

1
2 **Opportunities and limits of combining microbiome and genome data**
3 **for complex trait prediction**

4
5 **Miguel Pérez-Enciso^{1,2,4}, Laura M. Zingaretti^{2,4}, Yulixaxis Ramayo-Caldas³,**
6 **Gustavo de los Campos⁴**

7
8 1 ICREA, Passeig de Lluís Companys 23, 08010 Barcelona, Spain

9 2 Centre for Research in Agricultural Genomics (CRAG), CSIC-IRTA-UAB-UB,
10 08193 Bellaterra, Barcelona, Spain

11 3 Animal Breeding and Genetics Program, Institute for Research and Technology in
12 Food and Agriculture (IRTA), Torre Marimon, 08140 Caldes de Montbui, Barcelona,
13 Spain

14 4 Michigan State University, Dept. of Epidemiology & Biostatistics, and Dept. of
15 Statistics & Probability, East Lansing, MI 48824, USA.

16
17 Short title:

18 Combining microbiome and genome data for complex trait prediction

19
20 Correspondence:

21 M. Pérez-Enciso

22 CRAG

23 Campus UAB

24 08193 Bellaterra, Spain

25 miguel.perez@uab.es

26

27

28 **Abstract**

29 The analysis and prediction of complex traits using microbiome data combined with host
30 genomic information is a topic of utmost interest. However, numerous questions remain
31 to be answered: How useful can the microbiome be for complex trait prediction? Are
32 microbiability estimates reliable? Can the underlying biological links between the host's
33 genome, microbiome, and the phenome be recovered? Here, we address these issues by
34 (i) developing a novel simulation strategy that uses real microbiome and genotype data
35 as input, and (ii) proposing a variance-component approach which, in the spirit of
36 mediation analyses, quantifies the proportion of phenotypic variance explained by
37 genome and microbiome, and dissects it into direct and indirect effects. The proposed
38 simulation approach can mimic a genetic link between the microbiome and SNP data via
39 a permutation procedure that retains the distributional properties of the data. Results
40 suggest that microbiome data could significantly improve phenotype prediction accuracy,
41 irrespective of whether some abundances are under direct genetic control by the host or
42 not. Overall, random-effects linear methods appear robust for variance components
43 estimation, despite the highly leptokurtic distribution of microbiota abundances.
44 Nevertheless, we observed that accuracy depends in part on the number of
45 microorganisms' taxa influencing the trait of interest. While we conclude that overall
46 genome-microbiome-links can be characterized via variance components, we are less
47 optimistic about the possibility of identifying the causative effects, i.e., individual SNPs
48 affecting abundances; power at this level would require much larger sample sizes than
49 the ones typically available for genome-microbiome-phenome data.

50

51 **Author summary**

52 The microbiome consists of the microorganisms that live in a particular environment,
53 including those in our organism. There is consistent evidence that these communities play
54 an important role in numerous traits of relevance, including disease susceptibility or feed
55 efficiency. Moreover, it has been shown that the microbiome can be relatively stable
56 throughout an individual's life and that is affected by the host genome. These reasons
57 have prompted numerous studies to determine whether and how the microbiome can be
58 used for prediction of complex phenotypes, either using microbiome alone or in
59 combination with host's genome data. However, numerous questions remain to be
60 answered such as the reliability of parameter estimates, or which is the underlying
61 relationship between microbiome, genome, and phenotype. The few available empirical

62 studies do not provide a clear answer to these problems. Here we address these issues by
63 developing a novel simulation strategy and we show that, although the microbiome can
64 significantly help in prediction, it will be difficult to retrieve the actual biological basis
65 of interactions between the microbiome and the trait.

66

67 **Introduction**

68 The relevance of microbial ecosystems associated with humans and animals in health and
69 production is now widely recognized, e.g., [1–5]. To quantify its influence, the fraction
70 of variance of a given trait explained by the microbiome has been named ‘microbiability’
71 (b^2) [6], in symmetry with the classical ‘heritability’ (h^2) concept [7]. Previously, the term
72 "hologenome" had been coined to describe the joint action of genome and microbiome in
73 explaining an observed phenotype [8].

74

75 A consequence of microbiability being typically larger than zero is that it can be used to
76 predict complex phenotypes, be it a disease or productive traits. This is an important issue
77 since the use of microbiome data has the potential to alter how medical diagnosis in
78 humans or breeding decisions agricultural species are performed. Several studies have
79 demonstrated the potential value of microbiome data for complex-trait prediction. For
80 example, Rothschild et al. [9] showed that microbiome can be used to improve accuracy
81 in the prediction of obesity and many other phenotypes in humans. Likewise, Lloyd-Price
82 et al. showed that microbiome-data was predicted if future outbursts of bowel disease
83 [10]. In cattle, various studies have shown the predictive power of microbiome for
84 methane emission from rumen microbiome [4,11], feed efficiency and carcass traits in
85 pigs [12,13], and various plant phenotypes (e.g., crop yield and diseases predicted from
86 the microbiota data from the rhizosphere, [14]). On the other hand, since the
87 groundbreaking study of Meuwissen, Hayes and Goddard [15], the prediction of complex
88 traits using genome information has been embraced in both plant [16] and animal
89 breeding [17] as well as in human genetics [18]. Therefore, a natural step further is
90 combining host’s genome and microbiome information to improve complex-trait
91 prediction, a topic that is currently receiving much attention [12,19].

92

93 Importantly, microbiome composition can be affected by the host’s genome. For instance,
94 Wang et al. [20] argue that it is evolutionarily justified that the microbiome is under
95 partial host genetic control since a non-negligible fraction of cells in an adult body is

96 made up of microbes, especially in the gut. Beginning with the seminal work by Pomp's
97 team [21], several studies have confirmed the relationship between host's genotype and
98 microbiome composition, e.g., [20,22,23]. These microbiome genome-wide association
99 studies (mGWAS) suggest that microbiome abundances can be treated as any other
100 complex trait in humans or livestock [22]. For instance, Crespo-Piazuelo et al. [24] or
101 Ramayo-Caldas *et al.*, [25,26] identified several quantitative trait loci (QTL) that
102 modulate gut bacterial and eukaryotic communities. In general, although the 'heritability'
103 of each genera or OTU (Operational Taxonomic Unit) is typically weak, considering the
104 whole microbiome simultaneously should increase power [27].

105

106 Large scale studies in humans suggest a predominant role of the environment in shaping
107 the gut microbiome [9]. However, regardless of the relative importance of genetic and
108 environmental factors in shaping the microbiota, microbiome composition *per se* can
109 have predictive value. Yet, the use of microbiota for prediction of future
110 phenotypes/disease outcomes, require some level of stability of the microbiome
111 throughout time. In the case of the gastrointestinal tract, microbiota colonization starts at
112 birth, where vertical transmission through the mother's birth canal occurs. Afterward,
113 microbiota diversity and richness tend to increase as the host ages and reaches stability at
114 adulthood [28,29]. In ruminants, populations inhabiting the rumen progressively appear
115 after birth and partly persists throughout life [30].

116

117 As noted, the genome-microbiome-phenome is a complex system; understanding the
118 links between host-genome, microbiota, and phenotypes is an important step towards the
119 effective use of microbiome data for complex trait prediction. In all, despite published
120 reports, we still lack detailed guidelines on the joint usage of microbiome and genome
121 information for complex trait prediction, and on the reliability of parameter inferences.
122 We are ignorant of the number of genes affecting microorganism abundance that can be
123 confidently identified, or on how many microorganism taxa can influence a given
124 phenotype. With this work, we aim to contribute to this important topic focusing on three
125 inter-related questions:

126

- 127 1. How useful can the microbiome be for complex trait prediction?
- 128 2. Are microbiability estimates reliable?

129 3. Can the underlying biological genome-microbiome-links be inferred at a system-
130 level? In a more refined level, can microbiome groups (e.g., OTUs, genera) with
131 sizable causal effects on phenotypes be identified with the typical size of
132 microbiome data sets?

133

134 In this study, we address the questions mentioned above via a novel simulation strategy
135 that uses real microbiome and genotype data as input and proposing a variance-
136 component approach which, in the spirit of mediation analyses, quantifies the proportion
137 of phenotypic variance explained by genome and microbiome, and dissects it into direct
138 and indirect effects. Importantly, the approach allows simulating a partial genetic control
139 of host's genome on the microbiome. This is accomplished using a partial permutation
140 approach that preserves the distribution of the genome and microbiome. We use Bayesian
141 variable selection models to estimate parameters which contemplate the possibility that
142 some or all the features available in the genome and/or the microbiome, have no effects
143 on the trait of interests. We investigate the questions presented above across diverse
144 scenarios regarding the links between host genomes and microbiomes, and of their
145 relations with a complex trait.

146

147 **Results and Discussion**

148

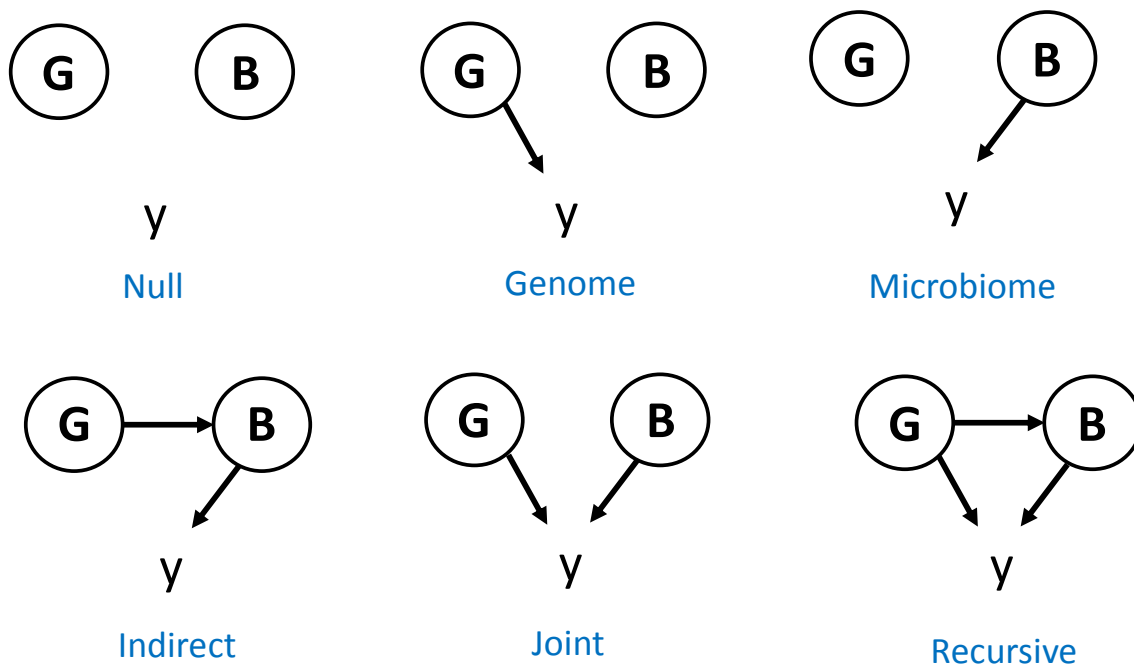
149 The exact nature of the links between genome (**G**), microbiome (**B**), and phenotype (**y**)
150 are largely unknown and will likely vary from case to case. However, we will use the six
151 generic causal models ('scenarios') depicted in Fig 1 to shed light on the nature of the
152 genome-microbiome-phenome links. In the 'Null' scenario, there is no link between any
153 of the data-layers; while this is unlikely, it serves as an 'overall null hypothesis' and it is
154 useful to assess potential biases in parameter estimates. Model 'Genome' assumes that **G**
155 only affects the phenotype. In turn, only **B** has a direct effect on phenotype in
156 'Microbiome' and 'Indirect' scenarios. The Indirect scenario, however, allows for some
157 of the causative abundances to be controlled genetically. This would be similar to a
158 scenario where a phenotype is directly controlled by gene expression levels and
159 expression in turn is controlled genetically [31,32]. The 'Joint' scenario is the simplest
160 configuration for a trait under the influence of both genes and microbiome. It assumes
161 microbiome and genome are independent and that their effects on the phenotype are also
162 independent. The Joint model is the most widely assumed, implicitly, or explicitly, in the

163 literature, e.g., [4,9,12]. The ‘Recursive’ model is similar to the Joint model; however,
164 the Recursive model contemplates the possibility that some causative OTU may be under
165 partial genetic control by the host. Therefore, in this case, the genome has both direct and
166 indirect (microbiome-mediated) effects on phenotypes. Note the Recursive model does
167 not assume that the same loci have simultaneously direct and indirect effects, neither it
168 assumes that all OTU abundances are under genetic control.

