

1 **SNP-based quantitative deconvolution of biological mixtures: application to the detection**
2 **of cows with subclinical mastitis by whole genome sequencing of tank milk.**

3

4 *Wouter Coppieters¹, Latifa Karim¹, Michel Georges².*

5

6 ¹Genomics Platform, GIGA Institute, University of Liège. ²Unit of Animal Genomics, GIGA
7 Institute & Faculty of Veterinary Medicine, University of Liège.

8

9 Correspondence: michel.georges@uliege.be

10

11 **Biological products of importance in food (f.i. milk) and medical (f.i. donor blood derived**
12 **products) sciences often correspond to mixtures of samples contributed by multiple**
13 **individuals. Identifying which individuals contributed to the mixture and in what**
14 **proportions may be of interest in several circumstances. We herein present a method that**
15 **allows to do this by shallow whole genome sequencing of the DNA in mixed samples from**
16 **hundreds of donors. We demonstrate the efficacy of the approach for the detection of cows**
17 **with subclinical mastitis by analysis of farms' tank mixtures containing milk from as many**
18 **as 500 cows.**

19

20 **Introduction**

21 Mastitis is the most important health issue in dairy cattle costing European farmers > 1 billion
22 € per year in treatment and milk loss¹. Mastitis is routinely managed by periodically counting
23 immune cells in milk samples to preemptively identify cows developing subclinical udder
24 inflammation. As profit margins decrease, farmers tend to forgo milk testing thereby
25 compromising health management. Cost-effective alternatives for rapid detection of cows
26 with subclinical mastitis are needed². We previously proposed that somatic cell counts (SCC)
27 in the milk of individual cows could be estimated if B allele frequencies were measured for
28 sufficient numbers of SNPs in the tank milk, provided that all cows contributing milk to the
29 tank be genotyped for the corresponding variants³. Thus, the proposed method would allow
30 to identify a minority of cows with subclinical mastitis by analyzing a single sample containing
31 a mixture of milk from all the cows on the farm, hence dramatically reducing costs. As
32 genomic selection (GS) is becoming routine (including for dams)⁴, herds that are fully

33 genotyped with low density SNP arrays (~15K) are becoming standard, and the proposed
34 method feasible. We herein demonstrate that by combining low density SNP genotyping or
35 shallow sequencing of the cows and tank milk's DNA with in silico genotype imputation,
36 individual SCC can be accurately determined and cows with subclinical mastitis effectively
37 identified even in the largest farms (≥ 500). The proposed method has the potential to
38 dramatically improve the monitoring of udder health in dairy farms, and to allow the tracing
39 of the origin of bulk animal food products other than milk.

40

41 **Results**

42 ***Principle of the proposed method.*** Milk of healthy cows typically contains $\leq 100,000$ somatic
43 cells per ml. Upon infection leucocytes migrate in the udder and SCC increase rapidly: SCC \geq
44 200,000 / ml are indicative of subclinical mastitis, while SCC into the millions are common for
45 cows with overt mastitis⁵. Assume that cows and tank (i.e. the reservoir in which the milk of
46 the cows is collected) milk are genotyped for a collection of SNPs. If all cows contribute
47 identical amounts of DNA to the milk, the expected "B" allele frequency in the tank milk
48 corresponds to the frequency of the "B" allele in the farm's cow population. The actual DNA
49 amount contributed by each cow depends on the volume of milk produced and its SCC.
50 Unequal DNA contributions will cause slight departures from the expected B allele
51 frequencies in the tank milk. Integrating these shifts over a large number of SNPs in
52 conjunction with the known genotypes of individual cows (using f.i. a linear model) allows for
53 the estimation of the relative DNA contribution of each cow. Accounting for individual milk
54 volumes and for the SCC in the tank milk allows for the estimation of SCC for individual cows
55 (Fig. 1 and Methods).

56 ***Evaluating the proposed method by simulation.*** We first evaluated the proposed method by
57 simulation (cfr. Methods). Genotyping the cows and the tank milk using 10K SNP arrays (i.e.
58 low-density (LD) arrays as generally used for GS) allowed for the accurate estimation of
59 individual SCC for farms with up to 100 cows ($r \geq 0.9$, where r is the correlation between
60 real and estimated SCC) (scheme A). However, farms with > 100 cows are increasingly
61 common. Medium- (MD, f.i. 50K) and high-density (HD, f.i. 700K) SNP arrays would be
62 needed for the approach to be effective in farms with ≥ 250 or ≥ 500 cows, respectively. Yet
63 – being too expensive - this is presently not a viable proposition (Fig. 2A). We therefore
64 envisaged a second scheme (B) in which the cows would still be genotyped with LD SNP arrays

65 (as done in practice) yet imputed⁶ to whole genome (8 million SNPs in the simulations) using
66 a sequenced reference population⁷, while the DNA of the tank milk would be genotyped by
67 shallow whole-genome sequencing (SWGS). We found that under this scenario sequencing
68 the tank milk at a depth of 0.25 was sufficient for farms with 100 cows, 0.5 for farms with 250
69 cows, and 2 for farms with 500 cows (Fig. 2B). Accuracies were not significantly affected by
70 the density of the SNP arrays, i.e. the method performed as well with LD as with MD arrays
71 (data not shown). Anticipating further advances in sequencing technology, we also envisaged
72 a scheme (C) in which both cows and tank milk would be genotyped by SWGS. We found that
73 a 1-fold sequencing depth of the tank milk would be sufficient when combined with a 0.25-
74 fold depth for 100 cows, while a 5-fold sequencing depth of the tank milk would be needed
75 in combination with 0.25-fold depth for 250 cows and 1-fold depth for 500 cows (Fig. 2C). In
76 scheme C, allelic dosage in the cows is directly measured from the number of alternative and
77 reference alleles in the sequence reads. We further explored the effectiveness of augmenting
78 the cow genotype information from SWGS by imputation (scheme D). This proved to be
79 effective, reducing the required sequence depth to 0.25-fold for tank milk and 0.25-fold for
80 100 cows, to 1-fold for tank milk and 0.25-fold for 250 cows, and to 5-fold for tank milk and
81 0.25-fold for 500 cows (Fig. 2D).

