

1 **SNP-based quantitative deconvolution of biological mixtures: application to the detection of cows**
2 **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 Institute &
7 Faculty of Veterinary Medicine, University of Liège.

8
9 Correspondence: michel.georges@uliege.be

10

11 **Abstract**

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

19

20 **Introduction**

21 Mastitis, [i.e. the inflammation of the udder](#), is the most important health issue in dairy cattle. [It is](#)
22 [estimated to cost European farmers > 1 billion € per year in treatment and milk loss \(Hogeveen et al.,](#)
23 [2011\).](#) [Upon inflammation, immune cells migrate in the udder and milk. While milk from healthy](#)
24 [cows typically contains > 100,000 cells per milliliter \(ml\) of milk, these numbers \(referred to as Somatic](#)
25 [Cell Counts or SCC\) typically increase into the millions in case of mastitis. Prior to the manifestation](#)
26 [of overt clinical symptoms, SCC progressively increase in the milk of cows developing mastitis: SCC ≥](#)
27 [200,000 are typically considered to be a sign of pre- or sub-clinical mastitis. Both yield and quality of](#)
28 [the milk of cows with subclinical mastitis is reduced \(Schukken et al., 2003\).](#) Mastitis is routinely
29 managed by periodically counting [SCC](#) in milk samples to preemptively identify cows developing
30 subclinical udder inflammation. As profit margins decrease, farmers tend to forgo milk testing thereby
31 compromising health management. Cost-effective alternatives for rapid detection of cows with
32 subclinical mastitis are [hence needed \(Viguier et al., 2009\).](#)

33 [The milk obtained from individual cows is typically collected in one or more large “milk tanks” on the](#)
34 [farm, before being shipped to dairy factories.](#) We previously proposed that somatic cell counts (SCC)

35 in the milk of individual cows could be estimated [by measuring the allelic frequencies in the tank milk](#)
36 for sufficient numbers of SNPs, provided that all cows contributing milk to the tank be genotyped [once](#)
37 for the corresponding variants. Thus, the proposed method would allow [the identification of a](#)
38 minority of cows with subclinical mastitis by [regularly](#) analyzing a single sample containing a mixture
39 of milk from all the cows on the farm, hence dramatically reducing costs. [Prior to ~2010 estimation](#)
40 [of breeding values to select the best dairy sires and dams used pedigree-based estimates of kinship.](#)
41 [Since then, selection methods increasingly use genome-wide SNP information in a process referred to](#)
42 [as “genomic selection” \(GS\) \(Georges et al., 2019\).](#) As GS is becoming routine [in dairy cattle](#) (including
43 for dams), herds that are fully genotyped with [genome-wide SNP](#) arrays are becoming standard, and
44 the proposed method feasible. We herein demonstrate that by combining low density SNP
45 genotyping or shallow sequencing of the cows and tank milk’s DNA with in silico genotype imputation,
46 individual SCC can be accurately determined and cows with subclinical mastitis effectively identified
47 even in the largest farms (≥ 500). The proposed method has the potential to dramatically improve
48 the monitoring of udder health in dairy farms, and to allow the tracing of the origin of bulk animal
49 food products other than milk.

50

51 Results

52 **Principle of the proposed method.** Assume that cows and tank (i.e. the reservoir in which the milk of
53 the cows is collected) milk are genotyped for a collection of SNPs. [Assume that the interrogated SNPs](#)
54 [are biallelic, each characterized by a A \(say the allele of the reference genome\) and a B allele \(say the](#)
55 [alternate allele\).](#) If all cows contribute identical amounts of DNA to the milk, the expected [proportion](#)
56 [of the B allele \(commonly referred to as “B-allele frequency” when analyzing SNP array data](#)
57 [particularly to search for Copy Number Variants\)](#) in the tank milk corresponds to the frequency of the
58 B allele in the farm’s cow population. The actual DNA amount contributed by each cow depends on
59 the volume of milk [that she](#) produced and its SCC. Unequal DNA contributions will cause slight
60 departures from the expected B allele frequencies in the tank milk. Integrating these shifts over a
61 large number of SNPs in conjunction with the known genotypes of individual cows allows for the
62 estimation of the relative DNA contribution of each cow. [This can for instance be achieved using a set](#)
63 [of \$m\$ linear equations in which the “B-allele frequency” of each SNP \$j\$ \(of \$m\$ \) is modelled as the sum](#)
64 [\(over \$n\$ cows\) of the products of the dosage of the B allele in the genotype of cow \$j\$ \(\$d_{ij}\$, known from](#)
65 [her SNP genotype\) multiplied by the proportion of DNA contributed by cow \$i\$ \(\$f_i\$ \) to the milk. The](#)
66 [proportions of DNA contributed by each cow can then be estimated using for instance least square](#)
67 [methods.](#) Accounting for individual milk volumes and for the SCC in the tank milk allows for the
68 estimation of SCC for individual cows (Fig. 1 and Methods).