169

170

171



173

174 **Fig 1.** Representation of the scenarios evaluated: **G**, genome, typically comprises marker
175 data; **B**, microbiome; **y**, phenotype of interest; arrows indicate causality. An arrow from
176 **G** to **y** indicates that there is a subset of **G** elements (causative SNPs) that influence **y**; an
177 arrow from **G** to **B** indicates there exists a subset of **G** that influences a subset of
178 abundances in **B** which, in turn, may also influence **y**. An arrow departing from **B**
179 indicates there is a subset of microbial abundances (the causative abundances) that
180 influence **y**. The SNPs affecting **B** need not necessarily be the same SNPs affecting **y**
181 directly in the Recursive scenario. Note **B** can contain one or more sets of abundances
182 such as archaea and bacteria communities, or different time or site sampling points.
183 Without loss in generality, we assume **B** is a single community.

184

185 We use the causal models depicted in Fig 1 to simulate genome-microbiome-phenotype
186 data using different configurations regarding the number of causative loci (QTN) and the
number of OTUs with effects on phenotypes, as well as the number of OTUs that were

187 affected by host's genome. Table 1 summarizes the simulation models and parameter
188 values.

189

190 **Table 1.** Definition of scenarios evaluated and parameters chosen: **G**, genome; **B**,
191 microbiome; **y**, phenotype of interest; N_{QTN} , number of SNPs with a direct causal effect
192 on **y**; N_{OTU} , number of OTUs with a direct effect on **y**; $N_{OTU(g)}$, number of OTUs with a
193 direct effect on **y** that are genetically determined, i.e., they are a subset of N_{OTU} ; h^2 is
194 heritability, b^2 is microbiability, and $r^2 = h^2 + b^2$.

195

Scenario	Abbreviation	N_{QTN}	N_{OTU}	$N_{OTU(g)}$	r^2	h^2	b^2
Null	0	-	-	-	0	0	0
Joint	J	100	25	0	0.25	0.125	0.125
					0.50	0.25	0.25
Genome	G	100	0	0	0.25	0.25	0.00
					0.50	0.50	0.00
Microbiome	M	0	25	0	0.25	0.00	0.25
					0.50	0.00	0.25
Recursive	R	100	25	25	0.25	0.125	0.125
					0.50	0.25	0.25
Indirect	I	0	25	25	0.25	0.00	0.25
					0.50	0.00	0.50

196

197 **Table 2.** Scenarios used to evaluate sensitivity to the number of causative OTUs. Symbols
198 as in Table 1.

199

Scenario	Abbreviation	N_{QTN}	N_{OTU}	$N_{OTU(g)}$	r^2	h^2	b^2
Joint	J10	100	10	0	0.50	0.25	0.25
	J100	100	100	0	0.50	0.25	0.25
	J250	100	250	0	0.50	0.25	0.25
Recursive	R10	100	10	5	0.50	0.25	0.25
	R100	100	100	50	0.50	0.25	0.25
	R250	100	250	125	0.50	0.25	0.25

200

201

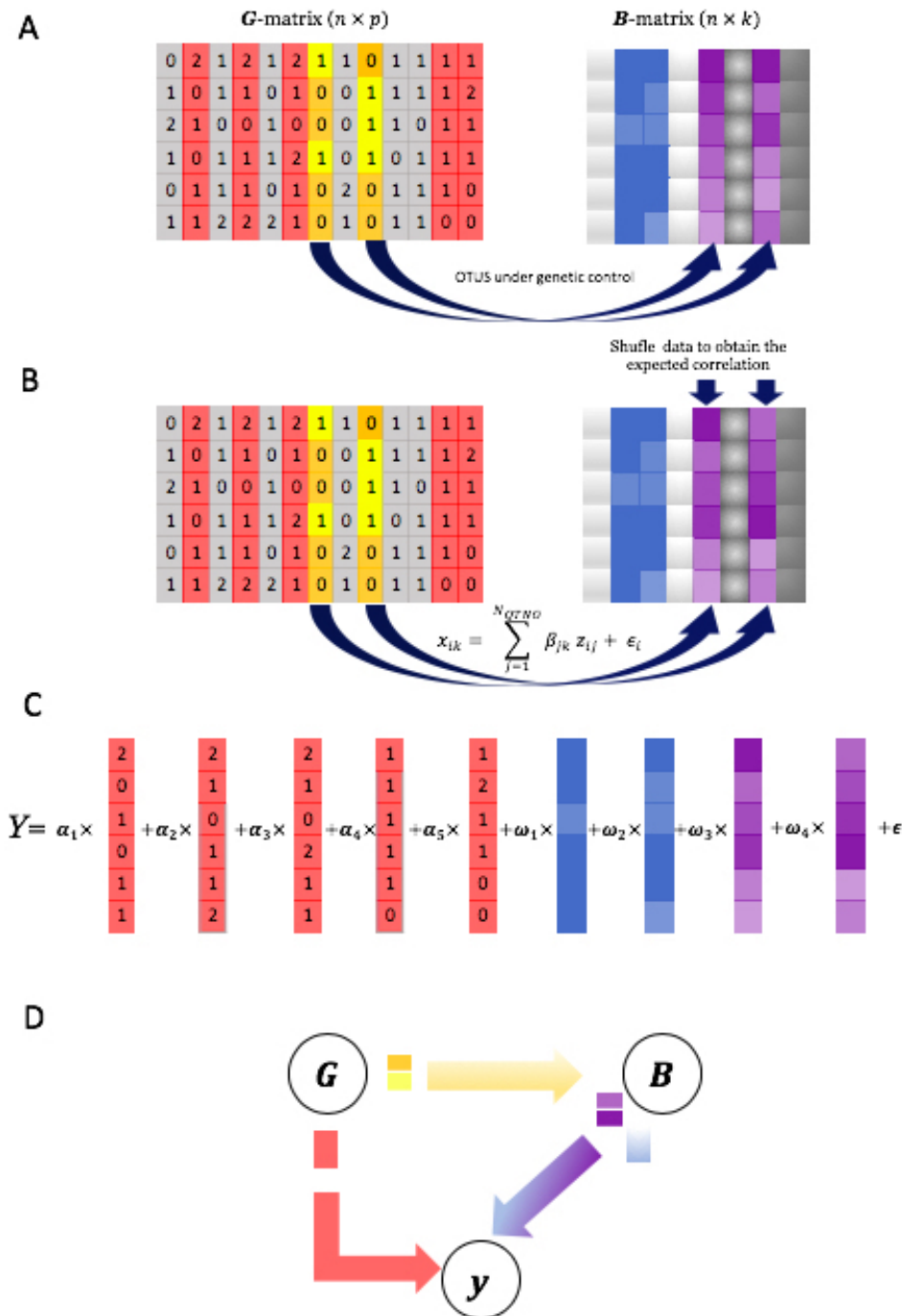
202 **A novel data-driven strategy to generate microbiome-genome-phenotype**
203 **experiments**

204 Two facts make it the simulation of scenarios in Fig 1 challenging: (i) microbiome data
205 follow zero-inflated highly leptokurtic multivariate distributions [33,34], it is not obvious
206 how to sample from these distributions *conditionally* on genome data as required in the
207 Recursive and Indirect scenarios; and (ii) it is difficult to obtain accurate estimates of key
208 parameters, such as microbiability values, in the absence of large scale published – and
209 public – datasets. To circumvent, or at least to alleviate, these constraints we use real data
210 for both **G** and **B**. Specifically, we used publicly available data from two of the largest
211 microbiome studies in livestock, genome data were downloaded from [11] and OTU
212 abundances from [4].

213

214 Fig 2 recapitulates the simulation strategy. Full details are given in Material and Methods
215 section, and R code to replicate the analyses are in
216 <https://github.com/miguelperezenciso/simubiome>). We assume the effects of the
217 causative microbiome abundances are additive on the log scale. Simulation under the
218 Joint scenario is straightforward, since **G** and **B** act independently: sample a list of
219 causative SNPs and abundances, simulate their effects, and apply Eqn. 1 (Material and
220 Methods) to generate phenotype values given observed genotypes and abundances. The
221 case of Recursive and Indirect scenarios is not that obvious because causative abundances
222 are under genetic control and a link must exist between **G** and **B** (Eqn. 2 in Material and
223 Methods). We solved this issue by rearranging abundances within individuals such that
224 the desired correlation between abundance and individual's genotypes is attained (see
225 Algorithm in Box 1 and R-code in
226 <https://github.com/miguelperezenciso/Simubiome/blob/master/sortCor.R>). This strategy
227 has the important advantage that the distribution of abundances is not changed.

228



229

230 **Fig 2:** Simulation scheme for the Recursive scenario, i.e., the most complex scenario (Fig
 231 1). **A)** Real input data comprises p genotypes (**G** matrix) and k taxa abundances (**B**
 232 matrix). SNPs in grey are neutral, those in red act directly on the phenotype y , and those
 233 in yellow/orange influence some OTU abundances (marked in magenta color in **B**
 234 matrix); abundances in blue are not genetically controlled. **B)** Given simulated effects, a
 235 genotypic value controlling the abundances is obtained via Eqn. 2. To fulfill the required
 236 heritability, abundances in magenta are reordered; high abundances (represented by a
 237 darker color) are associated with genotype '1' just to simplify visualization. A single SNP
 238 is shown as causative for each of the two OTUs but there is no limit in practice. **C)** The
 239 phenotype is simulated by adding the genome and the microbiome contributions plus a
 240 residual. **D)** The general causal diagram is shown.

241

242 **How useful can microbiome be for complex trait prediction?**

243 This depends on how much phenotypic variance is jointly explained by the genome (h^2)
244 and the microbiome (b^2), but also on how efficiently methods capture the relationship
245 between the microbiome and the phenotype, and on how stable the microbiome is. Note
246 prediction accuracy is conditionally independent of whether the microbiome itself is
247 heritable or not. This means that, *given* observed abundances **B** and observed genotypes
248 **G**, it does not matter whether the biological process generating **B** is affected by **G**. In
249 other words, prediction should not be affected by whether the Joint or Recursive scenarios
250 hold, for a constant $r^2 = h^2 + b^2$. The implications for breeding, however, could be
251 dramatically different. Breeding schemes targeting the microbiome could be designed
252 provided the Recursive scenario holds but make no sense under the Joint scenario.