82 ***Real-world application of the proposed method.*** To test the feasibility of our method in the
83 real world, we first collected cow (blood) and tank (milk) samples from a farm milking 133
84 Holstein-Friesian cows. When only using genotypes from the Illumina LD arrays (17K SNPs)
85 for both cows and tank milk (scheme A), correlations between predicted and measured SCC
86 were 0.91 (or 0.79 when ignoring one cow with SCC > 3 million). We then imputed the cows
87 to whole genome (13M SNPs) using a reference population of ~750 whole genome
88 sequenced Holstein-Friesian animals, and sequenced the tank milk at ~3.5-fold depth. The
89 corresponding correlations (scheme B) were 0.97 (0.95) when using all sequence information,
90 or 0.96 (0.92) when down-sampling sequence information as low as 0.1-fold depth (Fig. 3A).
91 We next performed a similar experiment on a farm milking 520 Holstein-Friesian cows. The
92 correlation between predicted and measured SCC was 0.78 (or 0.42 when ignoring 23 cows
93 with SCC > 3 million) when only using information from the LD array for both cows and tank
94 milk (scheme A). When imputing the cows to whole genome (13M SNPs) and sequencing
95 the milk at ~3.5-fold depth (scheme B), the correlation increased to 0.89 (0.83). Down-

96 sampling the sequence information to 0.1-fold depth reduced the correlation to 0.79 (0.57)
97 (Fig. 3B).

98 As shown in both farms, correlation estimates are affected by SCC spread: small numbers of
99 cows with very high SCC tend to inflate r . We therefore computed accuracies, computed as
100 the proportion of correctly classified cows for different SCC thresholds, which is how farmers
101 would likely use the information. It can be seen that for a threshold value of for example
102 500,000 SCC, accuracies > 0.85 were obtained when sequencing (scheme B) the tank milk at
103 respectively 0.1x (133 cows) and 3.5x depth (520 cows). Thus - as predicted by the simulations
104 - scheme A provided adequate precision for the farm with 133 cows, but not for the farm
105 with 520 cows. However, in this large farm, combining SWGS of the tank milk with whole
106 genome imputation of the cows (i.e. scheme B) was indeed effective (Fig. 3).

107 As costs per bp continue to decline, sequencing is likely to replace array-based genotyping in
108 the future. To test the feasibility of schemes C and D (i.e. genotype the cows by SWGS rather
109 than with SNP arrays, without (C) and with (D) imputation), we collected samples from a farm
110 with 120 Holstein-Friesian cows. All cows were genotyped with the Illumina LD array (17K) as
111 well as sequenced at average 1.08 -fold depth (range: 0.26-1.73). The milk was sequenced at
112 ~ 3.5 -fold depth. The correlation between predicted and measured SCC was 0.97 (or 0.96
113 when ignoring one cow with SCC > 3 million) under scheme A. Under scheme C, correlations
114 were 0.82 (0.83) when sequencing the milk at 3.5x and 0.75 (0.76) when down-sampling the
115 milk to 0.1x. We then imputed the sequenced cows to HD (770K SNPs) using a population of
116 800 reference animals genotyped with the HD array (scheme D). The correlation increased
117 to 0.93 (0.94) when sequencing the milk at 3.5x and to 0.83 (0.77) when down-sampling the
118 milk to 0.1x (Fig. 3C). Accuracies at SCC threshold of 500,000 were 0.96 (scheme A), 0.95
119 (3.5x) and 0.80 (0.1x) (scheme B), 0.82 (3.5x) and 0.81 (0.1x) (scheme C), and 0.95 (3.5x) and
120 0.88 (0.1x) (scheme D) (Fig. 3C). In summary, (i) combining cow genotyping using SNP arrays
121 with genome-wide imputation with SWGS of tank milk allows for cost-effective identification
122 of cows with subclinical mastitis even in farms with as many as 500 cows per milk tank, and
123 (ii) as sequencing costs continue to decline, arrays-based targeted SNP genotyping of the
124 cows could be replaced by genotyping by SWGS and yield comparable results.

125 **Monitoring SCC dynamics with the proposed method.** Farmers typically measure individual
126 SCC once a month or less. Yet, SCC may rapidly change. The SCC measured on the milk testing
127 date may not be a reliable indicator of the cow's udder health during the intervening period.

128 To examine the SCC dynamics over time, we collected 20 tank milk samples over a 100-day
129 period (day -84 to +17 from day of milk testing) for the farm with 120 cows. Milk samples
130 were genotyped using the Illumina LD array, and individual SCC estimated using scheme A.
131 Fig. 4A shows the SCC predicted every 5 days on average for the 120 cows, sorted by SCC
132 measured on day 0 (=milk testing day). Of note, the correlation between the SCC measured
133 on day 0 and the average of the SCC estimates for the 21 collection dates was low ($r =$
134 0.52)(Fig. 4B) and decreased rapidly with the number of days from milk testing day (Fig. 4C).
135

136 **Discussion**

137 We herein demonstrate that by combining array-based SNP genotyping and whole-genome
138 imputation for the cows with SWGS of the tank milk, it is possible to accurately estimate SCC
139 for individual cows and hence effectively identify animals with subclinical mastitis even for
140 tanks collecting milk for >500 cows, and this by performing a single analysis for the entire
141 herd. Reagent costs to sequence a mammalian genome at 1-fold depth are now <20€ thus
142 making this a cost-effective proposition. As a matter of fact, the method is being deployed
143 in the field in several countries.

144 Implementing the method requires all cows on the farm to be genotyped. This will
145 increasingly correspond to reality as genotyping costs continue to decrease and genomic
146 selection is more and more used for the selection of cows. In 2016 more than 1.2 million
147 dairy cows had been reportedly genotyped in the US alone⁸ and present worldwide numbers
148 are likely ≥ 3 million. In addition, a reference population of a few hundred animals of the
149 breed of interest that are either HD genotyped (700K) or better whole-genome sequenced
150 are required for accurate imputation. Such reference populations are already available for
151 the most important dairy cattle breeds^{7,9}, and could be easily generated for the remaining
152 ones.

153 We show that SCC are dynamic and rapidly change over time. SCC measured on day 0 are
154 poor indicators of SCC in previous and future weeks: cows with high SCC on the day of milk
155 testing may have low SCC a few days later (or earlier) and vice versa. The proposed method
156 would allow tighter monitoring of SCC hence improving udder health management. More
157 frequent monitoring of SCC for large number of cows may reveal interindividual differences
158 with regards to SCC dynamics that may be correlated with mastitis resistance, heritable and
159 hence amenable to selection including by GS.