69 **Evaluating the proposed method by simulation.** We first evaluated the proposed method by
70 simulation (cfr. Methods). Genotyping the cows and the tank milk using 10K SNP arrays (i.e. low-
71 density (LD) arrays as generally used [in the context of genomic selection](#)) allowed for the accurate
72 estimation of individual SCC for farms with up to 100 cows ($r \geq 0.9$, where r is the correlation
73 between real and estimated SCC) (scheme [I](#)). However, farms with > 100 cows are increasingly
74 common. Medium- (MD, f.i. 50K) and high-density (HD, f.i. 700K) SNP arrays would be needed for
75 the approach to be effective in farms with ≥ 250 or ≥ 500 cows, respectively. Yet – being too
76 expensive - this is presently not a viable proposition (Fig. [2A and Supplemental Table 1](#)). We therefore
77 envisaged a second scheme ([II](#)) in which the cows would still be genotyped with LD SNP arrays (as
78 done in practice) yet imputed ([Marchini & Howie, 2010](#)) to whole genome (8 million SNPs in the
79 simulations) using a sequenced reference population ([Daetwyler et al., 2014](#)), while the DNA of the
80 tank milk would be genotyped by shallow whole-genome sequencing (SWGS). We found that under
81 this scenario sequencing the tank milk at a depth of 0.25 was sufficient for farms with 100 cows, 0.5
82 for farms with 250 cows, and 2 for farms with 500 cows (Fig. [2B](#)). Accuracies were not significantly
83 affected by the density of the SNP arrays, i.e. the method performed as well with LD as with MD arrays
84 (data not shown). Anticipating further advances in sequencing technology, we also envisaged a
85 scheme ([III](#)) in which both cows and tank milk would be genotyped by SWGS. We found that a 1-fold
86 sequencing depth of the tank milk would be sufficient when combined with a 0.25-fold depth for 100
87 cows, while a 5-fold sequencing depth of the tank milk would be needed in combination with 0.25-
88 fold depth for 250 cows and 1-fold depth for 500 cows (Fig. [2C&D](#)). In scheme [III](#), allelic dosage in the
89 cows is directly measured from the number of alternative and reference alleles in the sequence reads.
90 We further explored the effectiveness of augmenting the cow genotype information from SWGS by
91 imputation (scheme [IV](#)). This proved to be effective, reducing the required sequence depth to 0.25-
92 fold for tank milk and 0.25-fold for 100 cows, to 1-fold for tank milk and 0.25-fold for 250 cows, and
93 to 5-fold for tank milk and 0.25-fold for 500 cows (Fig. 2). [The previous simulations make a number of](#)
94 [assumptions that may not apply in the real world: \(i\) SNPs were sampled from a uniform distribution](#)
95 [\(i.e. rare and common SNPs equally represented\), \(ii\) SNPs were assumed to be in linkage equilibrium,](#)
96 [\(iii\) cows on the farm were assumed to be unrelated, and \(iv\) milk volumes were assumed to be known](#)
97 [without error. To more accurately mimic real conditions we repeated the simulations by \(i\) sampling](#)
98 [genotypes from a phased dataset of 750 Holstein-Friesian whole genome sequences \(hence properly](#)
99 [accounting for true MAF distribution, true linkage disequilibrium \(LD\) and relatedness - many of the](#)
100 [sequenced animals were related as parent offspring, full- or half-sibs\), and \(ii\) adding a normally](#)
101 [distributed error with mean 0 and standard deviation of five liter to the simulated milk volumes](#)
102 [\(normally distributed with mean of 30 liter and standard deviation of 10 liter\). This error rate](#)

103 [corresponds approximately to that expected when having to estimate the daily milk volume from the](#)
104 [total lactation yield using a standard lactation curve \(Miel Hostens, personal communication\). We](#)
105 [assumed in these simulations that the genotypes of the cows were known without error and that the](#)
106 [milk was sequenced at a depth ranging from 0.25 to 5 as before. MAF, LD and relatedness jointly](#)
107 [had a relatively modest impact on the accuracy of the method, which could be compensated for by](#)
108 [increasing the sequencing depth of the milk to five-fold and still allowing for accurate estimates even](#)
109 [in farms with 500 cows. Estimating the milk volume with error had a more pronounced impact on](#)
110 [the accuracy making it possible difficult to reach a correlation reaching 0.9 in farms with 500 cows](#)
111 [\(Fig. 2\).](#)

112
113 **Real-world application of the proposed method.** To test the feasibility of our method in the real
114 world, we first collected cow (blood) and tank (milk) samples from a farm milking 133 Holstein-Friesian
115 cows. When only using genotypes from the Illumina LD arrays (17K SNPs) for both cows and tank milk
116 (scheme A), correlations between predicted and measured SCC were 0.91 (or 0.79 when ignoring one
117 cow with SCC > 3 million). We then imputed the cows to whole genome (13M SNPs) using a reference
118 population of ~750 whole genome sequenced Holstein-Friesian animals, and sequenced the tank milk
119 at ~3.5-fold depth. The corresponding correlations (scheme B) were 0.97 (0.95) when using all
120 sequence information, or 0.96 (0.92) when down-sampling sequence information as low as 0.1-fold
121 depth (Fig. 3A). We next performed a similar experiment on a farm milking 520 Holstein-Friesian cows.
122 The correlation between predicted and measured SCC was 0.78 (or 0.42 when ignoring 23 cows with
123 SCC > 3 million) when only using information from the LD array for both cows and tank milk (scheme
124 A). When imputing the cows to whole genome (13M SNPs) and sequencing the milk at ~3.5-fold
125 depth (scheme B), the correlation increased to 0.89 (0.83). Down-sampling the sequence information
126 to 0.1-fold depth reduced the correlation to 0.79 (0.57) (Fig. 3B).

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

136 As costs per bp continue to decline, sequencing is likely to replace array-based genotyping in the
137 future. To test the feasibility of schemes C and D (i.e. genotype the cows by SWGS rather than with
138 SNP arrays, without (C) and with (D) imputation), we collected samples from a farm with 120 Holstein-
139 Friesian cows. All cows were genotyped with the Illumina LD array (17K) as well as sequenced at
140 average 1.08 -fold depth (range: 0.26-1.73). The milk was sequenced at ~3.5-fold depth. The
141 correlation between predicted and measured SCC was 0.97 (or 0.96 when ignoring one cow with SCC
142 > 3 million) under scheme A. Under scheme C, correlations were 0.82 (0.83) when sequencing the
143 milk at 3.5x and 0.75 (0.76) when down-sampling the milk to 0.1x. We then imputed the sequenced
144 cows to HD (770K SNPs) using a population of 800 reference animals genotyped with the HD array
145 (scheme D). The correlation increased to 0.93 (0.94) when sequencing the milk at 3.5x and to 0.83
146 (0.77) when down-sampling the milk to 0.1x (Fig. 3C). Accuracies at SCC threshold of 500,000 were
147 0.96 (scheme A), 0.95 (3.5x) and 0.80 (0.1x) (scheme B), 0.82 (3.5x) and 0.81 (0.1x) (scheme C), and
148 0.95 (3.5x) and 0.88 (0.1x) (scheme D) (Fig. 3C). In summary, (i) combining cow genotyping using SNP
149 arrays with genome-wide imputation with SWGS of tank milk allows for cost-effective identification
150 of cows with subclinical mastitis even in farms with as many as 500 cows per milk tank, and (ii) as
151 sequencing costs continue to decline, arrays-based targeted SNP genotyping of the cows could be
152 replaced by genotyping by SWGS and yield comparable results.