253

254 We compared predictive performance of Bayes C [15] when both genome and
255 microbiome are employed in the model (*Bayes Cgb*) only genome (*Bayes Cg*), or only
256 microbiome data (*Bayes Cb*). First, we verified the null model resulted in no false
257 predictive accuracies (Fig S1A). Fig 3 shows simulated predictive accuracies for the two
258 r^2 values considered (0.25 and 0.50) and for each causative scenario (Fig 1). Predictive
259 accuracies using *Bayes Cgb* were consistently the best. As expected, this was especially
260 the case when both h^2 and b^2 are larger than zero, that is, when Joint or Recursive scenario
261 hold. In these scenarios, using both sources of variation clearly improved prediction
262 compared to using only genome (*Bayes Cg*) or microbiome data (*Bayes Cb*). Importantly,
263 predictive accuracy was somewhat lower in Joint and Recursive scenarios than in
264 Microbiome or Genome scenarios. This indicates that predictive accuracy does not
265 depend only on total r^2 , but also on how this variance is split between genome and
266 microbiome. Although this likely occurs because of the larger noise in Recursive or Joint
267 scenarios than in Microbiome or Genome scenarios, it also suggests that our analysis
268 strategy may not be optimum. There is room to develop more efficient tools, especially
269 when the Recursive scenario holds. Note that variance of prediction was larger in the
270 Recursive than in the Joint scenario, i.e., the fact that some abundances are inherited is
271 an additional source of noise.

272

273 It is noticeable that predictions were better when only the microbiome influenced the
274 phenotype than when the genome was the only source of variation, a phenomenon also

275 observed with real data [11,12,19]. In this simulation, this occurs likely because the
276 number of causative effects and of input variables (SNPs vs. OTUS) is smaller in the
277 Microbiome or Indirect scenarios than in the Genome scenario. In fact, we do observe a
278 consistent negative correlation between the number of causative OTUs and predictive
279 accuracy in both Joint and Recursive scenarios (Fig 4A).

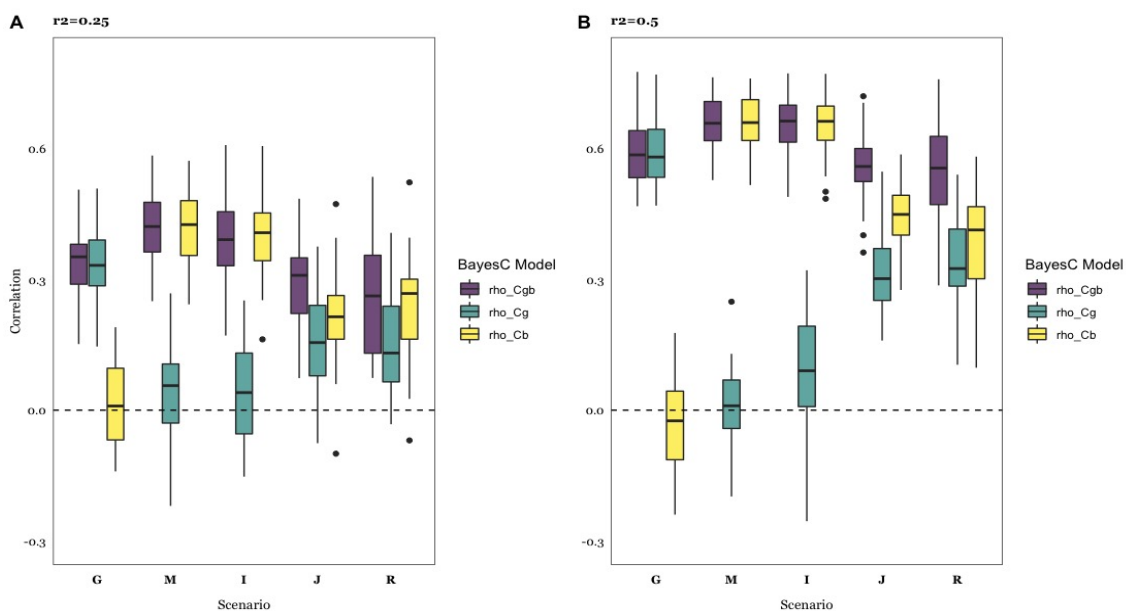
280

281 In all, our results suggest that predictive accuracy could be increased by ~ 50 % when
282 considering microbiome data, provided microbiability is of the same order as heritability
283 (Fig 3). We speculate that this is probably an upper limit, since it will be difficult to have
284 microbiome data collected homogeneously across time and in different locations. While
285 individuals can be genotyped at birth, the microbiome in early life is not representative
286 of adult or later stages. Maltecca et al., for instance, show that early life microbiota is not
287 a good proxy for carcass composition in pigs [35].

288

289 We observed, roughly, a two-fold increase in predictive accuracy when doubling
290 heritability for Genome, Joint and Recursive scenarios, and a 50% increase for
291 Microbiome or Indirect scenarios (Fig 3A vs. 3B).

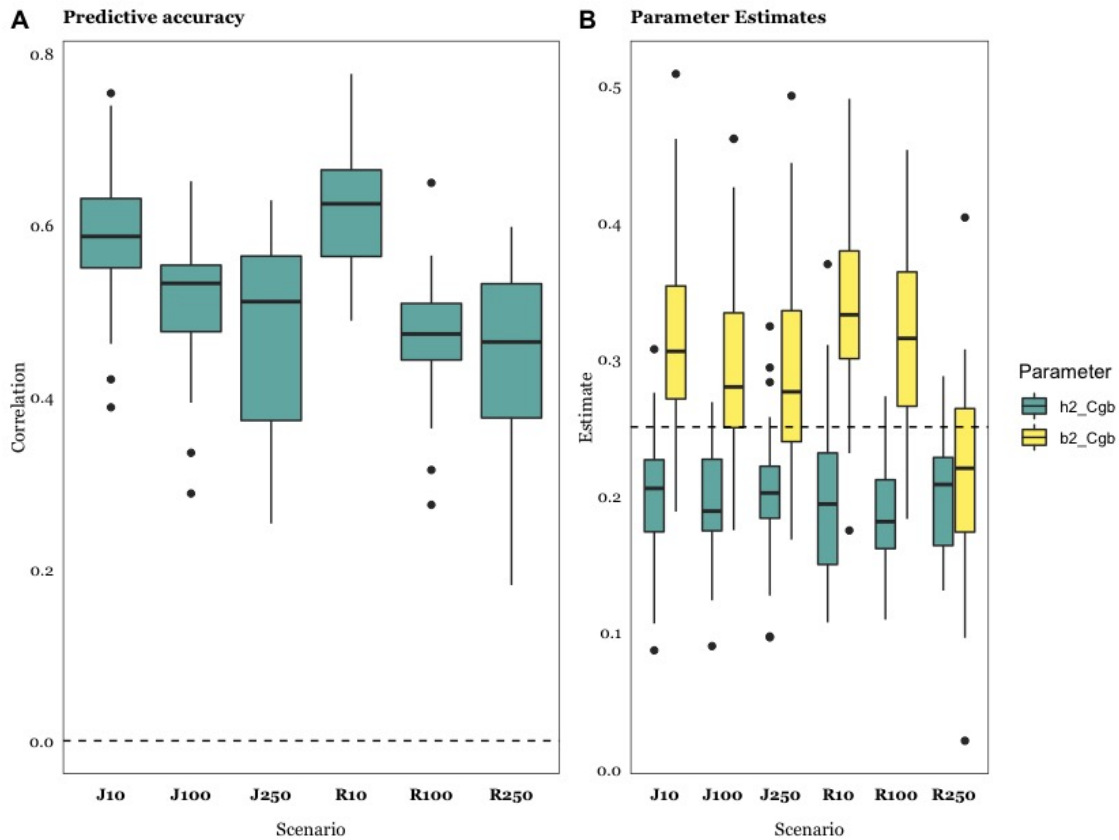
292



293

294 **Fig 3.** Predictive accuracy, computed as correlation between predicted and observed
295 phenotypes across causal scenarios (Fig 1), for each of the Bayes C models: Cgb
296 considers microbiome and genome; Cg includes genome data only, and Cb includes
297 microbiome data only. **A:** $r^2 = 0.25$; **B:** $r^2= 0.50$. Details of scenarios are in Table 1: G,
298 Genome; M, Microbiome; I, Indirect; J, Joint; R, Recursive. Results are average of 30
299 replicates per case.

300



301

302 **Fig 4.** Effect of varying number of causative OTUs, $r^2 = 0.5$. **A:** Predictive accuracy,
303 computed as correlation between predicted and observed phenotypes, using Bayes Cgb.
304 **B:** Heritability and microbiability estimates using Bayes Cgb. Details of scenarios are in
305 Table 2 and diagrams in Fig 1: Jx, Joint scenario; Rx, Recursive scenario, with x being
306 the number of causative OTUs ($x = 10, 100, 250$), half of them under genetic control.
307 Results shown are the average of 30 replicates.

308

309 Are microbiability estimates reliable?

310 Reliable parameter estimates are needed to optimize the design of breeding schemes or
311 microbiome wide association studies (MWAS) [36]. They are also needed for
312 understanding the biology behind the interaction of microbiome and complex phenotypes.
313 Thus far, microbiability has been usually estimated using 'standard' linear methods, e.g.,
314 [4,9,27], much as we have done here. It is of interest then to know how accurate these
315 estimates could be.

316

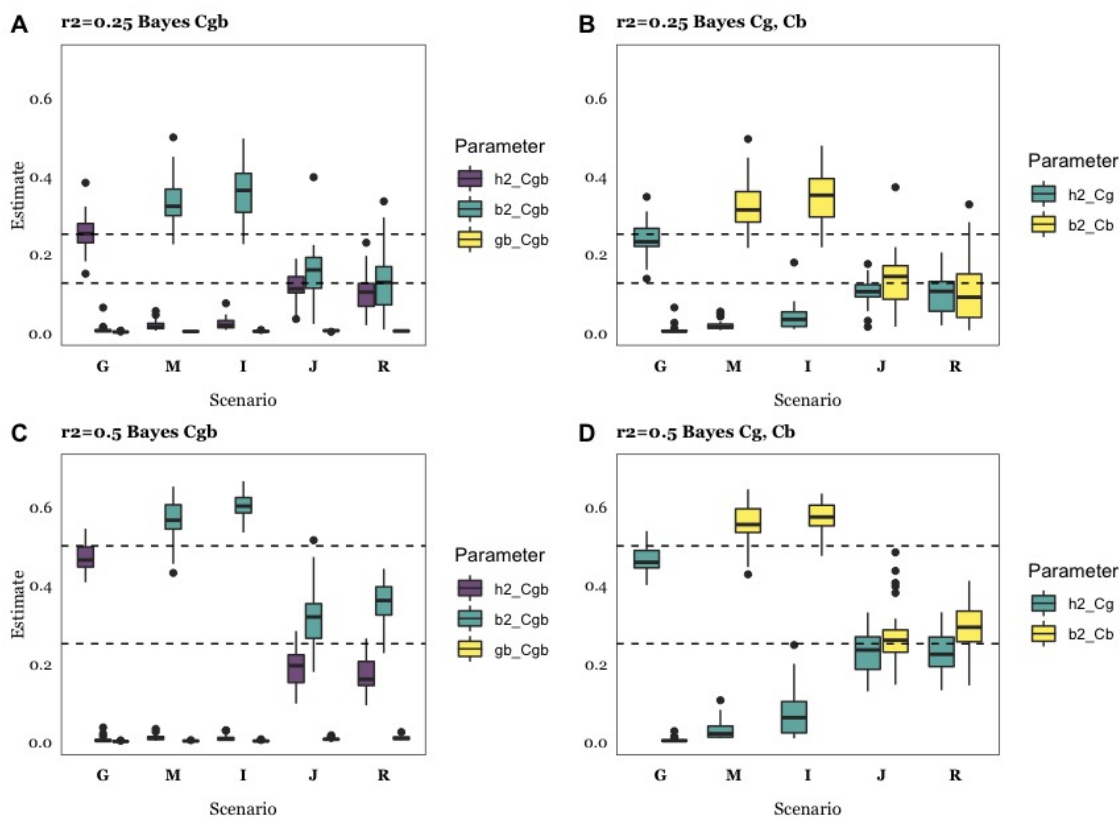
317 Fig 4 shows estimates of variance components for each of the scenarios in Table 1. Bayes
318 Cgb allows us to assess whether h^2 and/or b^2 are different from zero: microbiability
319 estimate is near zero when the data are simulated according to the Genome scenario and
320 heritability is zero when the Indirect or Microbiome scenarios hold, as it should.