160 Sequencing of the DNA in the tank milk allows simultaneous characterization of the tank's
161 microbiome. As a matter of fact, ~1% of reads in this study mapped to bacterial genomes
162 (data not shown). This information may be very useful both from a farm health management
163 point of view as well as from a downstream dairy processing point of view. Whole genome
164 sequence data of bulk milk also informs about the herd frequency of functional variants such
165 casein variants affecting consumer health or processing properties¹⁰, or variants causing
166 inherited defects or embryonic lethality in cows⁴. In many countries, it is not allowed to add
167 milk from cows being treated with antibiotics to the tank. As suggested before, the proposed
168 approach can be adapted to verify whether a specific cow did contribute milk to the tank or
169 not (f.i. by testing the significance of the corresponding cow effect in the linear model)³. The
170 described method may have applications in tracing the origins of bulk animal food products
171 other than milk, as well as in monitoring the composition of mixed-donor blood-derived
172 transfusion products.

173

174 **Acknowledgements**

175 This work was funded by the Unit of Animal Genomics and by the ERC DAMONA grant to
176 Michel Georges. We are grateful to Jean-Bernard Davière, Pierre Lenormand, Bonny Van
177 Ranst, Kristien Neyens and Miel Hostens for providing the samples and information needed
178 to conduct the experiments.

179

180 **References**

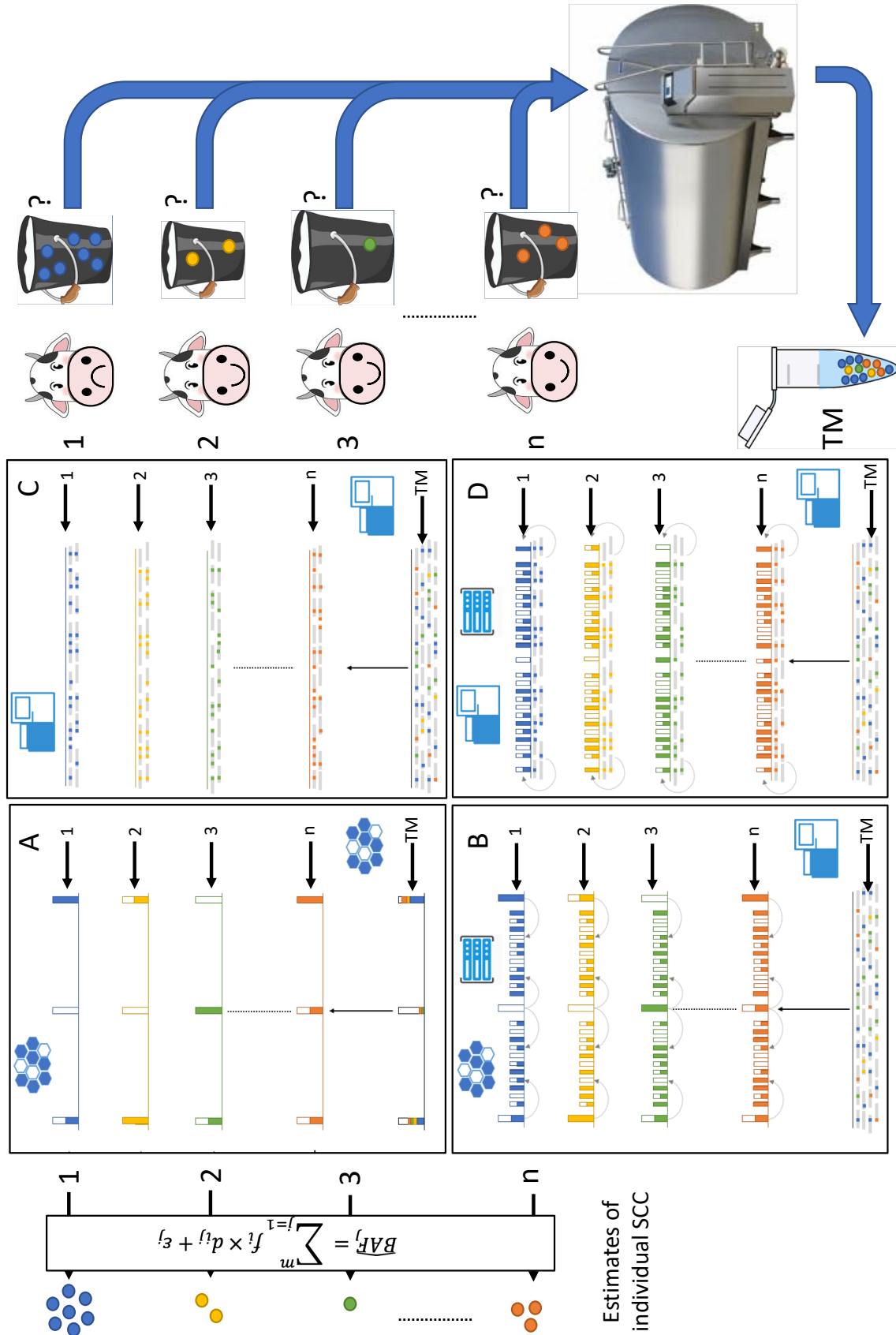
- 181 1. Hogeveen H, Huijps K, Lam TJGM. Economic aspects of mastitis: New developments. *N. Z.*
182 *Vet. J.* 59:16–23 (2011).
- 183 2. Viguié C, Arora S, Gilmartin N, Welbeck K, O’Kennedy R. Mastitis detection: current
184 trends and future perspectives. *Trends Biotechnol* 27: 486-493 (2009).
- 185 3. Blard G, Zhang Z, Coppieters W, Georges M. Identifying cows with subclinical mastitis by
186 bulk SNP genotyping of tank milk. *J Dairy Sci* 95:4109-4113 (2012).
- 187 4. Georges M, Charlier C, Hayes B. Genomics selection of livestock and beyond. *Nat Rev*
188 *Genet* 20:135-156 (2019).
- 189 5. Schukken YH, Wilson DJ, Welcome F, Garrison-Tikofsky L, Gonzales RN. Monitoring udder
190 health and milk quality using somatic cell counts. *Vet Res* 34: 579-596 (2003).