153 **Monitoring SCC dynamics with the proposed method.** Farmers typically measure individual SCC once
154 a month or less. Yet, SCC may rapidly change. The SCC measured on the milk testing date may not be
155 a reliable indicator of the cow's udder health during the intervening period. To examine the SCC
156 dynamics over time, we collected 20 tank milk samples over a 100-day period (day -84 to +17 from
157 day of milk testing) for the farm with 120 cows. Milk samples were genotyped using the Illumina LD
158 array, and individual SCC estimated using scheme A. Fig. 4A shows the SCC predicted every 5 days on
159 average for the 120 cows, sorted by SCC measured on day 0 (=milk testing day). Of note, the
160 correlation between the SCC measured on day 0 and the average of the SCC estimates for the 21
161 collection dates was low ($r = 0.52$)(Fig. 4B) and decreased rapidly with the number of days from milk
162 testing day (Fig. 4C).

163

164 Discussion

165 We herein demonstrate that by combining array-based SNP genotyping and whole-genome
166 imputation for the cows with SWGS of the tank milk, it is possible to accurately estimate SCC for
167 individual cows and hence effectively identify animals with subclinical mastitis even for tanks
168 collecting milk for >500 cows, and this by performing a single analysis for the entire herd. Reagent
169 costs to sequence a mammalian genome at 1-fold depth are now <20€ thus making this a cost-

170 effective proposition. As a matter of fact, the method is being deployed in the field in several
171 countries.

172 Implementing the method requires all cows on the farm to be genotyped. This will increasingly
173 correspond to reality as genotyping costs continue to decrease and genomic selection is more and
174 more used for the selection of cows. In 2016 more than 1.2 million dairy cows had been reportedly
175 genotyped in the US alone⁸ and present worldwide numbers are likely ≥ 3 million. In addition, a
176 reference population of a few hundred animals of the breed of interest that are either HD genotyped
177 (700K) or better whole-genome sequenced are required for accurate imputation. Such reference
178 populations are already available for the most important dairy cattle breeds^{7,9}, and could be easily
179 generated for the remaining ones.

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

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

198

199 **Acknowledgements**

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

203 experiments. [The proposed method is the subject of awarded \(WO/2013/079289\) and filed patents](#)
204 [\(PCT/EP2019/057628\).](#)

205

206 **References**

207 1. Hogeveen H, Huijps K, Lam TJGM. Economic aspects of mastitis: New developments. *N. Z. Vet. J.*
208 59:16–23 (2011).

209 2. Viguier C, Arora S, Gilmartin N, Welbeck K, O’Kennedy R. Mastitis detection: current trends and
210 future perspectives. *Trends Biotechnol* 27: 486-493 (2009).

211 3. Blard G, Zhang Z, Coppieters W, Georges M. Identifying cows with subclinical mastitis by bulk SNP
212 genotyping of tank milk. *J Dairy Sci* 95:4109-4113 (2012).

213 4. Georges M, Charlier C, Hayes B. [Harnessing genomic information for livestock improvement.](#) *Nat*
214 *Rev Genet* 20:135-156 (2019).

215 5. Schukken YH, Wilson DJ, Welcome F, Garrison-Tikofsky L, Gonzales RN. Monitoring udder health
216 and milk quality using somatic cell counts. *Vet Res* 34: 579-596 (2003).

217 6. Marchini J, Howie B. Genotype imputation for genome-wide association studies. *Nat Rev Genet*
218 11:499-511 (2010).

219 7. Daetwyler HD *et al.* Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and
220 complex traits in cattle. *Nat Genet* 46: 858-865 (2014).