321 Similarly, both h^2 and b^2 estimates are near zero when the null scenario holds (Fig S1B).
322 An overestimation of b^2 is nevertheless evident in Fig 4, and it does not vanish at higher
323 r^2 . This upward bias in b^2 estimate is accompanied by an underestimation of h^2 , indicating
324 that variance estimates are confounded when using *Bayes Cgb* model. This bias decreases
325 though when the number of causative OTUs increases. For instance, the bias in b^2
326 estimate is $\sim 40\%$ when $N_{\text{OTU}} = 10$ but is reduced to $\sim 10\%$ with $N_{\text{OTU}} = 250$ (Fig 4B).
327 Therefore, it is likely that the presence of a few causative OTUs, but of large effect,
328 combined with the presence of highly leptokurtic abundance distributions, may result in
329 biased parameter estimates. This should be considered when interpreting microbiability
330 estimates in real experiments. For instance, Difford et al. [4] report estimates $h^2 = 0.21$
331 and $b^2 = 0.13$ ($N = 750$), finding **G** and **B** to behave independently. Assuming the number
332 of causative OTUs is small compared to that of SNPs with an effect on abundances
333 (QTNs), we can presume Difford's estimate of b^2 to be inflated. This means that the actual
334 microbiome contribution may be too small to improve prediction over that obtained from
335 using marker data exclusively. Although authors focused on inference and not so much
336 in prediction, Difford et al reported that no bacteria genera were significantly associated
337 with methane emissions [4]. Other authors in turn have reported polymicrobial
338 associations, including members of bacterial, archaeal, fungal, and protozoan
339 communities, with methane emissions, e.g., [11,25,37–39].

340

341 For comparison, Fig 5B,D show the estimates obtained with *Bayes Cg*, when only h^2 is
342 estimated, or *Bayes Cb*, only b^2 is estimated. The most noticeable aspect is that bias in b^2
343 estimates is somewhat reduced relative to that found with *Bayes Cgb*, signaling again
344 some confounding between b^2 and h^2 . Bias was reduced overall at higher r^2 but did not
345 vanish.

346



347

348 **Fig 5:** Estimates of heritability (h^2), microbiability (b^2), and correlation between genome
 349 and microbiome (gb) for each of the three Bayes C analysis models: Cgb includes
 350 microbiome and genome in the model (left panels); Cg includes genome only, and Cb
 351 includes microbiome data only (right panels). Upper rows correspond to $r^2 = 0.25$ and
 352 lower rows to $r^2 = 0.50$. Details of simulation scenarios are in Table 1: G, Genome; M,
 353 Microbiome; I, Indirect; J, Joint; R, Recursive. Horizontal dashed lines indicate true h^2
 354 or b^2 parameter values (0.125, 0.25, 0.5 depending on the scenario and on r^2). Results
 355 are average of 30 replicates. **A:** $r^2=0.25$, Bayes Cgb estimates (h^2 , b^2 and gb); **B:** $r^2 =$
 356 0.25, Bayes Cg (h^2) and Cb (b^2) estimates; **C:** $r^2 = 0.50$, Bayes Cgb estimates; **D:** $r^2 =$
 357 0.50, Bayes Cg and Cb estimates. Data are average of 30 replicates per case.

358

359 **Can the underlying biological scenario be recovered? Can causative OTUs be**
 360 **identified?**

361 An important goal of many experiments is to dissect the biological basis of microbiome
 362 and genome interactions, even if this is not strictly needed for prediction. So far, our
 363 simulations suggest that standard statistical methods can be used to quantify – with some
 364 bias – microbiability contribution to phenotypic variance. It also seems feasible to
 365 distinguish whether the Microbiome or Genome scenario fits real data best. Similarly, it
 366 seems plausible to assess when **G** and **B** contribute to the phenotypic variance, i.e. when
 367 Recursive or Joint scenarios are plausible.

368

369 Could Joint vs. Recursive scenarios be distinguished? Can data point to which of Indirect
370 or Microbiome scenarios is more plausible, if any? Further, can causative OTUs be
371 identified? These are far more difficult questions to answer than assessing prediction
372 performance or estimating microbiability. Compare variance component estimates
373 obtained under the Joint or Recursive scenarios (Fig 5): they are nearly identical for the
374 same r^2 . The two scenarios differ in that at least some causative OTUs abundances can
375 be under partial genetic control in the Recursive scenario. The Recursive scenario should
376 result in a covariance between \mathbf{G} and \mathbf{B} . We conjectured that the two scenarios could be
377 distinguished by analyzing the covariance $Cov(\mathbf{u}^{(i)}, \mathbf{v}^{(i)}) / Var(\mathbf{y})$ (see methods).
378 Unfortunately, these estimates are close to zero irrespective of the true scenario (Fig 5A,
379 C). The likely reason is that the actual fraction of phenotypic variance explained by
380 indirect effects is *conditionally* negligible. Note there can be a genetic effect of \mathbf{G} on \mathbf{B}
381 but, for our purposes, we are interested only in those genes that affect causative OTUs
382 (i.e., those that affect the phenotype) and not on the whole microbial system.

383

384 An alternative approach to infer whether the Recursive causative scenario holds or not is
385 to run a genome-wide association study (GWAS) for each of the OTU abundances on
386 each SNP, where the SNP P-values can indicate a genetic basis for some of the
387 abundances. If we identify significant SNPs for OTUs likely influencing \mathbf{y} , we could
388 conclude that the Recursive scenario is plausible. Unfortunately, this analysis can be
389 doomed by the large number of tests to be realized, i.e., $N_{OTU} \times N_{SNP}$. To illustrate the
390 caveats of GWAS on abundances, Fig 6A shows the distribution of $-\log_{10}$ P-values of
391 neutral SNPs vs. SNPs with an effect on abundances. Assume we take the 5% empirical
392 threshold of the neutral P-value distribution as indicative of association. Simulations
393 suggest that only $\sim 3\%$ of causative SNP P-values will be above that threshold, i.e.,
394 approximately what is expected by chance. These P-values depend of course on the actual
395 number of causative SNPs and on abundance heritabilities, but most evidence so far
396 points to a weak relationship between genome and microbiome [22]. We warn it is going
397 to be very difficult to identify abundance causative SNPs using GWAS information alone
398 [9,20].

399

400 Another question of interest is how many of the OTUs affecting the phenotype can we
401 expect to discover. One option is to count the frequency of a given OTU entering into the
402 Bayes C model during sampling. Fig 6B shows the probability of including a causative

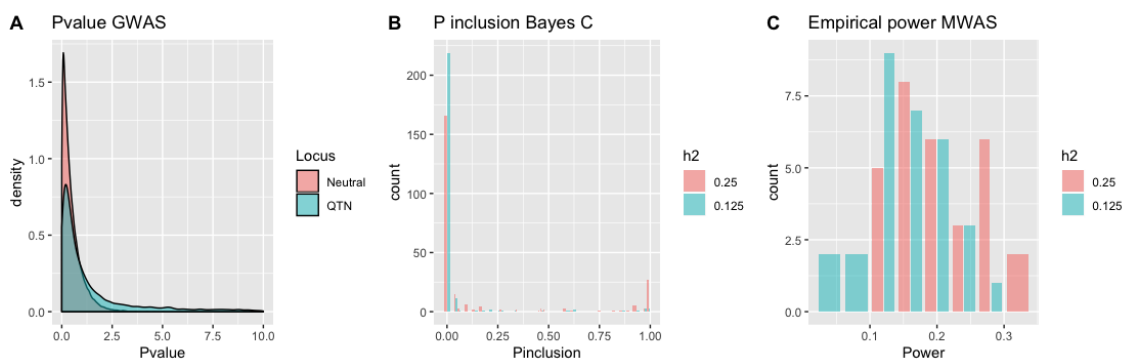
403 OTU in the Bayes C sampling chain, which varied between ~5% ($b^2 = 0.125$) to ~20%
404 ($b^2 = 0.25$). About 50% ($b^2 = 0.25$) or 30% ($b^2 = 0.125$) of causative OTUs were among
405 the 5% most frequently included OTUs in the Bayes C chain, on average. Since the
406 number of causative OTUs was 25, the rate of false positives was high nevertheless. We
407 can conjecture that only a few causative OTUs are likely to be identified in medium-sized
408 experiments, such as this one.

409

410 An alternative approach is a Microbiome Wide Association Study (MWAS), i.e., to
411 perform a linear regression of the phenotype on each of OTU abundances and then select
412 the significant results as potential causative OTUs [4]. Fig 6C shows the average power,
413 defined as the percentage of true causative OTUs within the 5% most significant results.
414 Power was ~15% and ~20% for $b^2 = 0.125$ and 0.25, respectively, in the Recursive
415 scenario. Again, this is not too satisfactory, as we expect a high fraction of false positives.
416 In this particular scenario, it is perhaps more useful to consider probabilities of inclusion
417 in the Bayes C chain rather than at P-values since the former are the result of a joint
418 analysis of all OTUs and can be used directly for prediction.

419

420



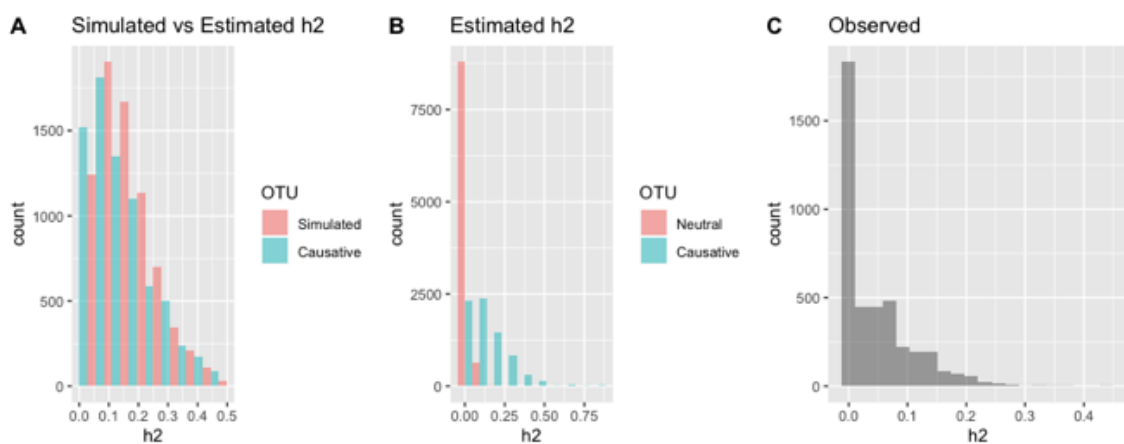
421

422 **Fig 6:** *A) Distribution of $-\log_{10}$ P-values of a GWAS of abundances on SNP data; B)*
423 *Probability of inclusion in the Bayes Cgb model of causative OTUs for the two levels of*
424 *microbiability considered. C) Power of identifying a causative OTU computed as the*
425 *probability of exceeding the 95% threshold of the empirical distribution of P-values in*
426 *an MWAS for the Recursive scenario.*

427

428 Finally, we investigated the pattern of abundance heritabilities. Fig 7A shows the
429 simulated heritabilities for the causative, inherited OTUs, which approximately follows a
430 gamma distribution, together with estimated heritabilities for the causative OTUs in the
431 Recursive scenario. We observe that both distributions are rather similar although
432 estimates are somewhat shrunk towards zero, a consequence of using a REML-like prior.

