having a detailed understanding of the microbiome becomes essential.

Although there is growing evidence of the benefits of considering the hologenome in genomic predictions (Hess et al., 2023; Saborio-Montero et al., 2021; Weishaar et al., 2020), the relative contribution of host genetics and microbiome profile is bound to vary according to the trait in question and the microbial sample site. For instance, Camarinha-Silva et al. (2017) reported higher estimates of microbiability based on pig gastrointestinal microbiota for feed conversion ratio ( $0.21 \pm 0.14$ ) and feed intake ( $0.16 \pm 0.10$ ) than their corresponding heritabilities ( $0.19 \pm 13$  and  $0.11 \pm 11$ , respectively). The opposite has also been observed. Using rumen microbial composition in dairy cattle, Difford et al. (2018) showed higher heritabilities for methane emissions ( $0.19 \pm 0.09$ ) compared to microbiability ( $0.15 \pm 0.08$ ). Importantly, in this example, a model containing both effects was able to explain 34% of the total phenotypic variation, demonstrating the value of the combined information. In our study, we have both cases. While for PNS we observed a much higher microbiability (0.87) compared to heritability (0.36), the opposite was true for PROLIF ( $h^2=0.41$ ,  $m^2=0.31$ ), although with a smaller difference. However, in both cases, the model containing both sources of information explained altogether a higher proportion of the phenotypic variation (i.e., 0.56 for PROLIF and 0.94 for PNS). By adding information about the microbiome profiles into our bi-variate analyses, we stated that the reproductive microbiome can significantly affect the reproduction function and performance of males and aid genomic selection.

It is worthwhile noting possible limitations in the study. Firstly, microbiability for PNS might be overestimated as a result of our approach to simulate the relationship between the trait for each individual and its respective microbial profile, which was the most challenging step of the simulation. Using bulls as a model, we retrieved genetic information from 1031 sires and 37,221 progeny from the highly curated Angus Australia database, all with 45K genotype and pedigree information. However, these bulls had no records of male fertility traits and seminal microbiome, which had to be simulated. Phenotype simulation based on genotype information is a widely acceptable strategy to explore difficult-to-measure

traits, and the methodology to do that has been used by many authors and implemented by different software (Sargolzaei & Schenkel, 2009). However, to date, simulation of the seminal microbiome is scarce in mammals and still unexplored in livestock. Simulated values for PNS were within one SE of literature reference values (Porto-Neto et al., 2023), demonstrating a good reproduction of real measurements. Similarly, from the availability of Koziol's microbiome data, we generated a trustworthy simulated seminal microbiome dataset, which accounted for the data's sparsity, overdispersion, and compositionality, and included an associated PNS value simulated based on real data. However, combining both datasets by sorting data based on genotype- and microbiome-simulated PNS does not preserve the exact association pattern between PNS and microbial profile in the original (real) dataset, even though the correlation between genotype- and microbiome-simulated PNS in the final dataset was 0.97. Nevertheless, based on the soundness of our simulation strategy and results, we are confident that there is an important contribution of the seminal microbiome on PNS and prolificacy, which needs to be validated through the generation of real data.

In terms of prediction accuracy, the models containing both GRM and MRM did not necessarily show higher accuracies compared to the other models, particularly for accuracies calculated based on LR methods, which the algebra does not allow for an accurate estimate of the combined accuracy based on both GRM and MRM. However, for both traits, expected phenotypic differences between animals predicted to be at the top and bottom quartile of predictions (Q1Q4) were higher using the hologenome model. For PNS, we saw a difference of 37.37 in average percentage of normal sperm using the hologenome model, while the model using only genomic information resulted in a difference of 9.43%. In terms of progeny, Q1Q4 using the hologenome model yielded a difference of 2.34 progeny, while the model based on GRM resulted in a difference of 2.16 progeny. These differences expected across multiple bulls in a herd can significantly improve productivity, resource utilization, and sustainability.