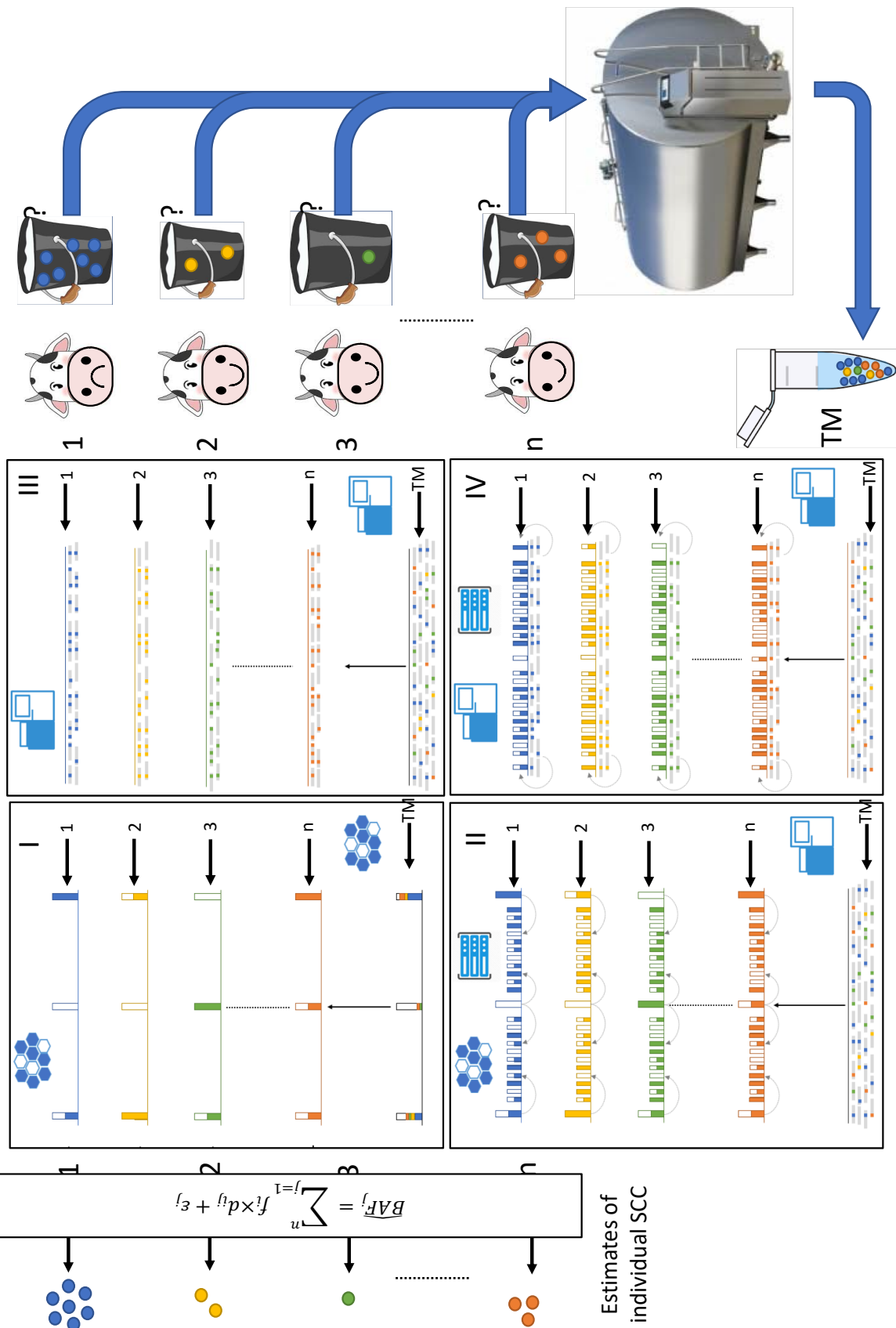
221 8. Wiggans GR, Cole JB, Hubbard SM, Sonstegard TS. Genomic selection in dairy cattle: the USDA
222 experience. *Annu Rev Anim Biosci* 5:309-327 (2017).

223 9. Charlier C. *et al.* Reverse genetic screen for embryonic lethal mutations comprising fertility in
224 cattle. *Genome Res* 26: 1-9 (2016).

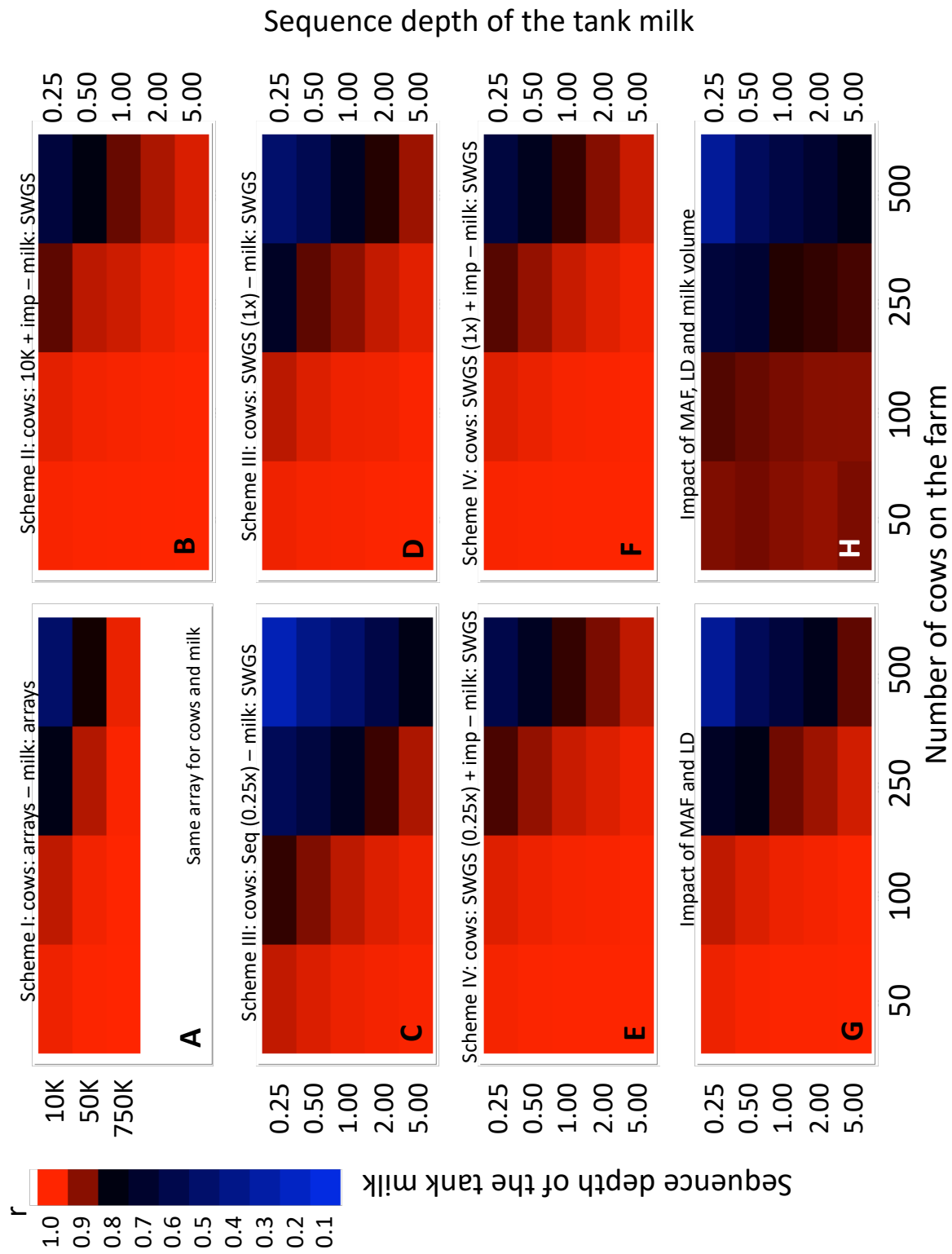
225 10. Brooke-Taylor S, Dwyer K, Woodford K, Kost N. Systematic review of the gastrointestinal effects
226 of A1 compared with A2 β -casein. *Adv Nutr* 8:739-748 (2017).