433 A problem of course is that we do not know which OTUs are inherited and which are not,
434 and the true distribution of OTU heritability estimates will be a mixture. Fig 7B illustrates
435 the heritability distributions of neutral (non-inherited) and causative (inherited) OTUs. In
436 Fig 7B, we mixed 1.7 neutral OTU per causative OTU. This is completely arbitrary since
437 we do not know the actual number of OTUs under genetic control, but we did so because
438 the resulting mixture is similar to the distribution of heritabilities observed by Difford et
439 al. (Fig 7C). If distributions in Fig 7B were representative of the true state of nature, this
440 would suggest that about $1/(1+1.7) \sim 40\%$ rumen OTUs could show some genetic additive
441 variance in the experiment reported by Difford et al.[4].
442



443 **Fig 7: A:** ‘True’ (simulated) and GBLUP estimated distribution of abundance
444 heritabilities for causative OTUs in the Recursive scenario. **B:** GBLUP estimated
445 distribution of abundance heritabilities for neutral and causative OTUs in the Recursive
446 scenario. **C:** Actual distribution of OTU abundance heritabilities reported by Difford et
447 al.[4].
448
449

450 Discussion

451 Fig 1 represents but highly simplified relationships between the genome, microbiome,
452 and phenotype. These scenarios are nevertheless important to interpret empirical data and
453 can help to identify limiting factors in prediction. Further, provided a good fit is found,
454 they will help in designing experiments that combine microbiome and genetic data. We
455 chose parameter combinations that represent extreme case scenarios and we found that
456 results were, qualitatively, robust to parameter choice such as r^2 . A parameter that can be
457 relevant though is the number of causative microbiome taxa, i.e., those with an effect on
458 the phenotype. This number seems to affect the bias of microbiability estimates (Fig 4).
459

460 Here, we have proposed a new simulation procedure that addresses some important
461 challenges. First, the algorithm avoids the need for actual phenotype simulation by using
462 real genotype and abundance data. Although we concede that this procedure may limit
463 the generality of the study, e.g., in terms of data size, we believe the advantages of using
464 real data are numerous, since no simulation procedure can accommodate all known and
465 unknown subtleties of the highly dimensional distributions at hand. Second, we develop
466 an ingenious permutation procedure (Box 1) that allows linking previously uncorrelated
467 data to fit a desired genetic hypothesis. By also permuting all OTUs within a given cluster,
468 we minimize disruption of the whole covariance structure (Fig S2).

469

470 Numerous studies have reported microbiability values for economically important traits
471 e.g. [4,12,25,39], but their actual reliability is not known. Estimates may be affected by
472 the estimation procedure. There are numerous alternatives to estimate b^2 , among them
473 Bayes C [15], GBLUP [40], Bayesian RKHS regression using either Bray–Curtis
474 dissimilarities as relationship matrix [25] or with the variance-covariance from the
475 log-transformed OTUs as kinship matrix[25,41]. Our results (Fig 5) indicate that BayesC
476 estimates may be biased upwards, especially when b^2 is higher than 0.25 and the number
477 of causative OTUs is small. However, we found that estimates of b^2 derived with Bayes
478 C were very close to zero in the null scenario (Fig S1B); therefore, we conclude that
479 models using priors from the Spike-Slab family, which contemplate a priori the
480 possibility of null effects, can be used to test whether heritability or microbiability is
481 substantial. Ramayo-Caldas et al. [25] report that estimates using Bray-Curtis based
482 kernels are higher than those using the log-transformed covariance matrix. The behavior
483 of estimation methods for microbiability merits further research.

484

485 One conclusion from this work is that it is going to be difficult to distinguish between
486 some underlying scenarios or to identify the causative OTUs and SNPs, at least using
487 standard linear models as was done here. The distinction between Joint and Recursive
488 scenarios is of special relevance for breeding. The latter assumes partial genetic control
489 of some causative OTUs. Yet, we found both scenarios result in very similar patterns
490 (Figs 3, 4, 5). Perhaps, a more powerful approach would be to use structural equation
491 models (SEM), which allow including a variable both as independent and dependent.
492 Saborio-Montero et al. [42] compared a linear bivariate (one OTU and the phenotype)
493 model with a SEM but found few differences. One restriction of their approach is that one

494 SEM was fitted for each abundance. A whole-genome approach seems in principle more
495 adequate; however, modeling recursive effects in this context is both statistically and
496 computationally challenging because of the large number of SNP-OTU combinations that
497 would need to be considered.

498

499 A line of research that we have not considered involves possible microbiome-DNA
500 interactions. Although the number of possible interactions to consider can be huge when
501 the number of SNPs and the number of OTUs is large, interactions between features in
502 two high-dimensional sets can be modeled in a Gaussian context using co-variance
503 functions. These functions are the Hadamard product of set-specific similarity matrices
504 such as the Hadamard product of a SNP-derived and an OTU-derived ‘relationship’
505 matrix. Such an approach has been used before to model, e.g. interactions between SNPs
506 or between SNPs and environmental covariates (e.g., 43).

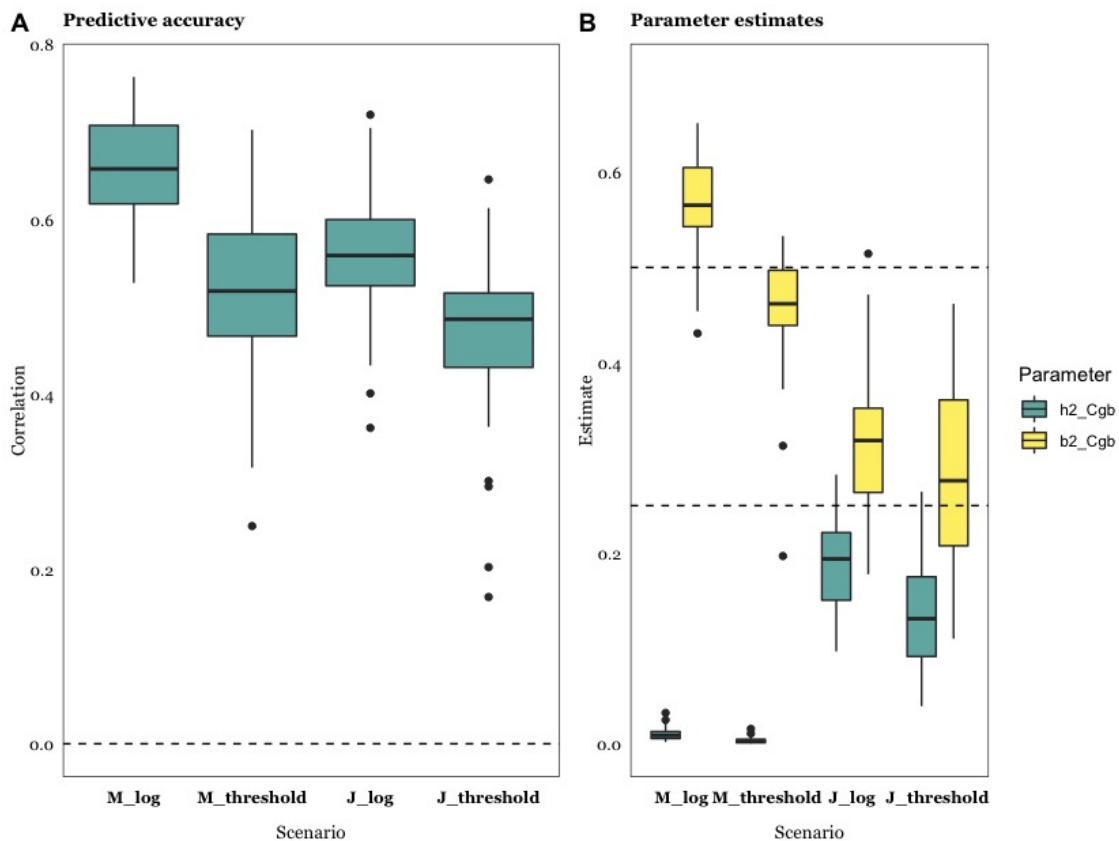
507

508 The usefulness of microbiome in prediction depends crucially on its stability in time and
509 space. For instance, although measures of gastrointestinal microbiome abundances are
510 known to be repeatable, it cannot be expected to remain stable throughout an individual’s
511 life span. After weaning and under standard management conditions, e.g., constant diet
512 and absence of antibiotic treatment, the diversity of monogastric gut microbiota increases
513 with host age until its composition remains stable. Rumen microbial communities are
514 highly resilient and host-specific [44,45] but change in early life. The transition towards
515 a more stable an adult-like ruminal ecosystem occurs between weaning and one year of
516 age [46]. Therefore, for prediction purposes, we recommend the inclusion of microbial
517 information at least after weaning, preferably at adulthood. This may limit the usefulness
518 of microbiota for prediction in breeding schemes as compared to genomic data in
519 livestock.

520

521 At present, modeling the influence of microbiome abundances on complex phenotypes is
522 an open area of research. Here we have presumed that the effects on abundances are
523 additive in the log scale. Similar models are widely used in a diversity of scenarios, e.g.,
524 multiplicative models are used to accommodate fitness effects in evolutionary genetics
525 [47] or to deal with highly leptokurtic distributions such as raw abundances, an effect that
526 is smoothed with the log transformation. In addition to the log-transformation, a widely
527 popular choice in genetics is the threshold model [7], which assumes the presence of a

528 continuous liability (here abundances) with an effect value ‘0’ below a given threshold
529 and ‘1’ otherwise. This model has the advantage of being independent on whether
530 abundances are log-transformed or not and is also biologically sound since it is
531 conceivable that a minimum microorganism abundance is required to trigger a particular
532 effect. To test the robustness of the log-transformation, we simulated phenotypes such
533 that 25% of causative abundance observations were above the threshold and the analysis
534 was performed on the log transformed abundances as before. As could be expected, using
535 a ‘wrong’ model for the analyses was detrimental to prediction but not dramatically (Fig
536 8A). Parameter estimates were affected downwards compared to the multiplicative model
537 (Fig 8B). We suggest that major conclusions from this work should hold even if the
538 relationship between variables and phenotype is not strictly multiplicative.
539



540
541 **Fig 8.** Comparison of multiplicative (log) and threshold Microbiome (M) and Joint (J)
542 scenarios ($r^2 = 0.5$). **A:** Predictive accuracy, computed as correlation between predicted
543 and observed phenotypes, using Bayes Cgb. **B:** Heritability (h^2) and microbiability (b^2)
544 estimates using Bayes Cgb. Results are average of 30 replicates. Scenarios M and J as
545 specified in Table 1; the log transformation results are shown for completeness and are
546 the same as in Figs 3 and 5. Data are average of 30 replicates.
547

548 Conclusion

549 This study suggests that microbiome data can significantly improve the prediction of
550 complex phenotypes, irrespective of whether some abundances are under direct genetic
551 control or not. For this strategy to be successful, though, medium to large-sized
552 experiments are required, the microbiome should be relatively stable and should be
553 available before the phenotype is collected. This limits the usefulness of microbiome for
554 prediction in breeding schemes as compared to genome data, which can be collected at
555 birth and remains unchanged. Important potential applications remain nevertheless, such
556 as predicting methane emission in cattle, feed efficiency, disease predisposition, or crop
557 production using soil metagenome. Overall, we can be rather confident that standard
558 linear methods can be used despite the highly leptokurtic distributions observed in OTU
559 abundances. There is room for specific theoretical developments though, perhaps along
560 the lines proposed by Saborio-Montero et al. [42], but these should be based on a better
561 understanding of the relation between microbiome and phenotype. It seems critical to
562 quantify, even approximately, the number of taxa affecting the phenotype and to
563 characterize the distribution of their effects. We are far less optimistic in what regards the
564 identification of causative OTUs, and in particular of the putative QTNs affecting relative
565 abundances.