The potential contribution of the semen microbiome on male fertility-related traits demonstrated here open several exciting possibilities. If proven accurate based on real data, current genetic selection models could include a holonomic scheme for selecting desirable reproductive traits. In addition, bulls with the most seminal beneficial microbiomes should be identified for microbiome harvesting and transplanting, facilitating response to selection. Like animal breeding programs, microbiome selection is focused on achieving phenotypic outcomes in terms of the traits and is agnostic to the specific microbiome

composition and function. A better understanding of the functional properties of semen microbiome and its implications for semen quality traits might improve the efficiency of microbiome breeding. For example, it is known that in humans, the adhesion of *Escherichia coli* to sperm cells leads to sperm agglutination and destruction of the sperm plasma membrane, with negative consequences for sperm motility and ultrastructure (Diemer et al., 2000). Alternatively, the release or active secretion of bacterial membrane proteins has been shown to impair sperm function, possibly by inhibiting macrophage function or induction of excessive reactive oxygen species (ROS) production (Eley et al., 2005). These findings are timely when considering new technological advances in microbiome studies. Gene-editing technologies (e.g., CRISPR-Cas9) can be used to engineer synthetic microbial communities (SynComs) or bacterial metabolites to enhance the quality of the semen with positive consequences for bull prolificacy and the overall efficiency and sustainability of livestock production.

## 5 | CONCLUSION

Our simulation study was designed to reflect real-world scenarios, demonstrating the promise of hologenomic approaches in enhancing productivity and sustainability in livestock farming systems. By incorporating both genomic and microbiome information, our results highlight the significant contribution of the hologenome to fertility indicator traits, emphasizing the importance of considering microbiome data, particularly in the presence of inbreeding. We need to acknowledge that, as with any simulation study, our results are based on a series of pre-defined parameters. For instance, the correlation between PNS and prolificacy that, although follows expectations according to the literature, was artificially generated in this study. That brings some insight into the results we can expect from real scenarios, but future validation of these findings with real data is crucial and could pave the way for the integration of microbiome-based selection strategies and the development of targeted microbiome interventions to improve reproductive performance in livestock.

## ACKNOWLEDGEMENTS

The authors acknowledge the contribution of Angus Australia for kindly providing the dataset used in this study. The authors also acknowledge the valuable contributions of Dr. James Kijas and Dr. Scott Rice in reviewing the text with a focus on the genetics and microbiome aspects, respectively. This research was funded by the CSIRO-INRAe Linkage Travel Grant awarded to Pâmela Alexandre and Silvia Rodríguez-Ramilo.

## CONFLICT OF INTEREST STATEMENT

None of the authors have a conflict of interest to disclose.

## DATA AVAILABILITY STATEMENT

Koziol et al. (2022) raw semen microbiome sequence data are available at the National Center for Biotechnology Information (NCBI) Sequence Read Archive (SRA), Bioproject PRJNA747921, Biosamples SAMN20300345–SAMN20300393. Results of downstream analysis of this data, including ASVs' raw counts per sample and taxonomic classification, as well as simulated microbiome data, are available as supplementary files in this publication. Angus cattle genotypes and pedigree information belong to the Angus Australia Breed Association and are available under specific collaboration agreements. No software codes were generated during this work, and all the analyses were done using publicly available software. However, Unix and R scripts used to run specific analyses are available upon request.