227

228 **Figure 1:** Estimating Somatic Cell Counts (SCC) in the milk of individual cows by analyzing a sample of
229 milk from the farm's tank. Cows 1 to n contribute different amounts of milk (buckets of various sizes
230 in the figure) to the farm's tank. The milk contains somatic cells (shown as small spheres in the milk
231 colored by cow) whose numbers reflect the health status of the cow's udder. Cow 1 has higher SCC,
232 an indicator of subclinical mastitis. SCC are unknown upon milking (indicated by the "?"). Cows are
233 individually SNP genotyped once. In scheme I this is done using SNP arrays (illustrated by the mesh)
234 yielding genotype information for the limited number of interrogated SNPs (high bars) that can be
235 summarized by the B-allele frequency as shown (white: 0, half colored: 0.5, fully colored: 1). SNP
236 genotypes of individual cows are coded in the same colors as the SCC. In scheme II, the genotypes of
237 the interrogated SNPs are augmented by imputation (illustrated by the computer rack), yielding
238 dosage information (B-allele frequency) for many more SNPs (small bars). In scheme III, cows are
239 genotyped individually by shallow whole genome sequencing (SWGS) (illustrated by the sequencer).
240 Sequence reads (gray lines) are aligned to the reference genome and alternate alleles at SNP positions
241 highlighted as color-coded tics. The B-allele frequency at specific SNP positions is measured as the
242 ratio of the number of reads with the B allele vs the total number of reads. In scheme IV, the genotype
243 information from SWGS is augmented by imputation improving the accuracy of the B-allele frequency
244 estimates for millions of SNPs (small bars). A small sample of milk (T(ank) M(ilk)) is periodically (f.i.
245 monthly or weekly) collected from the farm's tank. DNA is extracted from TM and genotyped using
246 SNP arrays (scheme I) or SWGS (schemes II, III and IV). B-allele frequency for SNP j in the milk (\widehat{BAF}_j)
247 is estimated from the ratio of fluorescence intensities when using SNP arrays, or from the proportion
248 of reads with B allele in SWGS. The SCC of individual cows are estimated from a set of linear equations
249 modelling \widehat{BAF}_j as the sum of B allele dosage (d_{ij}) multiplied by the proportion of the DNA in the tank
250 contributed by cow i (f_i). The estimated proportions of DNA contributed by each cow correspond to
251 the values of f_i 's that minimize the sum of squared errors (ϵ_j) over all SNPs. The SCC for individual
252 cows, *per se*, can be estimated as $SCC_i = SCC_{tank} \times V_{tank} \times f_i / V_i$, where SCC_{tank} is the SCC
253 measured in the farm's tank, and V_i / V_{tank} is the proportion of the milk volume contributed by cow i .
254
255