- 191 6. Marchini J, Howie B. Genotype imputation for genome-wide association studies. *Nat Rev*
192 *Genet* 11:499-511 (2010).
- 193 7. Daetwyler HD *et al.* Whole-genome sequencing of 234 bulls facilitates mapping of
194 monogenic and complex traits in cattle. *Nat Genet* 46: 858-865 (2014).
- 195 8. Wiggans GR, Cole JB, Hubbard SM, Sonstegard TS. Genomic selection in dairy cattle: the
196 USDA experience. *Annu Rev Anim Biosci* 5:309-327 (2017).
- 197 9. Charlier C. *et al.* Reverse genetic screen for embryonic lethal mutations comprising
198 fertility in cattle. *Genome Res* 26: 1-9 (2016).
- 199 10. Brooke-Taylor S, Dwyer K, Woodford K, Kost N. Systematic review of the gastrointestinal
200 effects of A1 compared with A2 β -casein. *Adv Nutr* 8:739-748 (2017).

201 **Figure 1:** Estimating Somatic Cell Counts (SCC) in the milk of individual cows by analyzing a
202 sample of milk from the farm's tank. Cows 1 to n contribute different amounts of milk
203 (buckets of various sizes in the figure) to the farm's tank. The milk contains somatic cells
204 (shown as small spheres in the milk colored by cow) whose numbers reflect the health status
205 of the cow's udder. Cow 1 has higher SCC, an indicator of subclinical mastitis. SCC are
206 unknown upon milking (indicated by the "?"). Cows are individually SNP genotyped once. In
207 scheme A this is done using SNP arrays (illustrated by the mesh) yielding genotype
208 information for the limited number of interrogated SNPs (high bars) that can be summarized
209 by the B-allele frequency as shown (white: 0, halve colored: 0.5, full colored: 1). SNP
210 genotypes of individual cows are coded in the same colors as the SCC. In scheme B, the
211 genotypes of the interrogated SNPs are augmented by imputation (illustrated by the
212 computer rack), yielding dosage information (B-allele frequency) for many more SNPs (small
213 bars). In scheme C, cows are genotyped individually by shallow whole genome sequencing
214 (SWGS) (illustrated by the sequencer). Sequence reads (gray lines) are aligned to the
215 reference genome and alternate alleles at SNP positions highlighted as color-coded ticks. The
216 B-allele frequency at specific SNP positions is measured as the ratio of the number of reads
217 with alternate (B) vs the total number of reads. In scheme D, the genotype information from
218 SWGS is augmented by imputation improving the accuracy of the B-allele frequency estimates
219 for millions of SNPs (small bars). A small sample of milk (T(ank) M(ilk)) is periodically (f.i.
220 monthly or weekly) collected from the farm's tank. DNA is extracted from TM and genotyped
221 using SNP arrays (scheme A) or SWGS (schemes B, C and D). B-allele frequency for SNP j in
222 the milk ($\overline{B\bar{A}F}_j$) is estimated from the ratio of fluorescence intensities when using SNP arrays,
223 or from the proportion of reads with B allele in SWGS. The SCC of individual cows are
224 estimated from a set of linear equations modelling $\overline{B\bar{A}F}_j$ as the sum of B allele dosage (d_{ij})
225 multiplied by the proportion of the DNA in the tank contributed by cow i (f_i). The estimated
226 proportions of DNA contributed by each cow correspond to the values of f_i 's that minimize
227 the sum of squared errors (ε_j) over all SNPs. The SSC for individual cows, *per se*, can be
228 estimated as $SCC_i = SCC_{tank} \times V_{tank} \times f_i / V_i$, where SCC_{tank} is the SCC measured in the
229 farm's tank, and V_i / V_{tank} is the proportion of the milk volume contributed by cow i .