## ORCID

Pâmela A. Alexandre  <https://orcid.org/0000-0002-0649-7033>

Silvia T. Rodríguez-Ramilo  <https://orcid.org/0000-0001-7150-0692>

Núria Mach  <https://orcid.org/0000-0002-8001-6314>

Antonio Reverter  <https://orcid.org/0000-0002-4681-9404>

## REFERENCES

- Alexandre, P. A., Li, Y., Hine, B. C., Duff, C. J., Ingham, A. B., Porto-Neto, L. R., & Reverter, A. (2021). Bias, dispersion, and accuracy of genomic predictions for feedlot and carcass traits in Australian Angus steers. *Genetics Selection Evolution*, *53*, 1–10.
- Alexandre, P. A., Porto-Neto, L. R., Karaman, E., Lehnert, S. A., & Reverter, A. (2019). Pooled genotyping strategies for the rapid construction of genomic reference populations. *Journal of Animal Science*, *97*, 4761–4769.
- Alexandre, P. A., Reverter, A., Lehnert, S. A., Porto-Neto, L. R., & Dominik, S. (2019). In silico validation of pooled genotyping strategies for genomic evaluation in Angus cattle. *Journal of Animal Science*, *98*, 1–5.
- Alkhawagah, A. R., Ricci, A., Banchi, P., Martino, N. A., Poletto, M. L., Donato, G. G., Nervo, T., & Vincenti, L. (2022). Effect of epidermal growth factor (EGF) on cryopreserved Piedmontese bull semen characteristics. *Animals*, *12*, 3179.
- Altmäe, S., Franasiak, J. M., & Mändar, R. (2019). The seminal microbiome in health and disease. *Nature Reviews. Urology*, *16*, 703–721.
- Antonios, S., Rodríguez-Ramilo, S. T., Aguilar, I., Astruc, J. M., Legarra, A., & Vitezica, Z. G. (2021). Genomic and pedigree estimation of inbreeding depression for semen traits in the Basco-Béarnaise dairy sheep breed. *Journal of Dairy Science*, *104*, 3221–3230.
- Baller, J. L., Kachman, S. D., Kuehn, L. A., & Spangler, M. L. (2020). Genomic prediction using pooled data in a single-step genomic best linear unbiased prediction framework. *Journal of Animal Science*, *98*, 1–12.
- Baller, J. L., Kachman, S. D., Kuehn, L. A., & Spangler, M. L. (2022). Using pooled data for genomic prediction in a bivariate framework with missing data. *Journal of Animal Breeding and Genetics*, *139*, 489–501.
- Bell, A. M., Henshall, J. M., Porto-Neto, L. R., Dominik, S., McCulloch, R., Kijas, J., & Lehnert, S. A. (2017). Estimating the genetic merit of sires by using pooled DNA from progeny of undetermined pedigree. *Genetics Selection Evolution*, *49*, 1–7.
- Ben Braiek, M., Fabre, S., Hozé, C., Astruc, J. M., & Moreno-Romieux, C. (2021). Identification of homozygous haplotypes carrying putative recessive lethal mutations that compromise fertility traits in French Lacaune dairy sheep. *Genetics Selection Evolution*, *53*, 41.
- Bennett, G. L., Keele, J. W., Kuehn, L. A., Snelling, W. M., Dickey, A. M., Light, D., Cushman, R. A., & McDaneld, T. G. (2021). Using genomics to measure Phenomics: Repeatability of bull prolificacy in multiple-bull pastures. *Agriculture*, *11*, 603.
- Bolormaa, S., Pryce, J. E., Kemper, K., Savin, K., Hayes, B. J., Barendse, W., Zhang, Y., Reich, C. M., Mason, B. A., Bunch, R. J., Harrison, B. E., Reverter, A., Herd, R. M., Tier, B., Graser, H. U., & Goddard, M. E. (2013). Accuracy of prediction of genomic breeding values for residual feed intake and carcass and meat quality traits in *Bos taurus*, *Bos indicus*, and composite beef cattle. *Journal of Animal Science*, *91*, 3088–3104.
- Bolyen, E., Rideout, J. R., Dillon, M. R., Bokulich, N. A., Abnet, C. C., al-Ghalith, G. A., Alexander, H., Alm, E. J., Arumugam, M., Asnicar, F., Bai, Y., Bisanz, J. E., Bittinger, K., Brejnrod, A., Brislawn, C. J., Brown, C. T., Callahan, B. J., Caraballo-Rodríguez, A. M., Chase, J., ... Caporaso, J. G. (2019). Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nature Biotechnology*, *37*, 852–857.
- Camarinha-Silva, A., Maushammer, M., Wellmann, R., Vital, M., Preuss, S., & Bennewitz, J. (2017). Host genome influence on gut microbial composition and microbial prediction of complex traits in pigs. *Genetics*, *206*, 1637–1644.
- Cojkic, A., Niazi, A., Guo, Y., Hallap, T., Padrik, P., & Morrell, J. M. (2021). Identification of bull semen microbiome by 16S sequencing and possible relationships with fertility. *Microorganisms*, *9*, 2431.
- Corte Pause, F., Crociati, M., Urli, S., Monaci, M., Degano, L., & Stradioli, G. (2022). Environmental factors affecting the reproductive efficiency of Italian Simmental young bulls. *Animals*, *12*, 2476.
- Diemer, T., Huwe, P., Michelmann, H. W., Mayer, F., Schiefer, H. G., & Weidner, W. (2000). *Escherichia coli* -induced alterations of human spermatozoa. An electron microscopy analysis. *International Journal of Andrology*, *23*, 178–186.
- Difford, G. F., Plichta, D. R., Løvendahl, P., Lassen, J., Noel, S. J., Højberg, O., Wright, A. D. G., Zhu, Z., Kristensen, L., Nielsen, H. B., Gulbrandtsen, B., & Sahana, G. (2018). Host genetics and the rumen microbiome jointly associate with methane emissions in dairy cows. *PLoS Genetics*, *14*, e1007580.
- Doekes, H. P., Bijma, P., & Windig, J. J. (2021). How depressing is inbreeding? A meta-analysis of 30 years of research on the effects of inbreeding in livestock. *Genes (Basel)*, *12*, 926.
- Dorado, J., Cid, R. M., Molina, A., Hidalgo, M., Ariza, J., Moreno-Millán, M., & Demyda-Peyrás, S. (2017). Effect of inbreeding depression on bull sperm quality and field fertility. *Reproduction, Fertility, and Development*, *29*, 712.