566

567 **Materials and Methods**

568 **Simulation Strategy**

569 There is ample literature and software available on the simulation of ‘standard’ complex
570 phenotypes, e.g., [48–51]. These algorithms, however, are not suited for some of the
571 scenarios posed in Fig 1. Here we propose simulating the joint influence of genome and
572 microbiome on a quantitative trait by adding their contributions plus a random noise:

573

574

$$575 \quad y_i = \sum_{j=1}^{N_{QTN}} \alpha_j z_{ij} + \sum_{k=1}^{N_{OTU}} \omega_k x_{ik} + \varepsilon_i, \quad (1)$$

576

577

578 where y_i is the i -th individual record, α_j is the genetic effect of j -th causal SNP (QTN),
579 with $j = 1, N_{QTN}$, the number of QTNs, z_{ij} is the genotype of the i -th individual for j -th
580 SNP coded say -1, 0 and 1 (strict additivity was assumed for all QTN), ω_k is the linear
581 effect of the k -th OTU abundance (x_{ik}), with $k = 1, N_{OTU}$, the number of abundances that

582 influence the phenotype and ϵ is a normally distributed residual. The OTU's coefficient
583 can be interpreted as the expected change in phenotype per OTU's abundance unit
584 increase. Since abundances are in the log scale, this is equivalent to a multiplicative effect
585 model. Equation (1) is valid for all scenarios in Fig 1, except that the term involving
586 markers $\sum_{j=1}^{N_{QTN}} \alpha_j z_{ij}$ is removed in the Microbiome and Indirect scenarios whereas the
587 term $\sum_{k=1}^{N_{OTU}} \omega_k x_{ik}$ is removed in the Genome scenario.

588

589 For the Indirect and Recursive scenarios, we also need to model the variation in
590 abundances (\mathbf{x}) that is explained by the genome (Fig 1). Again, we can resort to a linear
591 model where the abundance itself is treated as a standard complex phenotype:

592

$$593 \quad x_{ik} = \sum_{j=1}^{N_{QTN(k)}} \beta_{jk} z_{ij} + \epsilon_i, \quad (2)$$

594

595 where x_{ik} is the abundance level of the k-th OTU that is under partial genetic control for
596 i-th individual, β_j is the genetic effect of j-th QTN on abundance, and z_{ij} is the genotype
597 of the i-th individual for j-th SNP. The j-th sum is across the QTNs influencing k-th
598 abundance, $j = 1, N_{QTN(k)}$. Note abundances x_{ik} in Eqn. (2) are a subset of those in (1).
599 There may be other non-causative abundances under genetic control, but this is irrelevant
600 for our purposes. A phenotype following the Recursive scenario can then be simulated
601 via a two-step procedure: first, simulate abundances (\mathbf{x}) using Eqn. (2) followed by
602 phenotype simulation using (1) given the abundances obtained.

603

604 We used real genome and microbiome data as input for the simulation procedure. We
605 downloaded the rumen abundance table of 4,018 OTUs from dairy cattle rumen ($N = 750$,
606 [4]). A pseudo-count equal to one was added to zero abundances, which were next total-
607 sum scaled and log-transformed. This results in much less leptokurtic and less asymmetric
608 distributions than original raw abundances. In Eqns. 1 and 2, x_{ik} represent the already log-
609 transformed abundances. As for genotypes, high-density array genotypes from 750 dairy
610 cows among the total available were downloaded from [11]. To prune SNPs and facilitate
611 computation, 35% of all genotypes with a minimum allele frequency of 0.01 and a
612 maximum missing percentage of 1% were retained. A total of 32,204 autosomal SNPs

613 was finally retained. The few missing values were simply imputed with the mean. Thirty
614 simulation replicates per scenario were simulated.

615

616 Under the Joint scenario, which assumes independence between **G** and **B**, we can simply
617 sample the list of causative SNPs and abundances, simulate their effects, and apply Eqn.
618 1 to generate phenotype values given observed genotypes and abundances. The case of
619 Recursive and Indirect scenarios is not that obvious because we need to sample
620 abundances that are under genetic control and a link must exist between **G** and **B** (Eqn.
621 2). We solved this issue by rearranging abundances of a given OTU between individuals
622 such that the desired correlation between abundance and individual's genotypes is
623 attained. This strategy has the important advantage that the distribution of abundances is
624 not changed. Suppose $\gamma_{ik} = \sum_{j=1}^{N_{QTN0}} \beta_{jk} Z_{ij}$ is the simulated genetic effect of the *i*-th
625 individual for *k*-th abundance (Eqn. 2) and that the desired heritability for that abundance
626 is h_k^2 . The algorithm (Box 1) is based on the simple observation that, given any two
627 vectors **x** and **y**, correlation is maximum ($\rho \sim 1$) when observations in both vectors are
628 sorted and ρ is \sim zero when they are shuffled. Therefore, there must be some order \mathbf{y}_{sort}
629 then that fulfills, approximately, the constraint $\text{cor}(\mathbf{x}, \mathbf{y}_{sort}) = \rho$. For our purposes, we need
630 to rearrange the observed abundances \mathbf{x}_k such that the correlation between rearranged \mathbf{x}_k
631 and γ_k is h_k , the square root of heritability for *k*-th abundance. The algorithm is detailed
632 in the Box.

633

634 A drawback of this algorithm is that it locally breaks the covariance between abundances
635 of different OTUs. To alleviate this, we permuted all abundances that fell within the same
636 OTU cluster. We clustered abundances using R function `hclust(dist(.),`
637 `method="ward.D2")` and cut the tree in $K = 500$ clusters. We chose $K = 500$ because the
638 first quartile of intra-cluster average correlation was above the third quartile of the
639 average correlation between random abundances, that is, clusters were made up of highly
640 correlated abundances compared to average. We also explored $K = 200$ but we did not
641 find any difference neither in predictive accuracy nor in heritability estimates. To verify
642 that the shuffling algorithm did not alter the whole structure of the data, we show the
643 principal component analysis of the original and a few shuffled microbiome sets in Fig
644 S3.

645

646

Algorithm 1: Find a permutation of vectors \mathbf{x} and \mathbf{y} such that the correlation between permuted vectors is a predetermined value ρ

Take $\mathbf{x}, \mathbf{y}, \rho$, where \mathbf{x} and \mathbf{y} are arbitrary, uncorrelated vectors in R^n and $0 \leq \rho \leq 1$ is the desired correlation. The aim is to find a permutation of \mathbf{y} such that correlation $\text{cor}(\mathbf{x}, \mathbf{y}_{\text{sort}}) = \rho$, approximately. The algorithm can be equally applied when \mathbf{x} and / or \mathbf{y} are integer numbers, normality is not required either. Performance of the algorithm improves as n increases and when normality holds.

1. Sort the values of \mathbf{x} and \mathbf{y} in increasing or decreasing order. The correlation $\text{cor}(\mathbf{x}_{\text{sort}}, \mathbf{y}_{\text{sort}}) \cong 1$.
2. Generate a dummy variable $\mathbf{z} = \mathbf{y}_{\text{sort}} + \mathbf{e}$ where \mathbf{e} values are sampled from $\mathbf{e} \sim N(0, S_y^2 \frac{1-\rho^2}{\rho^2})$, S_y^2 is the sample variance of \mathbf{y} . The correlation $\text{cor}(\mathbf{x}_{\text{sort}}, \mathbf{z}) \sim \rho$.
3. Create an index variable \mathbf{iy} which indicates how \mathbf{y}_{sort} should be reordered according to \mathbf{z} order. This dummy index $\mathbf{iy} = \text{order}(\mathbf{y})[\text{order}(\mathbf{z})]$ contains the order of \mathbf{y} when values are back-sorted according to the order of \mathbf{z} .
4. Reorder $\mathbf{iy} = \mathbf{iy}[\text{rank}(\mathbf{x})]$ to match the index with positions \mathbf{y}_{sort} in the original vector \mathbf{x} . This is needed since \mathbf{x} remains unchanged and only \mathbf{y} is permuted.
5. The correlation $\text{cor}(\mathbf{x}, \mathbf{y}[\mathbf{iy}]) \cong \rho$.

Algorithm available at <https://github.com/miguelperezenciso/Simubiome>, see sortCorr function.

647

648 **Parameter fitting**

649 Little is known neither on the number of OTUs influencing a given phenotype nor on how
650 many of those are partly inherited. For that reason, we chose some extreme, yet ‘educated’
651 values for each of the five scenarios depicted in Fig 1. We considered $r^2 = h_g^2 + h_b^2 =$
652 0.25 and 0.50; $r^2 = 0.25$ is grossly the value reported by Difford et al. 2019 with $N = 750$,
653 whereas values closer to $r^2 = 0.50$ were reported by Wallace et al. in some farms. Overall,
654 augmenting r^2 values tries to mimic the effect of increasing sample size. We assumed h_g^2
655 $= h_b^2$ for Joint and Recursive scenarios, as also reported by Difford et al or Camarinha-
656 Silva et al. approximately. The number of QTNs was fixed to 100. This figure is
657 somewhat arbitrary, but the specific number of loci would not affect much the results.
658 Barton et al. [52] showed theoretically that most properties of the infinitesimal model

659 converge as fast as the inverse of the number of loci, or $\sim 1\%$ deviance with $N_{\text{QTN}} = 100$.
660 In general, genomic prediction is known to be relatively insensitive to the number of
661 QTNs [53]. As for individual genetic effects α , numerous empirical and theoretical works
662 show that they are not uniformly distributed and can be approximated by a gamma-like
663 distribution [54,55]. Here we sampled genetic effects $\alpha \sim \Gamma(\text{shape} = 0.2, \text{scale} = 5)$, as
664 suggested by Caballero et al. [56], and also used previously by us [57].