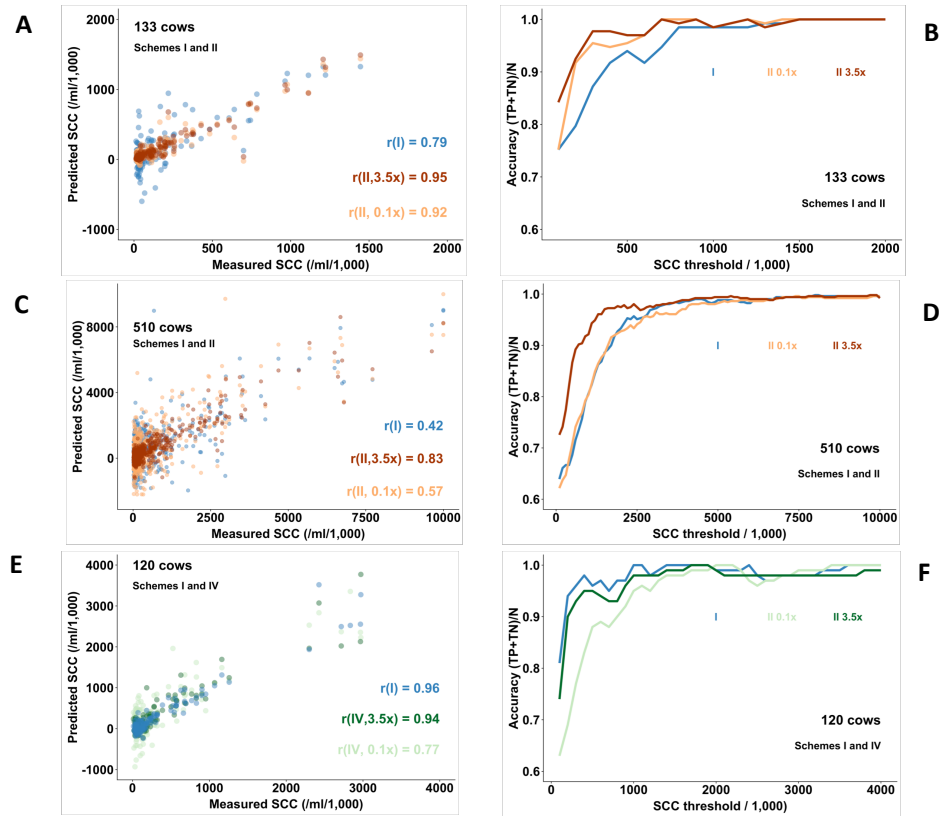


257 **Figure 2:** Evaluating the efficiency of the proposed approach by simulation. **(A)** Reference scheme I in
258 which individual cows and tank milk are genotyped with the same array interrogating 10K (LD), 50K
259 (MD) or 700F (HD) SNPs. **(B)** Scheme II in which individual cows are genotyped with a LD 10K SNP array
260 and imputed to whole-genome (8 million SNPs), while the tank milk is whole-genome sequenced at
261 depth ranging from 0.25x to 5x. **(C)** Scheme III in which individual cows (0.25x) and tank milk (range:
262 0.25x to 5x) are genotyped by shallow whole-genome sequencing (SWGS). **(D)** Same as C except that
263 individual cows are sequenced at 1x depth. **(E)** Scheme IV in which individual cows are genotyped by
264 SWGS (0.25x) followed by imputation to whole genome (8M SNPs), and tank milk is genotyped by
265 SWGS (range: 0.25x to 5x). **(F)** Same as E except that individual cows are sequenced at 1x depth. **(G)**
266 Scheme in which the cow genotypes are sampled from a real dataset hence conform to reality with
267 regards to distribution of MAF, LD and relatedness. Genotypes of the cows are assumed to be known
268 (very similar to II and IV) and tank milk genotyped by SWGS (range: 0.25x to 5x). **(H)** Same as G except
269 that the milk volume is estimated with error. The color code used to quantify the correlations
270 between predicted and real SCC is shown. Corresponding numerical values are provided in Suppl.
271 Table 1
272



273
274
275
276

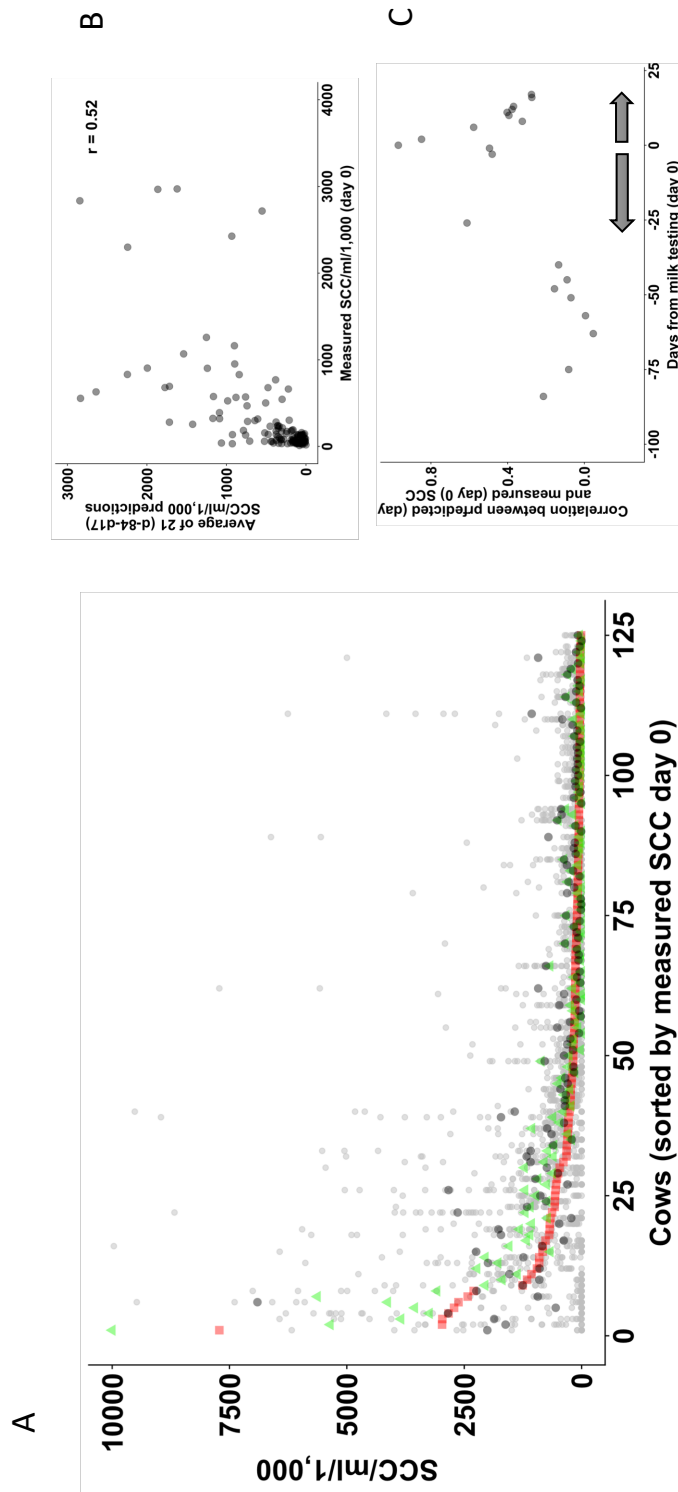
277 **Figure 3:** Correlation between predicted and measured SCC in the milk of individual cows (A,C,E), as
278 well as accuracies in classifying cows with SCC above and below a chosen threshold value (B,D,F), in
279 farms with 133 (A,B), 520 (C,D) and 120 (E,F) cows, using scheme I (blue), scheme II (red), or scheme
280 IV (green). Scheme I: cows and tank milk genotyped with LD SNP arrays (17K), no imputation. Scheme
281 II: cows genotyped with LD array and imputed to 13M SNPs, tank milk sequenced 3.5x (red) or 0.1x
282 (orange). Scheme IV: cows genotyped by whole-genome sequencing (1x) and imputation to HD, and
283 tank milk sequenced at 3.5x (dark green) or 0.1x (light green).
284



285

286

287 **Figure 4: (A)** SCC predicted using scheme A for 21 tank milk samples collected over a 100-day period
288 from 125 cows total. Small grey circles: 20 predictions per cow. Large grey circles: average of 21
289 measurements per cow. Red square: SCC measured on day 0. Green triangle: SCC predictions on day
290 0. **(B)** Relationship between SCC values measured on day 0 and average of 21 predictions sampled
291 over a 100-day period (days -84 to +17). **(C)** Correlations between measured (day 0) and predicted
292 (day x) SCC as a function of the number of days from day 0.