- Eley, A., Pacey, A. A., Galdiero, M., Galdiero, M., & Galdiero, F. (2005). Can chlamydia trachomatis directly damage your sperm? *The Lancet Infectious Diseases*, 5, 53–57.
- Fitzpatrick, L. A., Fordyce, G., McGowan, M. R., Bertram, J. D., Doogan, V. J., de Faveri, J., Miller, R. G., & Holroyd, R. G. (2002). Bull selection and use in northern Australia part 2. Semen traits. *Animal Reproduction Science*, 71, 39–49.
- García-Baccino, C. A., Legarra, A., Christensen, O. F., Misztal, I., Pocrnic, I., Vitezica, Z. G., & Cantet, R. J. C. (2017). Metafounders are related to F<sub>st</sub> fixation indices and reduce bias in single-step genomic evaluations. *Genetics Selection Evolution*, 49, 34.
- Ghoreishifar, M., Vahedi, S. M., Salek Ardestani, S., Khansefid, M., & Pryce, J. E. (2023). Genome-wide assessment and mapping of inbreeding depression identifies candidate genes associated with semen traits in Holstein bulls. *BMC Genomics*, 24, 230.
- Gowane, G. R., Lee, S. H., Clark, S., Moghaddar, N., al-Mamun, H. A., & van der Werf, J. H. J. (2019). Effect of selection and selective genotyping for creation of reference on bias and accuracy of genomic prediction. *Journal of Animal Breeding and Genetics*, 136, 390–407.
- Hess, M. K., Zetouni, L., Hess, A. S., Budel, J., Dodds, K. G., Henry, H. M., Brauning, R., McCulloch, A. F., Hickey, S. M., Johnson, P. L., Elmes, S., Wing, J., Bryson, B., Knowler, K., Hyndman, D., Baird, H., McRae, K. M., Jonker, A., Janssen, P. H., ... Rowe, S. J. (2023). Combining host and rumen metagenome profiling for selection in sheep: Prediction of methane, feed efficiency, production, and health traits. *Genetics Selection Evolution*, 55, 53.
- Holroyd, R. G., Doogan, V. J., De Faveri, J., Fordyce, G., McGowan, M. R., Bertram, J. D., Vankan, D. M., Fitzpatrick, L. A., Jayawardhana, G. A., & Miller, R. G. (2002). Bull selection and use in northern Australia 4. Calf output and predictors of fertility of bulls in multiple-sire herds. *Animal Reproduction Science*, 71, 67–79.
- Kim, B. R., Shin, J., Guevarra, R. B., Lee, J. H., Kim, D. W., Seol, K. H., Lee, J. H., Kim, H. B., & Isaacson, R. E. (2017). Deciphering diversity indices for a better understanding of microbial communities. *Journal of Microbiology and Biotechnology*, 27, 2089–2093.
- Kozioł, J. H., Sheets, T., Wickware, C. L., & Johnson, T. A. (2022). Composition and diversity of the seminal microbiota in bulls and its association with semen parameters. *Theriogenology*, 182, 17–25.
- Kurtz, Z. D., Müller, C. L., Miraldi, E. R., Littman, D. R., Blaser, M. J., & Bonneau, R. A. (2015). Sparse and compositionally robust inference of microbial ecological networks. *PLoS Computational Biology*, 11, e1004226.
- Legarra, A., & Reverter, A. (2018). Semi-parametric estimates of population accuracy and bias of predictions of breeding values and future phenotypes using the LR method. *Genetics Selection Evolution*, 50, 53.
- Mainali, K. P., Bewick, S., Thielen, P., Mehoke, T., Breitwieser, F. P., Paudel, S., Adhikari, A., Wolfe, J., Slud, E. V., Karig, D., & Fagan, W. F. (2017). Statistical analysis of co-occurrence patterns in microbial presence-absence datasets. *PLoS One*, 12, e0187132.