665

666 Much less is known on the number of causative OTUs (N_{OTU}), although we can presume
667 that N_{OTU} should be smaller than the number of QTNs. For instance, Duvallet et al. [36]
668 found in a large meta-analysis that the human diseases studied were affected on average
669 by 10 - 15 changes in abundances at the genus level. Here we considered $N_{\text{OTU}} = 25$ (0.6%
670 of all OTUs), although we also evaluated $N_{\text{OTU}} = 10, 100$ and 250. Similarly, for the
671 Recursive and Indirect scenarios, we took the extreme scenario where all causative OTUs
672 are genetically determined, i.e., $N_{\text{OTU}} = N_{\text{OTU}(\text{g})}$. The genetic effects β on abundances
673 (Eqn. 2) were sampled from the same distribution $\beta \sim \Gamma(\text{shape} = 0.2, \text{scale} = 5)$ as direct
674 genetic effects α . We are much more ignorant regarding the distribution of abundances'
675 effects ω on the phenotype (Eqn. 1). We took as proxy the regression coefficients of
676 methane emission on abundances published by Difford et al. [4], in their supplementary
677 information S4, which can be approximated by a $\Gamma(\text{shape}=1.4, \text{scale}=3.8)$. Fig S3
678 compares both gamma distributions and the fit to empirical data. This model predicts that
679 the variance of OTUs' effects is wider and of larger individual effect on average than that
680 of SNPs. Although this is speculative at this point, it is sensible to assume that only a few
681 taxa do have a sizeable influence on the phenotype, say methane emission.

682

683 **Analysis**

684 We used Bayes C algorithm [15] as implemented in BGLR [58] to assess prediction
685 performance and reliability of parameter estimates. We also tested Bayesian RKHS
686 regression, equivalent to GBLUP [40], but results were similar or worse and are not
687 presented. Three models were used to analyze the data:

688

$$689 \text{ Bayes Cgb: } \mathbf{y} = \mathbf{Z} \mathbf{a} + \mathbf{b} \mathbf{W} + \mathbf{e} \quad (3a)$$

690

$$691 \text{ Bayes Cg: } \mathbf{y} = \mathbf{Z} \mathbf{a} + \mathbf{e} \quad (3b)$$

692

693 Bayes Cb: $\mathbf{y} = \mathbf{W} \mathbf{b} + \mathbf{e}$ (3c)

694

695 where \mathbf{y} is the vector containing the simulated phenotypes, \mathbf{a} contains the marker effect
696 estimates, \mathbf{Z} contains the observed genotypes for the 33k markers, \mathbf{b} contains the OTU
697 abundance effects, \mathbf{W} is a matrix with all 4,018 abundances in the 750 individuals, and \mathbf{e}
698 is the residual. Prior to the analyses, phenotypes, abundances, and genotypic values were
699 standardized to mean zero and SD = 1. As priors π for SNPs or abundances probabilities
700 to enter into the model, we used $\pi \sim \text{Beta}(p_0 = 5, \pi_0 = 0.001)$, which has expectation π_0
701 and variance $\pi_0(1 - \pi_0) / (p_0 + 1)$. We also considered a much more liberal, flat prior for π
702 $\sim \text{Beta}(p_0 = 2, \pi_0 = 0.01)$, but we did not observe strong differences. A total of 50k
703 iterations were run per Bayes C chain, a plot of the residual variances along iterations
704 indicated convergence was attained with this number of iterations. To assess predictive
705 accuracy, 75 (10% of N) phenotypes were randomly removed and predicted with the fitted
706 model. Correlation between observed and predicted phenotypes was used as measure of
707 predictive accuracy.

708

709 The ‘heritability’ is not explicitly defined in a Bayes C framework, and here we used the
710 proposal by [58] (<https://github.com/gdlc/BGLR-R/blob/master/inst/md/heritability.md>).
711 In short, at each iteration i , the algorithm samples SNPs and OTUs effects:

712

713 $\mathbf{u}^{(i)} = \mathbf{Z} \hat{\mathbf{a}}^{(i)}$

714 $\mathbf{v}^{(i)} = \mathbf{W} \hat{\mathbf{b}}^{(i)}$

715

716 where $\mathbf{u}^{(i)}$ and $\mathbf{v}^{(i)}$ are genome and microbiome effects at i -the iteration for the set of
717 individuals, respectively, $\hat{\mathbf{a}}^{(i)}$ and $\hat{\mathbf{b}}^{(i)}$ are current SNP and OTU abundances solutions;
718 therefore, $\text{Var}(\mathbf{u}^{(i)}) / \text{Var}(\mathbf{y})$ and $\text{Var}(\mathbf{v}^{(i)}) / \text{Var}(\mathbf{y})$ are i -th iterate heritability and
719 microbiability estimates wherefrom posterior means can be estimated by averaging over
720 iterations. For Bayes Cgb, we also sampled the absolute covariance between \mathbf{u} and \mathbf{v} , i.e.,
721 $|\text{Cov}(\mathbf{u}^{(i)}, \mathbf{v}^{(i)})| / \text{Var}(\mathbf{y})$.

722

723 To assess how likely is to identify causative OTUs, we retained the probability of a given
724 OTU entering into the model, averaged over Gibbs sampling iterations. We run a GWAS

725 of abundances (\mathbf{x}_k , $k=1, N_{\text{OTU}}$) on SNP genotypes (\mathbf{z}_j , $j=1, N_{\text{SNP}}$) using R function $\text{lm}(\mathbf{x}_k$
726 $\sim \mathbf{z}_j)$ and we computed the P-value of both causative QTNs, i.e., affecting abundances,
727 and neutral SNPs. This was done in the Recursive scenario only. In this scenario, we also
728 computed the heritabilities of all abundance levels using GBLUP via a RKHS strategy
729 (<https://github.com/gdlc/BGLR-R/blob/master/inst/md/GBLUP.md#RKHS>) using
730 BGLR. Weakly informative priors for variances were used to mimic a REML-like
731 estimator.

732

733 **Author contributions**

734 MPE, GDLC and LMZ conceived research. MPE and LMZ performed research. All
735 authors discussed research. MPE wrote the manuscript with help from the rest of authors.

736

737 **Acknowledgments**

738 This work was developed while MPE and LMZ spent a research stay at Michigan State
739 University kindly funded by GDLC. LMZ is supported by a Ph.D. grant from the Ministry
740 of Economy and Science (MINECO, Spain), MPE is funded by MINECO grants
741 AGL2016-78709-R and PID2019-108829RB-I00, from the EU through BFU2016-
742 77236-P (MINECO/AEI/FEDER, EU) and “Centro de Excelencia Severo Ochoa 2016-
743 2019” award SEV-2015-0533. YRC was funded by Marie Skłodowska-Curie grant (P-
744 Sphere) agreement No 6655919 (EU).

745

746 **References**

- 747 1. Ruff WE, Greiling TM, Kriegel MA. Host–microbiota interactions in immune-
748 mediated diseases [Internet]. Nature Reviews Microbiology. Nature Research;
749 2020. pp. 521–538. doi:10.1038/s41579-020-0367-2
- 750 2. Zhang Q, Difford G, Sahana G, Løvendahl P, Lassen J, Lund MS, et al. Bayesian
751 modelling reveals host genetics associated with rumen microbiota jointly
752 influence methane emission in dairy cows. ISME J. 2020;14: 2019–2033.
753 doi:10.1038/s41396-020-0663-x
- 754 3. Maltecca C, Bergamaschi M, Tiezzi F. The interaction between microbiome and
755 pig efficiency: A review. J Anim Breed Genet. Blackwell Publishing Ltd;
756 2020;137: 4–13. doi:10.1111/jbg.12443
- 757 4. Difford GF, Plichta DR, Løvendahl P, Lassen J, Noel SJ, Højberg O, et al. Host
758 genetics and the rumen microbiome jointly associate with methane emissions in

- 759 dairy cows. *PLoS Genet. Public Library of Science*; 2018;14: e1007580.
760 doi:10.1371/journal.pgen.1007580
- 761 5. Kundu P, Blacher E, Elinav E, Pettersson S. Our Gut Microbiome: The Evolving
762 Inner Self [Internet]. *Cell. Cell Press*; 2017. pp. 1481–1493.
763 doi:10.1016/j.cell.2017.11.024
- 764 6. Difford GF, Lassen J, Løvendahl P. Genes and microbes, the next step in dairy
765 cattle breeding. *EAAP—67th Annual Meeting*. 2016.
- 766 7. Falconer D, Mackay T. *Introduction to Quantitative Genetics*. essex: Longman
767 Publishing Group; 1996.
- 768 8. Zilber-Rosenberg I, Rosenberg E. Role of microorganisms in the evolution of
769 animals and plants: the hologenome theory of evolution. *FEMS Microbiol Rev.*
770 *Oxford Academic*; 2008;32: 723–735. doi:10.1111/j.1574-6976.2008.00123.x
- 771 9. Rothschild D, Weissbrod O, Barkan E, Kurilshikov A, Korem T, Zeevi D, et al.
772 Environment dominates over host genetics in shaping human gut microbiota.
773 *Nature. Nature Publishing Group*; 2018;555: 210–215. doi:10.1038/nature25973
- 774 10. Lloyd-Price J, Arze C, Ananthakrishnan AN, Schirmer M, Avila-Pacheco J, Poon
775 TW, et al. Multi-omics of the gut microbial ecosystem in inflammatory bowel
776 diseases. *Nature. Nature Publishing Group*; 2019;569: 655–662.
777 doi:10.1038/s41586-019-1237-9
- 778 11. Wallace JR, Sasson G, Garnsworthy PC, Tapio I, Gregson E, Bani P, et al. A
779 heritable subset of the core rumen microbiome dictates dairy cow productivity
780 and emissions. *Sci Adv. American Association for the Advancement of Science*;
781 2019;5: eaav8391. doi:10.1126/sciadv.aav8391
- 782 12. Camarinha-Silva A, Maushammer M, Wellmann R, Vital M, Preuss S,
783 Bennewitz J. Host genome influence on gut microbial composition and microbial
784 prediction of complex traits in pigs. *Genetics. Genetics Society of America*;
785 2017;206: 1637–1644. doi:10.1534/genetics.117.200782
- 786 13. Khanal P, Maltecca C, Schwab C, Fix J, Tiezzi F. Microbiability of meat quality
787 and carcass composition traits in swine. *bioRxiv. Cold Spring Harbor*
788 *Laboratory*; 2019; 833731. doi:10.1101/833731
- 789 14. Pereyra MA, Creus CM. *Modifying the Rhizosphere of Agricultural Crops to*
790 *Improve Yield and Sustainability: Azospirillum as a Model Rhizotroph.*
791 *Rhizotrophs: Plant Growth Promotion to Bioremediation. Springer Singapore*;
792 2017. pp. 15–37. doi:10.1007/978-981-10-4862-3_2