294 **Methods**

295 **Simulated data.** Reference scheme (A): We simulated farms with n (25, 50, 100, 250 and 500) cows
296 contributing milk to the tank. Cows were genotyped with SNP arrays for m (10K, 50K, or 750K) markers
297 without error. Minor Allele Frequencies (MAFs) were sampled from a uniform $]0,0.5]$ distribution, and
298 genotypes from the corresponding Hardy-Weinberg distributions. SCC of individual cows (SCC_i) were
299 simulated by sampling values from a Weibull distribution with scale parameter $\alpha=1$ and shape
300 parameter $\beta=2$, and multiplying the ensuing value by 200,000. Exact B-allele frequencies of individual
301 SNPs (BAF_j) in the milk were determined for each SNP j based on the combination of cellular
302 contribution of the n cows to the milk, and their genotype. It was assumed that B-allele frequencies
303 were estimated with a normally distributed error $N(0, 0.0025)$ (i.e. SE = 0.05), yielding $m \widehat{BAF}_j$.
304 Scheme B: Same setting as in the reference scheme with the following additions. For cows genotyped
305 for 10K or 50K SNPs, we simulated imputation by augmenting the data to 8 million (M) genotypes
306 using an error model mimicking real, MAF-dependent imputation accuracy. The error model was
307 constructed using a real data set for 800 unrelated Holstein-Friesian individuals that were genotyped
308 for the Illumina 777K array. This data set was split into a set of 200 and a set of 600 individuals. The
309 set of 200 was reduced first to the genotypes interrogated by the Illumina 10K (LD) array and then to
310 the genotypes interrogated by the Illumina 50K SNP arrays. The reduced SNP sets were imputed back
311 to the content of the Illumina 777K (HD) SNP array using the 600 individuals as reference population.
312 The frequencies of imputing a given genotype depending on the real genotype, were scored for MAF
313 bins of 0.01 separately for the LD and 50K array data. We simulated genotyping-by-sequencing of tank
314 milk as follows. For each of the 8M SNP positions, we sampled local read depth ($r \in \text{integers}$) from a
315 Poisson distribution with mean C , where C is the average genome-wide coverage (0.25, 0.5, 1, 2 or 5).
316 We then sampled r reads, each with a probability = BAF_j (computed as above) of being the B-allele.
317 Scheme C: Individual SNP genotypes and tank B-allele frequencies (BAF_j) were generated as in
318 scheme A (genotypes at 8 M SNP positions). It was assumed that milk tank was genotyped by SWGS
319 at average coverage of C (0.25, 0.5, 1, 2 or 5) and cows were genotyped by SWGS at average coverage
320 of C (0.25, 0.5, or 1). Genotyping-by-sequencing of individual cows was simulated by (i) sampling, for
321 each of 8M SNP positions, local read depth ($r \in \text{integers}$) from a Poisson distribution with mean C ,
322 and (ii) sampling r reads with probability 0, 0.5 or 1 to be the alternate allele (**B**) depending on the
323 genotype of the cow (**AA**, **AB** or **BB**). Genotyping-by-sequencing of the tank milk was done as in
324 Scheme A. Scheme D: Identical to scheme C except that cow genotypes were generated at 8M SNP
325 position using a MAF- and sequence-depth dependent imputation error model. The error model was
326 constructed using available SWGS data down sampled to 1x (176 cows) or 0.25x coverage (192 cows).
327 The cows were imputed to HD (777K SNPs) using a reference population of 800 unrelated Holstein-

328 Friesian individuals that were genotyped with the Illumina 777K array. At each of the 777K SNP
329 positions, the likelihood of the sequence data under the three possible genotypes (AA, AB and BB),
330 were computed following Chan et al.³, as:

$$331 \quad L(nr_A, nr_B | "AA", \varepsilon) = \binom{nr_A + nr_B}{nr_B} \times (1 - \varepsilon)^{nr_A} \times \varepsilon^{nr_B}$$

$$332 \quad L(nr_A, nr_B | "AB", \varepsilon) = \binom{nr_A + nr_B}{nr_B} \times 0.5^{(nr_A + nr_B)}$$

$$333 \quad L(nr_A, nr_B | "BB", \varepsilon) = \binom{nr_A + nr_B}{nr_B} \times (1 - \varepsilon)^{nr_B} \times \varepsilon^{nr_A}$$

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

340 **Real data.** Data set 1: We obtained a sample of tank milk from a farm in France milking 133 Holstein-
341 Friesian cows. All had been genotyped with an Illumina LD array interrogating 17K SNPs using standard
342 procedures. For all cows, genotypes were imputed to whole genome using a reference population of
343 743 Holstein-Friesian animals sequenced at average depth of 15x (range: 4-48) and the Beagle
344 software (v5.0)¹ yielding allelic dosages for a total of 13 million SNPs. Individual milk records, including
345 volume and SCC (cells/ml) measured on the day of the sample collection, were obtained for all cows
346 that had contributed milk to the tank. DNA was isolated from 1.5 ml tank milk using the NucleoMag
347 kit (Macherey-Nagel). The tank milk DNA was first genotyped using the Illumina LD array interrogating
348 17K SNPs. An Illumina compatible NGS library was then prepared with 50ng of genomic DNA using
349 the KAPA HyperPlus kit (Roche). Sequencing was performed on a NextSeq500 instrument (Illumina),
350 yielding 63 million paired end reads of 2*75 bp, corresponding to a genome coverage of 3.5x. Reads
351 were mapped to the bosTau8 genome build using BWA mem. Reference (R) and alternate (A) alleles
352 were counted at 13M SNP positions of the HD array using the Bam-ReadCount tool
353 (<https://github.com/genome/bam-readcount.git>) for reads with a minimum mapping quality of 30.