- Makanjuola, B. O., Maltecca, C., Miglior, F., Schenkel, F. S., & Baes, C. F. (2020). Effect of recent and ancient inbreeding on production and fertility traits in Canadian Holsteins. *BMC Genomics*, 21, 605.
- McDonald, D., Jiang, Y., Balaban, M., Cantrell, K., Zhu, Q., Gonzalez, A., Morton, J. T., Nicolaou, G., Parks, D. H., Karst, S. M., Albertsen, M., Hugenholtz, P., DeSantis, T., Song, S. J., Bartko, A., Havulinna, A. S., Jousilahti, P., Cheng, S., Inouye, M., ... Knight, R. (2023). Greengenes2 unifies microbial data in a single reference tree. *Nature Biotechnology*, 42, 715–718. <https://doi.org/10.1038/s41587-023-01845-1>
- Misztal, I., Tsuruta, S., Lourenco, D., Masuda, Y., Aguilar, I., Legarra, A., & Vitezica, Z. (2018). *BLUPF90*. [Preprint].
- Mueller, U. G., & Linksvayer, T. A. (2022). Microbiome breeding: Conceptual and practical issues. *Trends in Microbiology*, 30, 997–1011.
- Pérez-Enciso, M., & Misztal, I. (2011). Qxpk: 5: Old mixed model solutions for new genomics problems. *BMC Bioinformatics*, 12, 202.
- Pérez-Enciso, M., Zingaretti, L. M., Ramayo-Caldas, Y., & de los Campos, G. (2021). Opportunities and limits of combining microbiome and genome data for complex trait prediction. *Genetics Selection Evolution*, 53, 65.
- Porto-Neto, L. R., Alexandre, P. A., Hudson, N. J., Bertram, J., McWilliam, S. M., Tan, A. W. L., Fortes, M. R. S., McGowan, M. R., Hayes, B. J., & Reverter, A. (2023). Multi-breed genomic predictions and functional variants for fertility of tropical bulls. *PLoS One*, 18, e0279398.
- Rowe, M., Veerus, L., Trosvik, P., Buckling, A., & Pizzari, T. (2020). The reproductive microbiome: An emerging driver of sexual selection, sexual conflict, mating systems, and reproductive isolation. *Trends in Ecology & Evolution*, 35, 220–234.
- Saborío-Montero, A., Gutiérrez-Rivas, M., López-García, A., García-Rodríguez, A., Atxaerandio, R., Goiri, I., Jiménez-Montero, J. A., & González-Recio, O. (2021). Holobiont effect accounts for more methane emission variance than the additive and microbiome effects on dairy cattle. *Livestock Science*, 250, 104538.
- Sargolzaei, M., & Schenkel, F. S. (2009). QMSim: A large-scale genome simulator for livestock. *Bioinformatics*, 25, 680–681.
- VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *Journal of Dairy Science*, 91, 4414–4423.
- Venegas, L., López, P., Derome, N., & Yáñez, J. M. (2023). Leveraging microbiome information for animal genetic improvement. *Trends in Genetics*, 39, 721–723.
- Weishaar, R., Wellmann, R., Camarinha-Silva, A., Rodehutschord, M., & Bennewitz, J. (2020). Selecting the hologenome to breed for an improved feed efficiency in pigs—A novel selection index. *Journal of Animal Breeding and Genetics*, 137, 14–22.

## SUPPORTING INFORMATION

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

**How to cite this article:** Alexandre, P. A., Rodríguez-Ramilo, S. T., Mach, N., & Reverter, A. (2024). Combining genomics and semen microbiome increases the accuracy of predicting bull prolificacy. *Journal of Animal Breeding and Genetics*, 00, 1–14. <https://doi.org/10.1111/jbg.12899>