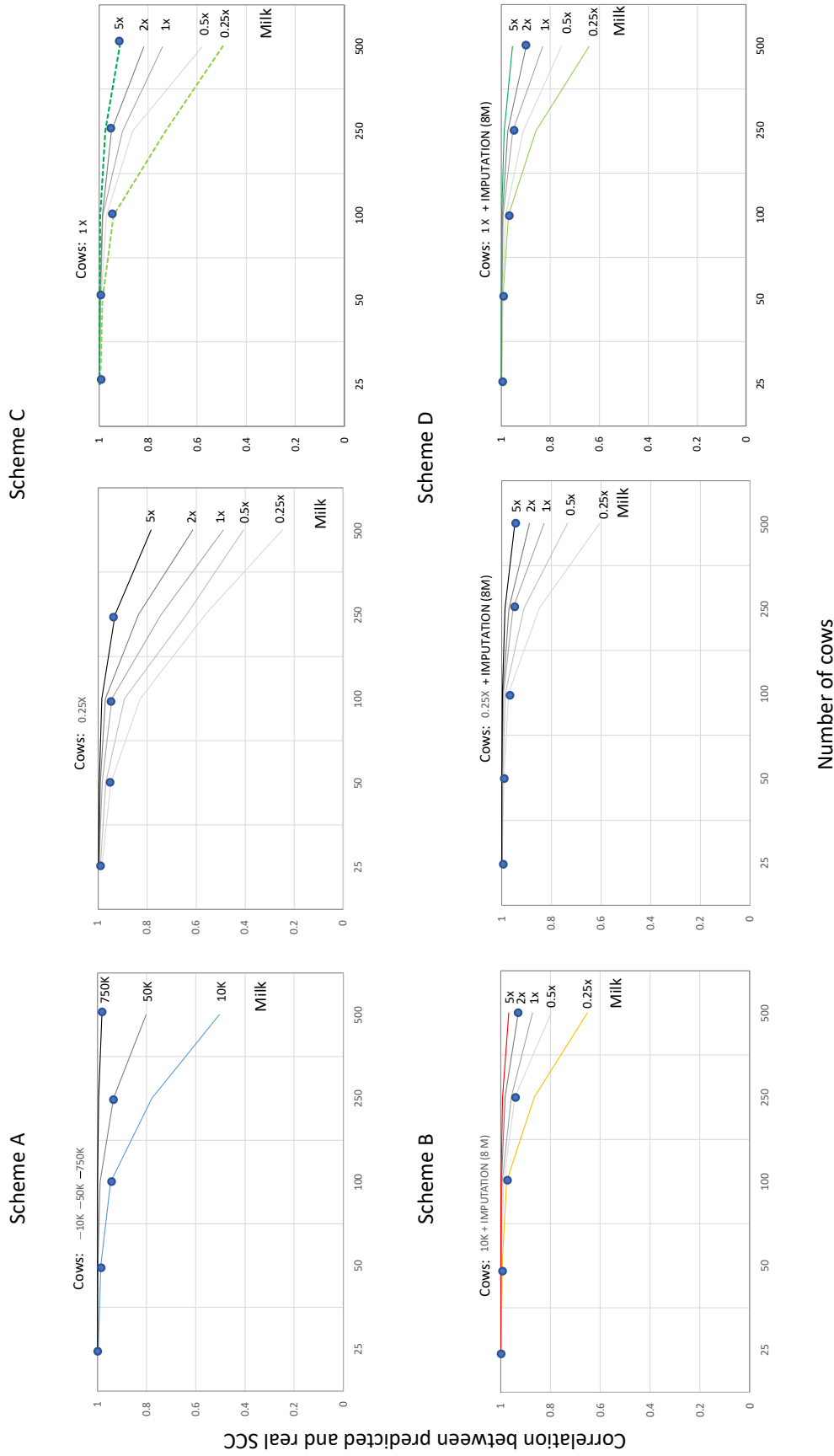
230
231

232
233



234

235 **Figure 2:** Evaluating the efficiency of the proposed approach for the estimation of SCC in the
236 milk of individual cows by genotyping the tank milk, by simulation. **(A)** Reference scheme in
237 which individual cows and tank milk are genotyped with the same array interrogating 10K
238 (LD), 50K (MD) or 700F (HD) SNPs. **(B)** Scheme in which individual cows are genotyped with a
239 LD 10K SNP array and imputed to whole-genome (8 million SNPs), while the tank milk is
240 whole-genome sequenced at depth ranging from 0.25x to 5x. **(C)** Scheme in which individual
241 cows (0.25x and 1x) and tank milk (range: 0.25x to 5x) are genotyped by shallow whole-
242 genome sequencing (SWGS). **(D)** Scheme in which individual cows are genotyped by SWGS
243 (0.25x and 1x) followed by imputation to whole genome (8M SNPs), and tank milk is
244 genotyped by SWGS (range: 0.25x to 5x). In all graphs, the X axis corresponds to the number
245 of cows contributing milk to the tank. The dots mark parameter combinations that yield
246 satisfactory correlations ($r \geq 0.9$). Colored lines correspond to conditions that were used
247 with the real data as shown in Fig. 2.
248