354 Data set 2: We obtained samples of tank milk from a Belgian farm including milk from 520 Holstein-
355 Friesian cows. Milk volume and SCC (cells/ml) measured on the same day, were obtained for all cows
356 that had contributed milk to the tank. All cows were genotyped with the Illumina LD array
357 interrogating 17K SNPs using standard procedures, and imputed to whole genome using whole
358 genome sequence data (average depth: 15x; range: 4x-48x) from 743 Holstein-Friesian animals as
359 reference (M. Georges, unpublished) and the Beagle software (v5.0)² yielding allelic dosages for a total
360 of 13 million SNPs. DNA extraction from the tank milk samples and genotyping with the Illumina LD

361 (17K) array were conducted as for dataset 1. For sequencing of the tank milk, an illumina compatible
362 sequencing library was prepared using 12 ng of DNA and the Riptide High Throughput Rapid Library
363 Prep Kit (iGenomx). The library was sequenced on an Illumina NextSeq500 2*150 paired end flow cell
364 at 4X coverage. Data set 3: We obtained samples of tank milk from a Belgian farm including milk from
365 120 Holstein-Friesian cows. Milk volume and SCC (cells/ml) measured on the same day, were obtained
366 for all cows that had contributed milk to the tank. All cows were genotyped with the Illumina LD array
367 interrogating 17K SNPs using standard procedures, and imputed to whole genome using whole
368 genome sequence data (average depth: 15x; range: 4x-48x) from 743 Holstein-Friesian animals as
369 reference (M. Georges, unpublished) and the Beagle software (v5.0)² yielding allelic dosages for a total
370 of 13 million SNPs. We additionally prepared Illumina compatible NGS library for each cow, using 12
371 ng of genomic DNA and the Riptide High Throughput Rapid Library Prep Kit (iGenomx). Libraries were
372 sequenced on an Illumina Novaseq S4 2*150 paired end flow cell at average 1.08x depth (range: 0.26x-
373 1.73x). Cow genotype-by-sequencing data were imputed to HD (777K) density using a reference
374 population of 800 Holstein-Friesian animals genotyped with the bovine HD Illumina array (777K SNPs)
375 and the Beagle software (v5.0)² yielding allelic dosages for a total of 777K SNPs. DNA extraction
376 from the tank milk samples, genotyping with the Illumina LD (17K) array, and sequencing (coverage
377 4x) were conducted as for datasets 1&2. Data set 4: In addition to obtaining a sample of tank milk on
378 the day of the milk recording (i.e. yielding the SCC measured using with a cell counter) for the Belgian
379 farm with 120 cows, we weekly collected an additional 11 tank milk samples before and 9 samples
380 after, spanning a total period of ~3 months. The corresponding DNA samples were genotyped using
381 the Illumina LD (17K) array.

382

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

384
$$\widehat{BAF}_j = \sum_{i=1}^n f_i \times d_{ij} + \varepsilon_j$$

385 in which f_i is the proportion of the DNA in the tank milk contributed by cow i , d_{ij} is the “dosage” of
386 the alternate allele A for cow i and marker j , and ε_j is the error term for marker j . When genotyping
387 the tank milk with arrays, \widehat{BAF}_j corresponds to the B-allele frequency estimated by Genome Studio
388 (Illumina). When genotyping the tank milk by SWGS, \widehat{BAF}_j corresponds to the proportion of A reads
389 at the corresponding genome position. For cow genotypes obtained with arrays, d_{ij} corresponds to
390 0, 0.5 or 1 for genotypes AA, AB and BB, respectively. For cow genotypes obtained by imputation, d_{ij}
391 is the dosage of the B allele estimated by Beagle. For cow genotypes obtained by SWGS, $d_{ij} =$
392 $0.5 \times P("AB" | nr_A, nr_B, q_j) + P("BB" | nr_A, nr_B, q_j)$ where nr_A (respectively nr_B) is the number of A

393 (respectively B reads) for marker j and cow i , and q_j is the population frequency of the B allele of
 394 marker j .

395

$$396 \quad P("AB" | nr_A, nr_B, q_j) = \frac{2q_j(1-q_j) \times 0.5^{nr_A} \times 0.5^{nr_B} \times \frac{(nr_A + nr_B)!}{nr_A! \times nr_B!}}{(1-q_j)^2 \times 1^{nr_A} \times 0^{nr_B} + 2q_j(1-q_j) \times 0.5^{nr_A} \times 0.5^{nr_B} \times \frac{(nr_A + nr_B)!}{nr_A! \times nr_B!} + q_j^2 \times 0^{nr_A} \times 1^{nr_B}}$$

397

$$398 \quad P("BB" | nr_A, nr_B, q_j)$$

$$399 \quad = \frac{q_j^2 \times 0^{nr_A} \times 1^{nr_B}}{(1-q_j)^2 \times 1^{nr_A} \times 0^{nr_B} + 2q_j(1-q_j) \times 0.5^{nr_A} \times 0.5^{nr_B} \times \frac{(nr_A + nr_B)!}{nr_A! \times nr_B!} + q_j^2 \times 0^{nr_A} \times 1^{nr_B}}$$

400

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

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

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

$$407 \quad SCC_i = SCC_{tank} \times V_{tank} \times f_i / V_i$$

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

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

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

415

416 References

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

423