- 793 15. Meuwissen THE, Hayes BJ, Goddard ME. Prediction of total genetic value using
794 genome-wide dense marker maps. *Genetics*. 2001;157: 1819–1829. Available:
795 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1461589&tool=pmce](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1461589&tool=pmcentrez&rendertype=abstract)
796 [ntrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1461589&tool=pmcentrez&rendertype=abstract)
- 797 16. Crossa J, Pérez-Rodríguez P, Cuevas J, Montesinos-López O, Jarquín D, de los
798 Campos G, et al. Genomic Selection in Plant Breeding: Methods, Models, and
799 Perspectives [Internet]. *Trends in Plant Science*. Elsevier Current Trends; 2017.
800 pp. 961–975. doi:10.1016/j.tplants.2017.08.011
- 801 17. Meuwissen T, Hayes B, Goddard M. Accelerating Improvement of Livestock
802 with Genomic Selection. *Annu Rev Anim Biosci*. 2013;1: 221–237.
803 doi:10.1146/annurev-animal-031412-103705
- 804 18. De Los Campos G, Gianola D, Allison DB. Predicting genetic predisposition in
805 humans: The promise of whole-genome markers. *Nat Rev Genet*. Nature
806 Publishing Group; 2010;11: 880–886. doi:10.1038/nrg2898
- 807 19. Tierney BT, He Y, Church GM, Segal E, Kostic AD, Patel CJ. The predictive
808 power of the microbiome exceeds that of genome-wide association studies in the
809 discrimination of complex human disease. *BioRxiv*. Cold Spring Harbor
810 Laboratory; 2020; 2019.12.31.891978. doi:10.1101/2019.12.31.891978
- 811 20. Wang J, Chen L, Zhao N, Xu X, Xu Y, Zhu B. Of genes and microbes: solving
812 the intricacies in host genomes. *Protein Cell*. Higher Education Press; 2018;9:
813 446–461. doi:10.1007/s13238-018-0532-9
- 814 21. Benson AK, Kelly SA, Legge R, Ma F, Low SJ, Kim J, et al. Individuality in gut
815 microbiota composition is a complex polygenic trait shaped by multiple
816 environmental and host genetic factors. *Proc Natl Acad Sci U S A*. National
817 Academy of Sciences; 2010;107: 18933–18938. doi:10.1073/pnas.1007028107
- 818 22. Goodrich JK, Davenport ER, Clark AG, Ley RE. The Relationship Between the
819 Human Genome and Microbiome Comes into View. *Annu Rev Genet*. Annual
820 Reviews; 2017;51: 413–433. doi:10.1146/annurev-genet-110711-155532
- 821 23. Blekhman R, Goodrich JK, Huang K, Sun Q, Bukowski R, Bell JT, et al. Host
822 genetic variation impacts microbiome composition across human body sites.
823 *Genome Biol*. BioMed Central Ltd.; 2015;16: 191. doi:10.1186/s13059-015-
824 0759-1
- 825 24. Crespo-Piazuelo D, Migura-Garcia L, Estellé J, Criado-Mesas L, Revilla M,
826 Castelló A, et al. Association between the pig genome and its gut microbiota

- 827 composition. *Sci Rep. Nature Publishing Group*; 2019;9: 8791.
828 doi:10.1038/s41598-019-45066-6
- 829 25. Ramayo-Caldas Y, Zingaretti L, Popova M, Estellé J, Bernard A, Pons N, et al.
830 Identification of rumen microbial biomarkers linked to methane emission in
831 Holstein dairy cows. *J Anim Breed Genet.* 2019; 49–59. doi:10.1111/jbg.12427
- 832 26. Ramayo-Caldas Y, Prenafeta-Boldú F, Zingaretti LM, Gonzalez-Rodriguez O,
833 Dalmau A, Quintanilla R, et al. Gut eukaryotic communities in pigs: diversity,
834 composition and host genetics contribution. *Anim Microbiome. BioMed Central*;
835 2020;2: 18. doi:10.1186/s42523-020-00038-4
- 836 27. Weissbrod O, Rothschild D, Barkan E, Segal E. Host genetics and microbiome
837 associations through the lens of genome wide association studies. *Current*
838 *Opinion in Microbiology.* Elsevier Ltd; 2018. pp. 9–19.
839 doi:10.1016/j.mib.2018.05.003
- 840 28. Bäckhed F, Roswall J, Peng Y, Feng Q, Jia H, Kovatcheva-Datchary P, et al.
841 Dynamics and stabilization of the human gut microbiome during the first year of
842 life. *Cell Host Microbe. Cell Press*; 2015;17: 690–703.
843 doi:10.1016/j.chom.2015.04.004
- 844 29. Jakobsson HE, Abrahamsson TR, Jenmalm MC, Harris K, Quince C, Jernberg C,
845 et al. Decreased gut microbiota diversity, delayed Bacteroidetes colonisation and
846 reduced Th1 responses in infants delivered by Caesarean section. *Gut. BMJ*
847 *Publishing Group*; 2014;63: 559–566. doi:10.1136/gutjnl-2012-303249
- 848 30. Furman O, Shenhav L, Sasson G, Kokou F, Honig H, Jacoby S, et al.
849 Stochasticity constrained by deterministic effects of diet and age drive rumen
850 microbiome assembly dynamics. *Nat Commun. Nature Research*; 2020;11: 1–13.
851 doi:10.1038/s41467-020-15652-8
- 852 31. Gamazon ER, Segrè A V., Van De Bunt M, Wen X, Xi HS, Hormozdiari F, et al.
853 Using an atlas of gene regulation across 44 human tissues to inform complex
854 disease- and trait-associated variation. *Nat Genet. Nature Publishing Group*;
855 2018;50: 956–967. doi:10.1038/s41588-018-0154-4
- 856 32. Hormozdiari F, van de Bunt M, Segrè A V., Li X, Joo JWJ, Bilow M, et al.
857 Colocalization of GWAS and eQTL Signals Detects Target Genes. *Am J Hum*
858 *Genet. Cell Press*; 2016;99: 1245–1260. doi:10.1016/j.ajhg.2016.10.003
- 859 33. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome
860 Datasets Are Compositional: And This Is Not Optional. *Front Microbiol.*

- 861 Frontiers; 2017;8: 2224. doi:10.3389/fmicb.2017.02224
- 862 34. Weiss S, Xu ZZ, Peddada S, Amir A, Bittinger K, Gonzalez A, et al.
863 Normalization and microbial differential abundance strategies depend upon data
864 characteristics. *Microbiome*. BioMed Central; 2017;5: 27. doi:10.1186/s40168-
865 017-0237-y
- 866 35. Maltecca C, Lu D, Schillebeeckx C, McNulty NP, Schwab C, Shull C, et al.
867 Predicting Growth and Carcass Traits in Swine Using Microbiome Data and
868 Machine Learning Algorithms. *Sci Rep*. Nature Publishing Group; 2019;9: 1–15.
869 doi:10.1038/s41598-019-43031-x
- 870 36. Power RA, Parkhill J, De Oliviera T. Microbial Genome-Wide Association
871 Studies: Lessons from Human GWAS. doi:10.1101/093211
- 872 37. Roehe R, Dewhurst RJ, Duthie CA, Rooke JA, McKain N, Ross DW, et al.
873 Bovine Host Genetic Variation Influences Rumen Microbial Methane Production
874 with Best Selection Criterion for Low Methane Emitting and Efficiently Feed
875 Converting Hosts Based on Metagenomic Gene Abundance. Leeb T, editor.
876 *PLoS Genet*. 2016;12: e1005846. doi:10.1371/journal.pgen.1005846
- 877 38. Huws SA, Creevey CJ, Oyama LB, Mizrahi I, Denman SE, Popova M, et al.
878 Addressing global ruminant agricultural challenges through understanding the
879 rumen microbiome: Past, present, and future. *Frontiers in Microbiology*.
880 Frontiers Media S.A.; 2018. p. 2161. doi:10.3389/fmicb.2018.02161
- 881 39. Wallace RJ, Rooke JA, McKain N, Duthie CA, Hyslop JJ, Ross DW, et al. The
882 rumen microbial metagenome associated with high methane production in cattle.
883 *BMC Genomics*. BioMed Central Ltd.; 2015;16: 839. doi:10.1186/s12864-015-
884 2032-0
- 885 40. VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci*.
886 Elsevier; 2008;91: 4414–4423. doi:10.3168/jds.2007-0980
- 887 41. Ross EM, Moate PJ, Maret LC, Cocks BG, Hayes BJ. Metagenomic predictions:
888 from microbiome to complex health and environmental phenotypes in humans
889 and cattle. White BA, editor. *PLoS One*. Public Library of Science; 2013;8:
890 e73056. doi:10.1371/journal.pone.0073056
- 891 42. Saborío-Montero A, Gutiérrez-Rivas M, García-Rodríguez A, Atxaerandio R,
892 Goiri I, López de Maturana E, et al. Structural equation models to disentangle the
893 biological relationship between microbiota and complex traits: Methane
894 production in dairy cattle as a case of study. *J Anim Breed Genet*. Blackwell

- 895 Publishing Ltd; 2020;137: 36–48. doi:10.1111/jbg.12444
- 896 43. Muñoz PR, Resende MFR, Gezan SA, Resende MDV, de los Campos G, Kirst
897 M, et al. Unraveling additive from nonadditive effects using genomic relationship
898 matrices. *Genetics*. Genetics Society of America; 2014;198: 1759–1768.
899 doi:10.1534/genetics.114.171322
- 900 44. Cole NA. Effects of animal-to-animal exchange of ruminal contents on the feed
901 intake and ruminal characteristics of fed and fasted lambs. *J Anim Sci*. Oxford
902 Academic; 1991;69: 1795. doi:10.2527/1991.6941795x
- 903 45. Weimer PJ. Redundancy, resilience, and host specificity of the ruminal
904 microbiota: Implications for engineering improved ruminal fermentations.
905 *Frontiers in Microbiology*. Frontiers Media S.A.; 2015. p. 296.
906 doi:10.3389/fmicb.2015.00296
- 907 46. Dill-Mcfarland KA, Breaker JD, Suen G. Microbial succession in the
908 gastrointestinal tract of dairy cows from 2 weeks to first lactation. *Sci Rep*.
909 Nature Publishing Group; 2017;7: 1–12. doi:10.1038/srep40864
- 910 47. Crow J, Kimura M. *An Introduction to Population Genetics Theory*. The
911 Blackburn Press; 1970.
- 912 48. Peng B, Kimmel M. simuPOP: A forward-time population genetics simulation
913 environment. *Bioinformatics*. 2005;21: 3686–3687.
914 doi:10.1093/bioinformatics/bti584
- 915 49. Messer PW. SLiM: simulating evolution with selection and linkage. *Genetics*.
916 2013;194: 1037–1039. doi:10.1534/genetics.113.152181
- 917 50. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide
918 complex trait analysis. *Am J Hum Genet*. The American Society of Human
919 Genetics; 2011;88: 76–82. doi:10.1016/j.ajhg.2010.11.011
- 920 51. Pérez-Enciso M, Ramírez-Ayala LC, Zingaretti LM. SeqBreed: A python tool to
921 evaluate genomic prediction in complex scenarios. *Genet Sel Evol*. BioMed
922 Central Ltd.; 2020;52. doi:10.1186/s12711-020-0530-2
- 923 52. Barton NH, Etheridge AM, Véber A. The infinitesimal model: Definition,
924 derivation, and implications. *Theor Popul Biol*. Elsevier Inc.; 2017;118: 50–73.
925 doi:10.1016/j.tpb.2017.06.001
- 926 53. Daetwyler HD, Pong-Wong R, Villanueva B, Woolliams JA. The impact of
927 genetic architecture on genome-wide evaluation methods. *Genetics*. 2010;185:
928 1021–1031. doi:10.1534/genetics.110.116855

- 929 54. Eyre-Walker A, Keightley PD. The distribution of fitness effects of new
930 mutations. *Nat Rev Genet.* 2007;8: 610–618. doi:10.1038/nrg2146
- 931 55. Barton NH, Keightley PD. Understanding Quantitative Genetic Variation. *Nat*
932 *Rev Genet.* 2002;3: 11–21. doi:10.1038/nrg700
- 933 56. Caballero A, Tenesa A, Keightley PD. The nature of genetic variation for
934 complex traits revealed by GWAS and Regional Heritability Mapping analyses.
935 *Genetics.* 2015;201: 1601–1613.
- 936 57. Pérez-Enciso M, Forneris N, de Los Campos G, Legarra A. Evaluating Sequence-
937 Based Genomic Prediction with an Efficient New Simulator. *Genetics.* *Genetics*;
938 2017;205: 939–953. doi:10.1534/genetics.116.194878
- 939 58. Pérez P, de Los Campos G. Genome-Wide Regression & Prediction with the
940 BGLR Statistical Package. *Genetics.* 2014;198: 483–95.
941 doi:10.1534/genetics.114.164442
942