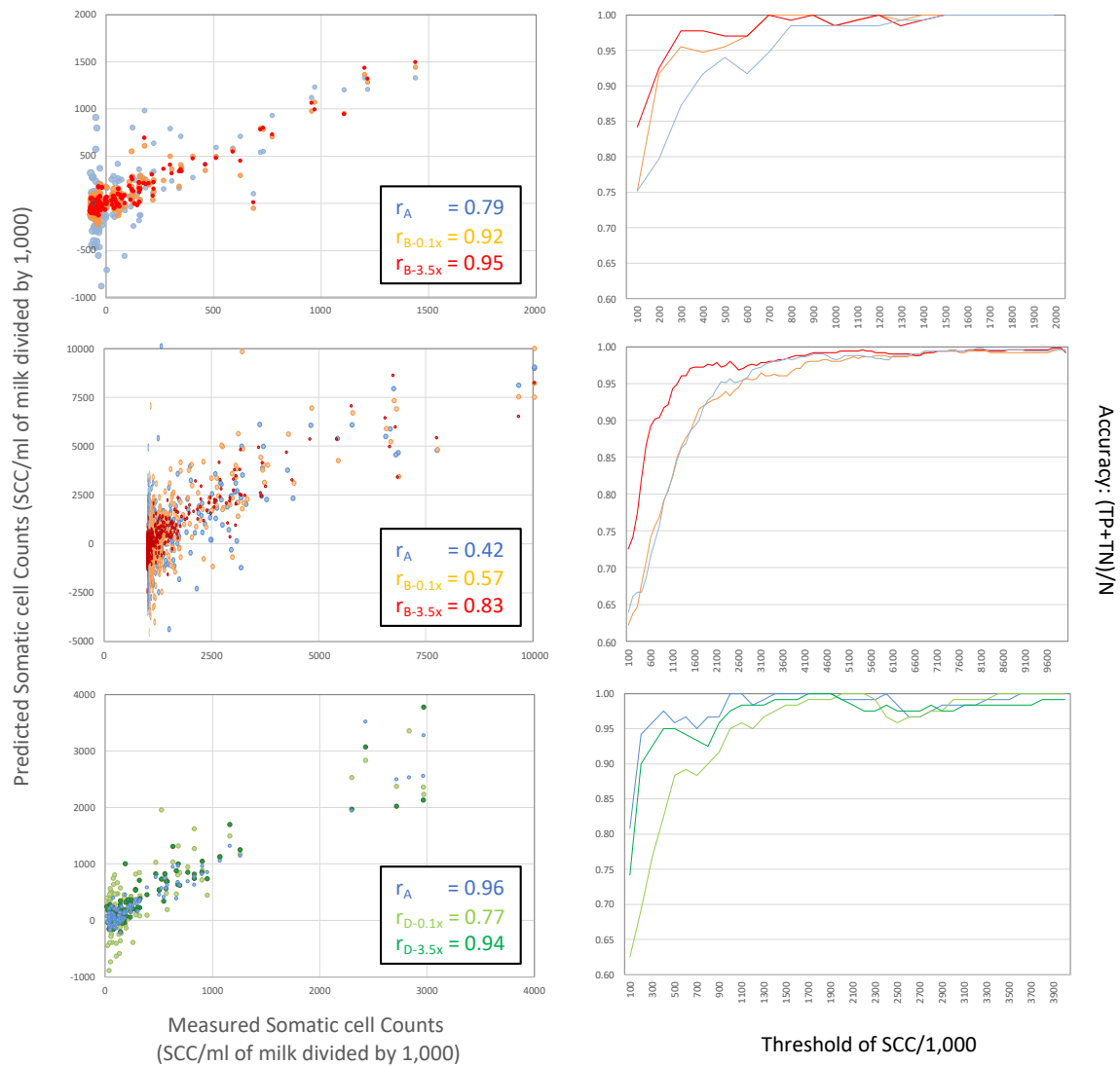


249
250

Correlation between predicted and real SCC

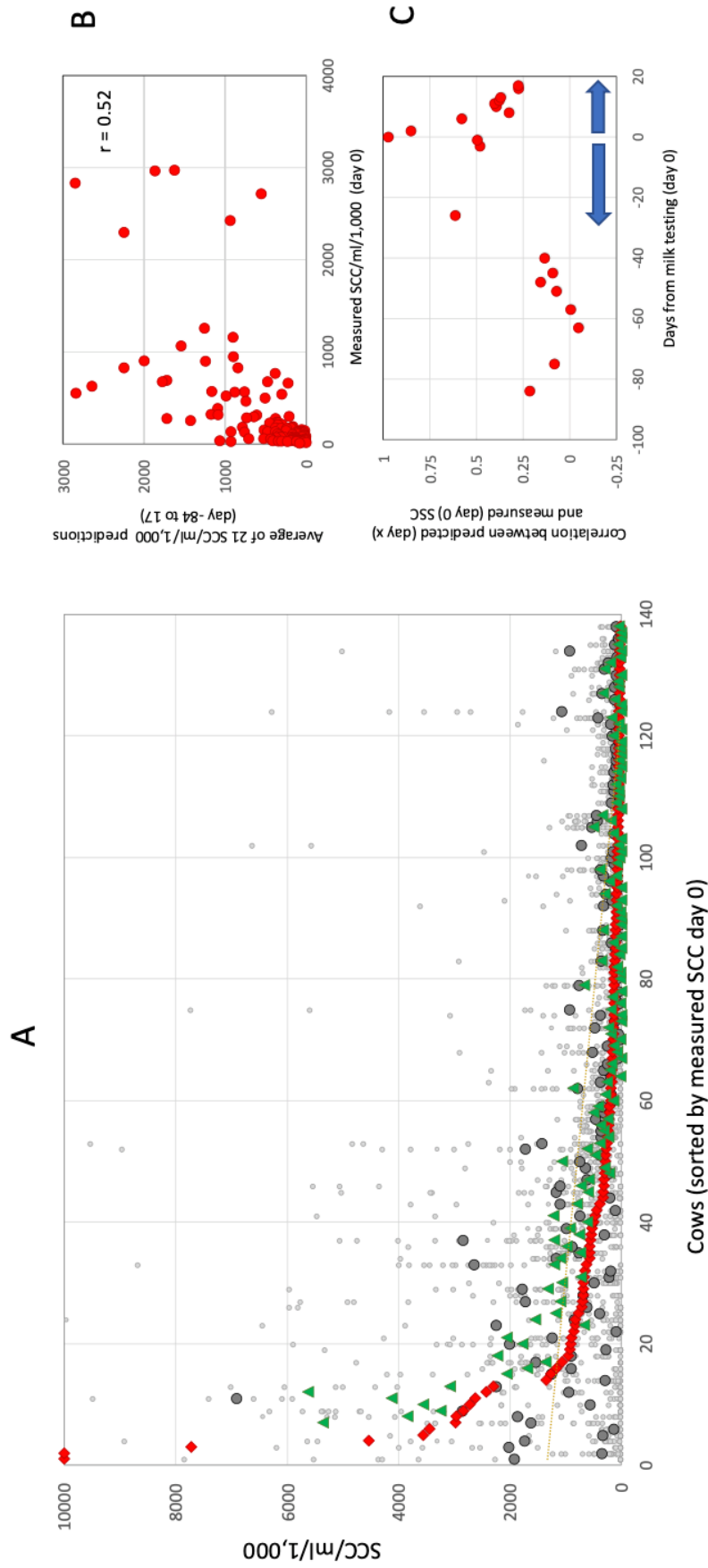
Number of cows

251 **Figure 3:** Correlation between predicted and measured SCC in the milk of individual cows (left
 252 column), as well as accuracies in classifying cows with SCC above and below a chosen
 253 threshold value (right column), in farms with 133 (top row), 520 (middle row) and 120
 254 (bottom row) cows, using scheme A (blue), scheme B (red), or scheme D (green). Scheme A:
 255 cows and tank milk genotyped with LD SNP arrays (17K), no imputation. Scheme B: cows
 256 genotyped with LD array and imputed to 13M SNPs, tank milk sequenced 3.5x (red) or 0.1x
 257 (orange). Scheme D: cows genotyped by whole-genome sequencing (1x) and imputation to
 258 HD, and tank milk sequenced at 3.5x (dark green) or 0.1x (light green).



259
 260

261 **Figure 4: (A)** SCC predicted using scheme A for 21 tank milk samples collected over a 100-day
262 period from 138 cows total. Small grey circles: 20 predictions per cow. Large grey circles:
263 average of 21 measurements per cow. Red diamond: SCC measured on day 0. Green triangle:
264 SSC predictions on day 0. **(B)** Relationship between SCC values measured on day 0 and
265 average of 21 predictions sampled over a 100-day period (days -84 to +17). **(C)** Correlations
266 between measured (day 0) and predicted (day x) SCC as a function of the number of days from
267 day 0.



269 **Methods**

270 **Simulated data.** Reference scheme (A): We simulated farms with n (25, 50, 100, 250 and 500)
271 cows contributing milk to the tank. Cows were genotyped with SNP arrays for m (10K, 50K,
272 or 750K) markers without error. Minor Allele Frequencies (MAFs) were sampled from a
273 uniform $]0,0.5]$ distribution, and genotypes from the corresponding Hardy-Weinberg
274 distributions. SCS of individual cows (SCS_i) were simulated by sampling values from a
275 Weibull distribution with scale parameter $\alpha=1$ and shape parameter $\beta=2$, and multiplying the
276 ensuing value by 200,000. Exact B-allele frequencies of individual SNPs (BAF_j) in the milk
277 were determined for each SNP j based on the combination of cellular contribution of the n
278 cows to the milk, and their genotype. It was assumed that B-allele frequencies were
279 estimated with a normally distributed error $N(0, 0.0025)$ (i.e. SE = 0.05), yielding $m \widehat{BAF}_j$.

280 Scheme B: Same setting as in the reference scheme with the following additions. For cows
281 genotyped for 10K or 50K SNPs, we simulated imputation by augmenting the data to 8 million
282 (M) genotypes using an error model mimicking real, MAF-dependent imputation accuracy.
283 The error model was constructed using a real data set for 800 unrelated Holstein-Friesian
284 individuals that were genotyped for the Illumina 777K array. This data set was split into a set
285 of 200 and a set of 600 individuals. The set of 200 was reduced first to the genotypes
286 interrogated by the Illumina 10K (LD) array and then to the genotypes interrogated by the
287 Illumina 50K SNP arrays. The reduced SNP sets were imputed back to the content of the
288 Illumina 777K (HD) SNP array using the 600 individuals as reference population. The
289 frequencies of imputing a given genotype depending on the real genotype, were scored for
290 MAF bins of 0.01 separately for the LD and 50K array data. We simulated genotyping-by-
291 sequencing of tank milk as follows. For each of the 8M SNP positions, we sampled local read
292 depth ($r \in \text{integers}$) from a Poisson distribution with mean C , where C is the average genome-
293 wide coverage (0.25, 0.5, 1, 2 or 5). We then sampled r reads, each with a probability = BAF_j
294 (computed as above) of being the B-allele. Scheme C: Individual SNP genotypes and tank B-
295 allele frequencies (BAF_j) were generated as in scheme A (genotypes at 8 M SNP positions).
296 It was assumed that milk tank was genotyped by SWGS at average coverage of C (0.25, 0.5, 1,
297 2 or 5) and cows were genotyped by SWGS at average coverage of C (0.25, 0.5, or 1).
298 Genotyping-by-sequencing of individual cows was simulated by (i) sampling, for each of 8M
299 SNP positions, local read depth ($r \in \text{integers}$) from a Poisson distribution with mean C , and

300 (ii) sampling r reads with probability 0, 0.5 or 1 to be the alternate allele (A) depending on the
301 genotype of the cow (RR, RA or AA). Genotyping-by-sequencing of the tank milk was done as
302 in Scheme A. Scheme D: Identical to scheme C except that cow genotypes were generated
303 at 8M SNP position using a MAF- and sequence-depth dependent imputation error model.
304 The error model was constructed using available SWGS data down sampled to 1x (176 cows)
305 or 0.25x coverage (192 cows). The cows were imputed to HD (777K SNPs) using a reference
306 population of 800 unrelated Holstein-Friesian individuals that were genotyped with the
307 Illumina 777K array. At each of the 777K SNP positions, the likelihood of the sequence data
308 under the three possible genotypes (RR, AR and AA), were computed following Chan et al.³,
309 as:

$$310 \quad L(nr_R, nr_A | "RR", \varepsilon) = \binom{nr_R + nr_A}{nr_A} \times (1 - \varepsilon)^{nr_{AR}} \times \varepsilon^{nr_{RA}}$$

$$311 \quad L(nr_R, nr_A | "RA", \varepsilon) = \binom{nr_R + nr_A}{nr_A} \times 0.5^{(nr_R + nr_A)}$$

$$312 \quad L(nr_R, nr_A | "AA", \varepsilon) = \binom{nr_R + nr_A}{nr_A} \times (1 - \varepsilon)^{nr_A} \times \varepsilon^{nr_R}$$

313 where nr_R (respectively nr_A) is the number of R (respectively A reads) and ε is the sequencing
314 error rate set at 0.01. The corresponding $\log_{10} L$ were used as input for Beagle4¹. Variant
315 positions without sequence coverage in any of the 176 (192) cows (hence not imputed by
316 Beagle4) were dealt with in a second round of imputation using Beagle5². The imputation
317 accuracy was evaluated in 0.01 MAF-bins by comparing imputed and real genotypes at the
318 ~17K variant positions interrogated by the Illumina LD array.

319 **Real data.** Data set 1: We obtained a sample of tank milk from a farm in France milking 133
320 Holstein-Friesian cows. All had been genotyped with an Illumina LD array interrogating 17K
321 SNPs using standard procedures. For all cows, genotypes were imputed to whole genome
322 using a reference population of 743 Holstein-Friesian animals sequenced at average depth of
323 15x (range: 4-48) and the Beagle software (v5.0)¹ yielding allelic dosages for a total of 13
324 million SNPs. Individual milk records, including volume and SCC (cells/ml) measured on the
325 day of the sample collection, were obtained for all cows that had contributed milk to the tank.
326 DNA was isolated from 1.5 ml tank milk using the NucleoMag kit (Macherey-Nagel). The tank
327 milk DNA was first genotyped using the Illumina LD array interrogating 17K SNPs. An Illumina
328 compatible NGS library was then prepared with 50ng of genomic DNA using the KAPA
329 HyperPlus kit (Roche). Sequencing was performed on a NextSeq500 instrument (Illumina),

330 yielding 63 million paired end reads of 2*75 bp, corresponding to a genome coverage of 3.5x.
331 Reads were mapped to the bosTau8 genome build using BWA mem. Reference (R) and
332 alternate (A) alleles were counted at 13M SNP positions of the HD array using the Bam-
333 ReadCount tool (<https://github.com/genome/bam-readcount.git>) for reads with a minimum
334 mapping quality of 30. Data set 2: We obtained samples of tank milk from a Belgian farm
335 including milk from 520 Holstein-Friesian cows. Milk volume and SCC (cells/ml) measured on
336 the same day, were obtained for all cows that had contributed milk to the tank. All cows were
337 genotyped with the Illumina LD array interrogating 17K SNPs using standard procedures, and
338 imputed to whole genome using whole genome sequence data (average depth: 15x; range:
339 4x-48x) from 743 Holstein-Friesian animals as reference (M. Georges, unpublished) and the
340 Beagle software (v5.0)² yielding allelic dosages for a total of 13 million SNPs. DNA extraction
341 from the tank milk samples and genotyping with the Illumina LD (17K) array were conducted
342 as for dataset 1. For sequencing of the tank milk, an illumina compatible sequencing library
343 was prepared using 12 ng of DNA and the Riptide High Throughput Rapid Library Prep
344 Kit (iGenomx). The library was sequenced on an Illumina NextSeq500 2*150 paired end flow
345 cell at 4X coverage. Data set 3: We obtained samples of tank milk from a Belgian farm
346 including milk from 120 Holstein-Friesian cows. Milk volume and SCC (cells/ml) measured on
347 the same day, were obtained for all cows that had contributed milk to the tank. All cows were
348 genotyped with the Illumina LD array interrogating 17K SNPs using standard procedures, and
349 imputed to whole genome using whole genome sequence data (average depth: 15x; range:
350 4x-48x) from 743 Holstein-Friesian animals as reference (M. Georges, unpublished) and the
351 Beagle software (v5.0)² yielding allelic dosages for a total of 13 million SNPs. We additionally
352 prepared Illumina compatible NGS library for each cow, using 12 ng of genomic DNA and the
353 Riptide High Throughput Rapid Library Prep Kit (iGenomx). Libraries were sequenced on an
354 Illumina Novaseq S4 2*150 paired end flow cell at average 1.08x depth (range: 0.26x-1.73x).
355 Cow genotype-by-sequencing data were imputed to HD (777K) density using a reference
356 population of 800 Holstein-Friesian animals genotyped with the bovine HD Illumina array
357 (777K SNPs) and the Beagle software (v5.0)² yielding allelic dosages for a total of 777K SNPs.
358 DNA extraction from the tank milk samples, genotyping with the Illumina LD (17K) array, and
359 sequencing (coverage 4x) were conducted as for datasets 1&2. Data set 4: In addition to
360 obtaining a sample of tank milk on the day of the milk recording (i.e. yielding the SCC
361 measured using with a cell counter) for the Belgian farm with 120 cows, we weekly collected

362 an additional 11 tank milk samples before and 9 samples after, spanning a total period of ~3
 363 months. The corresponding DNA samples were genotyped using the Illumina LD (17K) array.
 364

365 **Statistical model.** We defined a set of m linear equations of the form:

$$366 \quad \widehat{BAF}_j = \sum_{i=1}^m f_i \times d_{ij} + \varepsilon_j$$

367 in which f_i is the proportion of the DNA in the tank milk contributed by cow i , d_{ij} is the
 368 “dosage” of the alternate allele A for cow i and marker j , and ε_j is the error term for marker
 369 j . When genotyping the tank milk with arrays, \widehat{BAF}_j corresponds to the B-allele frequency
 370 estimated by Genome Studio (Illumina). When genotyping the tank milk by SWGS, \widehat{BAF}_j
 371 corresponds to the proportion of A reads at the corresponding genome position. For cow
 372 genotypes obtained with arrays, d_{ij} corresponds to 0, 0.5 or 1 for genotypes RR, RA and AA,
 373 respectively. For cow genotypes obtained by imputation, d_{ij} is the dosage of the A allele
 374 estimated by Beagle. For cow genotypes obtained by SWGS, $d_{ij} = 0.5 \times$
 375 $P("RA" | nr_R, nr_A, q_j) + P("AA" | nr_R, nr_A, q_j)$ where nr_R (respectively nr_A) is the number of R
 376 (respectively A reads) for marker j and cow i , and q_j is the population frequency of the A allele
 377 of marker j .

$$379 \quad P("RA" | nr_R, nr_A, q_j) = \frac{2q_j(1-q_j) \times 0.5^{nr_R} \times 0.5^{nr_A} \times \frac{(nr_R + nr_A)!}{nr_R!}}{(1-q_j)^2 \times 1^{nr_R} \times 0^{nr_A} + 2q_j(1-q_j) \times 0.5^{nr_R} \times 0.5^{nr_A} \times \frac{(nr_R + nr_A)!}{nr_R!} + q_j^2 \times 0^{nr_R} \times 1^{nr_A}}$$

380

$$381 \quad P("AA" | nr_R, nr_A, q_j) = \frac{q_j^2 \times 0^{nr_R} \times 1^{nr_A}}{(1-q_j)^2 \times 1^{nr_R} \times 0^{nr_A} + 2q_j(1-q_j) \times 0.5^{nr_R} \times 0.5^{nr_A} \times \frac{(nr_R + nr_A)!}{nr_R!} + q_j^2 \times 0^{nr_R} \times 1^{nr_A}}$$

382

383 For SNPs j without usable information for cow i (f.i. genotyping failure or no covering reads)
 384 d_{ij} was set at \widehat{BAF}_j .

385 The f_i 's were estimated by least square analysis, i.e. by minimizing $\sum_{j=1}^m \varepsilon_j^2$. When the tank
 386 milk was genotyped by SWGS, we also performed a weighted least square analysis, i.e. we
 387 estimated f_i 's by minimizing $\sum_{j=1}^m w_j \varepsilon_j^2$, where w_j is the coverage ($nr_R + nr_A$).

388 The SCC_i 's were calculated from the f_i 's

$$389 \quad SCC_i = SCC_{tank} \times V_{tank} \times f_i / V_i$$

390 Where V_{tank} and V_i are the volumes of milk in the tank and contributed by cow i , respectively.

391 The accuracies of the predictions were measured by the (i) correlation (r) between real and
392 estimated SCC_i , and/or (ii) the ability to discriminate animals with SCC above versus below a
393 certain threshold value measured as $(T_P + T_N)/m$, where T_P stands for the number of true
394 positives, T_N for the number of true negatives, and m for the total number of cows.

395 To test the effect of sequence depth on accuracy we sampled reads overlapping SNP positions
396 with probability x , such that $E(C \times x) = D$, where D is the desired sequence depth.

397

398 **References**

- 399 1. Browning BL, Browning SR. A unified approach to genotype imputation and haplotype
400 phase inference for large data sets of trios and unrelated individuals. *Am J Hum Genet*
401 84:210-223 (2009).
- 402 2. Browning BL, Zhou Y, Browning SR. A one-penny imputed genome from next generation
403 reference panels. *Am J Hum Genet* 103:338-348 (2018).
- 404 3. Chan AW, Hamblin MT, Jannink J-L. Evaluating imputation algorithms for low depth
405 genotyping-by-sequencing (GBS) data. *PLoS ONE* 11:e0160733 (2016).

406

407